Heterogeneous Endogenous Effects in Networks *

Sida Peng[†]

November 18, 2016

Abstract

Spatial econometrics has been widely used to study endogenous effects in network structures. Existing spatial autoregression models (SARs) implicitly assume that each individual in the network has the same endogenous effects on others. However, some individuals are more influential than others. For example, Banerjee et al. (2013) documents that individuals directly connected with some village leaders are more likely to join the micro-finance program than those connected to someone else. I develop a SAR model that allows for individual-specific endogenous effects and propose a two-stage LASSO (2SLSS) procedure to identify influential individuals in a network. Under an assumption of sparsity: only a subset of individuals (which can increase with sample size n) is influential, I show that my 2SLSS estimator for individual-specific endogenous effects is consistent and achieves asymptotic normality. I also develop robust inference including uniformly valid confidence intervals. These results also carry through to scenarios where the influential individuals are not sparse. I extend the analysis to allow for multiple types of connections (multiple networks), and I show how to use the square-root sparse group LASSO to detect which of the multiple connection types is more influential. Simulation evidence shows that my estimator has good finite sample performance. I further applied my method to the data in Banerjee et al. (2013) and my proposed procedure is able to identify leaders and effective networks.

key words: endogenous effects, spillovers, high-dimensional models, LASSO, model selection, robust inference,

^{*}I would like to thank Francesca Molinari, Matthew Backus, David Easley, Marten Wegkamp and Yanlei Ma. All remaining errors are mine.

[†]Cornell University, sp947@cornell.edu

1 Introduction

How an individual's behavior is affected by the behavior of her neighbors in an exogenously given network is an important research question in applied economics. With the increasing availability of detailed data documenting connections among individuals, spatial autoregression models (SARs) have been widely applied in empirical networks literature to estimate endogenous effects.

In SARs, an individual's behavior depends on the weighted average of other individuals' behaviors (see Anselin, 1988; Kelejian and Prucha, 1998). Standard SARs assume that the endogenous effects are the same across individuals in a network. Each individual influences her neighbors *at the same rate* regardless of who she is. However, in many contexts, some individuals are clearly more influential than others. For example, Mas and Moretti (2009) finds that the magnitude of spillovers varies dramatically among workers with different skill levels. Clark and Loheac (2007) also note that popular teenagers in a school have much stronger influence on their classmates' smoking decisions than their less popular peers.

I propose a novel SAR model which allows for *heterogeneous* endogenous effects. Each individual in a network simultaneously generates an outcome that takes into account all her neighbors' behaviors. Unlike standard SARs, each individual has an individual-specific effect on her neighbors. As a result, there are as many coefficients for individual-specific endogenous effects as there are individuals in the network. To achieve identification, I assume that "truly-influential" individuals only constitute a small fraction of the total population. In other words, individual-specific coefficients are assumed to be sparse. This assumption allows me to estimate the model via the least absolute shrinkage and selection operator (LASSO). The LASSO procedure penalizes the l_1 norm for the coefficients of heterogeneous endogenous effects. The geometry of the l_1 norm enforces the sparsity in the LASSO estimators. If a coefficient is selected by LASSO (i.e. the estimated coefficient is nonzero), the individual associated with this coefficient can influence all her neighbors at her specific rate. Otherwise the LASSO estimator will indicate that the individual has no influence on her neighbors. With some restrictions on the network structures, I show that the LASSO estimates for heterogeneous endogenous effects have near oracle performance (see Bühlmann and van de Geer, 2011). In other words, the selection of influential individuals is consistent and the convergence rate of non-zero LASSO estimates is the same as the convergence rate that would have been achieved if the truly influential individuals were known.

One challenge in my estimation process is the presence of endogeneity in the spatial lag and error term. As with standard SARs, the dependent variable in my model is used to construct spatial lags as an independent variable. As a result, the regressors are correlated with the error term and estimates would be biased if we were to apply LASSO directly.

First I propose a set of novel instruments to address the endogeneity. Following Kelejian and Prucha (1998), I express the dependent variable as an infinite sum of functions consisting of independent variables and an adjacency matrix. These functions are used as instruments. Then I design a two-stage estimation process for heterogeneous endogenous effects using LASSO at each stage. In the first stage, I use LASSO to estimate the coefficients for the instruments. These estimated coefficients and instruments are then used to create a synthetic dependent variable. In the second stage, I replace the dependent variable in the spatial lags with the synthetic variable to perform the LASSO estimation. Unlike in the standard two stage least square estimation process, the synthetic dependent variable in the first stage suffers from a shrinkage bias due to the LASSO fitting. However, I show that with certain restrictions on the network structure, the shrinkage bias is negligible (i.e. $o(1/\sqrt{n})$).

The next challenge is to construct robust confidence intervals for my LASSO type two-stage estimator. As pointed out in Leeb and Potscher (2008), it is impossible to construct uniformly valid confidence intervals for estimates based on model selection. Consistent model selection by LASSO is only guaranteed when all non-zero coefficients are large enough to be distinguished from zero in a finite sample (i.e. usually called the "beta-min" condition). LASSO may fail to select regressors with very small coefficients, resulting in omitted variable bias in the post LASSO inference.

I propose a bias correction for my two-stage estimator following the recent LASSO inference literature (see Belloni et al., 2015; van de Geer et al., 2014). The idea is to correct the first order bias and make the estimators independent from the model selection. Heuristically, shrinkage bias due to the l_1 penalty in LASSO can be expressed as a function of the LASSO estimators. Normality can still be achieved after adding back this bias. I show that this strategy also works in a twostage estimation process. I derive the asymptotic normality for my "de-sparse" two-stage LASSO estimator and conduct robust inference including confidence intervals.

My model can be extended to allow for more flexible network structures. One real world scenario is a network which consists of multiple cliques. Each clique has its local leaders, who only influence individuals within their own cliques but have no influence on individuals outside their cliques. One identification difficulty in this setting is that the number of leaders increases with the number of cliques. Hence, the sparsity assumption can potentially be violated.

To solve the problem in this scenario, I modify my model by bringing back the classical SAR model. I assume that there are both local leaders and global leaders in the network. In contrast to local leaders, global leaders can influence individuals across different cliques. I assume global leaders are sparse and show that identification can be achieved for this modified model. The endogenous effects of local leaders will be captured by the classical SAR model, which becomes an average endogenous effects in the network. The endogenous effects of global leaders, whose influence remains individualspecific, can be identified in the same way as it was in the previous model. If there is no global leader in the network, the model is effectively just the standard SAR model.

Another real world scenario is the existence of multiple types of connections among individuals. For example, connections among individuals can be classified as social (e.g. friendship, kinship) or economic (e.g. lending, employment). It is also important to identify which networks are more efficient at transmitting the endogenous effects.

I model different types of connections as multiple networks. I propose the use of square-root sparse group LASSO to estimate a heterogeneous endogenous effects model with multiple networks. The standard sparse group LASSO penalizes both the l_1 norm and the l_2 norm for each coefficient in each type of connection. I modify the sparse group LASSO by taking the square-root of the mean square error and thus make the estimation process pivotal. I derive the convergence rate and prove the consistency of selection. To the best of my knowledge, my paper is the first to show statistical properties for square-root sparse group LASSO.

I provide simulation evidence for networks of different sizes and different generating algorithms. The empirical coverage of my proposed estimators is close to the nominal level in all scenarios. Similar results are also found in models with multiple networks and with cliques.

I apply my method to study villagers' decisions to participate in micro-finance programs in rural areas of India as in Banerjee et al. (2013). Among different social and economic networks, my method shows that some networks such as "visit go-come" and "borrow money", are much more effective at influencing villagers' decisions than other networks such as "temple company" and "medical help". I further show that individuals in certain careers such as agricultural workers, Anganwadi teachers and small business owners are more likely to influence villagers.

My proposed methodology can be applied to detect influential individuals in empirical work when there are both leaders and followers. It is important to identify such individuals because we can then study why certain people are more influential than others. On the one hand, we can examine individuals exogenous characteristics and see if any of them contribute to an individual's influence. On the other hand, we can study how the position of an individual within a network may impact her influence by further introducing network formation into the model.

Being able to identify influential individuals could also lead to more effective policy outcomes. If

individuals with certain characteristics are found to be more influential than others, policy makers could potentially implement policies solely targeting influential individuals rather than the entire population. Since more resources are directed to the small group of highly influential individuals, one would expect much more effective policies. For example, online opinion leaders have influence on what people tweet and share on the Internet. In an election, instead of advertising on television and trying to influence every voter, a candidate could invest in these online opinion leaders and let them influence the public in a more efficient way. This technique could also work in employment contexts. Union leaders are often those workers who have the strongest influence on their fellow workers' opinions. Instead of reading through complaints from every worker, employers could identify those union leaders and make sure their complaints were addressed to prevent any ongoing strike. When studying peer effects in smoking behavior, my method can identity a group of teenagers who have a strong influence on their peers. A policy can target this group of students and encourage them to quit smoking.

1.1 Literature Review

This paper brings together literature on spatial autoregression model, LASSO and networks.

SARs:

SARs have been used widely applied in empirical studies. For instance, they have been used to study peer effects in labor productivity (see Mas and Moretti, 2009; Guryan et al., 2009; Bandiera et al., 2009), smoking behavior among teenagers (see Krauth, 2005; Clark and Loheac, 2007; Nakajima, 2007), educational achievements among different student groups (see Sacerdote, 2001; Neidell and Waldfogel, 2010), systemic risk in finance (see Bonaldi et al., 2015; Denbee et al., 2015), and the adoption of new agricultural technologies (see Coelli et al., 2002; Conley and Udry, 2010). My paper proposes a novel extension of standard SAR models that could be used to identify influential individuals in any given network. My methodology for estimating such a model could easily be adopted in existing empirical SAR analyses to identify influential individuals who influence their peers productivity, smoking decisions, or financial holdings.

More specifically, my model extends existing SARs literature by introducing *heterogeneous* endogenous effects. Until very recently, SARs always assume a constant rate of dependence for endogenous effects across different individuals (see Cliff and Ord (1973), the first monograph on the topic, and the later studies, Upton and Fingleton (1985); Anselin (1988); Cressie (1993); Lee and Liu (2010); Lee and Yu (2010); Jin and Lee (2016)). Recent developments in social interaction literature incorporate individual characteristics into SARs, essentially modeling the heterogeneity through a

linear combination of exogenous effects (see Manski, 1993; Bramoullé et al., 2009). In contrast, my model considers the heterogeneity in the endogenous effects. Heterogeneous endogenous effects can be identified from individuals' outcomes instead of being pre-specified through individuals' characteristics. To my knowledge, my proposed model is the first to capture the direct impact of an individuals neighbors' decisions on her own decision.

To estimate the heterogeneous endogenous effects in my model, I propose a methodology that is different from standard SARs literature. In classic SAR models, there is only one endogenous variable and hence it is sufficient to identify the model through only one instrument. In my model, the number of potentially endogenous variables increases as the number of observations increases. As a result, I propose a set of instruments that contain the same number of instruments as the total number of individuals. Moreover, each instrument is different from the standard SARs instrument as in Kelejian and Prucha (1998), Lee (2002), Lee (2003) and Lee (2004).

This paper also contributes to literature that models multiple networks through SARs. In standard SARs, multiple networks are modeled as higher order spatial lags (see Lee and Liu, 2010). Even though different networks are assumed to have different constant rates for endogenous effects in these models each individual in a given network faces the same constant rate. In contrast, my model allows for the a more realistic scenario where each individual has her own specific endogenous effects in each network. Moreover, my methodology allows some networks to be classified as completely irrelevant to decision-making *ex ante* and these networks can be consistently identified.

LASSO:

My paper extends LASSO literature by deriving statistical bounds and consistency of selection for the square-root sparse group LASSO estimator. This estimator builds on the group LASSO, squareroot LASSO, square-root group LASSO, and sparse group LASSO. Belloni et al. (2011) introduced the square-root LASSO, which does not require a pre-estimation of an unknown standard deviation σ . Yuan and Lin (2006) proposes the group LASSO, in which explanatory variables are represented by different groups. The group LASSO assumes that sparsity exists only among groups, i.e. some groups of variables are relevant while other groups are not. Simon et al. (2013) proposes the sparse group LASSO, which further allows sparsity within each group, i.e. some regressors within the relevant groups can also be irrelevant. citeBunea2013 derives statistical properties for the squareroot group LASSO, which combines group LASSO and square-root LASSO. When estimating a heterogeneous endogenous effects model with multiple networks, I provide proof for both statistical bounds and consistency of selection for the square-root sparse group LASSO estimator. To the best of my knowledge, this paper is the first to show asymptotic statistical properties for the square-root sparse group LASSO estimator. This paper also contributes to the growing literature on endogenous regressors in LASSO estimators. For instance, Belloni et al. (2014a) proposes the double selection mechanism to study confounded treatment effects. Fan and Liao (2014) proposes a GMM type estimator to deal with many endogenous regressors. Gautier and TsyBakov (2014) proposes a Self Tuning Instrumental Variables (STIV) estimator. The paper that is closest to mine is Zhu (2016), which studies the statistical properties of two-stage least square procedure with high-dimensional endogenous regressors. She studied a case when there exists p endogenous regressors. For each regressor j, she assumed that one can find d_j instruments. Both p and d_j may grow as n increases. I consider a case that is tailored to my SAR model. There are n endogenous regressors and each regressor shares the same n instruments. I show that a modified "de-sparse" LASSO estimator can be constructed for my estimator in a manner similar to Zhang and Zhang (2011), Bühlmann (2013), van de Geer et al. (2014) and Zhu (2016). I derive its asymptotic distributions and show how to perform inference.

Network:

My paper shares similar microfoundations with SARs as discussed in Blume et al. (2015), where the individual utility function can be written as a linear summation of the private and social components. The private component is a quadratic loss function on individual's efforts. The social components depend on the network structure as well as the efforts of one's neighbors. While the marginal rate of substitution between the private and social components of utility is assumed fixed in SARs, I assume this rate is individual-specific and depends on one's neighbors. My paper applies and extends LASSO approaches to deal with a high-dimensional problem in networks. The total number of possible edges in a network is n^2 , however, the social interaction networks we often observe are far more sparse. This is an ideal setting where LASSO could be applied. Manresa (2013) studies the heterogeneous exogenous effects in a network using LASSO. de Paula et al. (2015) explore the use of LASSO to recover network structures. Both these two papers consider panel data and rely on repeated observations of the same network to identify their models. My model considers cross-sectional data. To identify an individual's endogenous effects, I rely on the variations in her neighbors' outcome.

My paper also relates to the literature on identifying the key players in the network following Ballester et al. (2006), Calvó-Armengol et al. (2009), and Horracea et al. (2016). Under the framework of SARs, every individual is assumed to have the same endogenous effects. As a result, individuals who are well-connected in the network (with high centrality measure) become the key players in the network. However, this is not necessarily the case in my model, as well connected individuals can have zero endogenous effects on her neighbors. Indeed, as shown in the empirical application, well connected villagers such as tailors, hotel workers, veterans, and barbers are not influential in other villagers' decisions to join the micro-finance program. The rest of this paper is organized as follows: in Section 2, I introduce the model; in Section 3, I discuss identification assumptions; in Section 4, I design estimation procedures; in Section 5, I derive consistency and asymptotic properties; in Section 6, I show finite sample performance using Monte Carlo simulations; in Section 7, I apply my proposed model to study influential individuals and effective networks in promoting micro-finance programs in rural India; and in Section 8, I conclude.

2 Models

In this section, I first lay out the benchmark endogenous effects model and introduce the central model of this paper the heterogeneous endogenous effects model. Then I discuss two extensions of the heterogeneous endogenous effects model: a model for networks consisting of multiple cliques and a model for multiple networks. Finally, I provide two examples and illustrate how my model fits into these settings.

2.1 Benchmark Endogenous Effects Model

In this paper, I first introduce the standard spatial autoregression model (SAR) as the benchmark endogenous effects model. Let n denote the total number of observed individuals in a network. The outcome of individual i is denoted as d_i and is the variable of interest. Here d_i can represent any outcome associated with individual i, such as whether to smoke, whether to join a program, or whether to tweet a message from a friend. It is assumed that the outcome of each neighbor of individual i impacts her outcome homogeneously through a constant rate λ_0 :

$$d_i = \lambda_0 \sum_{j \in N_i} d_j + x_i \beta_0 + \epsilon_i, \tag{1}$$

where the set N_i is defined as individual *i*'s neighbors. The matrix form of this model is expressed as follows:

$$D_n = \lambda_0 M_n D_n + X_n \beta_0 + \epsilon_n, \tag{2}$$

where $D_n = (d_1, d_2, \dots, d_n)'$ is the *n*-dimensional vector of observable outcomes. The *n* by *k* matrix X_n represents the observable exogenous characteristics of individuals. When ϵ_n is specified as an *n*-dimensional vector of independent and identically distributed disturbances with zero mean and a constant variance σ^2 , equation (2) is also called a mixed regression model.

The spatial weight matrix M_n is of size n by n, where the (i, j)th entry represents the connection

between individual i and individual j. In empirical studies, the spatial weight matrix is often replaced by the adjacency matrix (see Ammermuller and Pischke, 2009; Acemoglu et al., 2012; Banerjee et al., 2013): the (i, j)-th entry of the matrix M_n takes value 1 if individual i and individual j are connected and takes value 0 otherwise.

In this model, endogenous effects (see Manski, 1993) or network effects (see Bramoullé et al., 2009) are captured by the scalar λ_0 . An implicit assumption in equation (2) is that λ_0 , the rate of endogenous effects, is identical across all individuals in the network. Every individual affects her neighbors at this same rate λ_0 no matter who she is, how many neighbors she has and where is she in the network. This limitation has been noted in various studies (see Ammermuller and Pischke, 2009; de Paula et al., 2015). I relax this assumption by proposing a more flexible model that allows individual-specific endogenous effects as discussed below.

2.2 Heterogeneous Endogenous Effects Model

I propose the following model to allow for heterogeneous endogenous effects:

$$d_i = \sum_{j \in N_i} d_j \eta_j + x_i \beta + \epsilon_i \tag{3}$$

where N_i represents the set of individual *i*'s neighbors and η_j represents the endogenous effects of individual *j* on the outcome of all her neighbors $i \in N_j$. the model can be rewritten in matrix form as:

$$D_n = \left(M_n \circ D_n\right)\eta_0 + X_n\beta_0 + \epsilon_n,\tag{4}$$

where $\eta_0 = (\eta_1, \eta_2, \dots, \eta_n)'$ is a vector of parameter of size n by 1. The *i*th entry in η_0 represents the endogenous effects of individual i on her neighbors. This model allows for individual heterogeneity to interact with endogenous effects so that every individual is allowed to have her own coefficient η_i . My model allows some $\eta_j = 0$. In other word, there are individuals that have no endogenous effects on their neighbors. I define those individuals with $\eta_j \neq 0$ as influential.

The operator \circ is defined between a *n* by *n* matrix M_n and a *n* by 1 vector D_n as

$$M_n \circ D_n = M_n \cdot \operatorname{diag}(D_n) = C_n$$

where $diag(\cdot)$ is the diagonalization operator and $C_{i,j} = M_{i,j}d_j$.

Note that in contrast to fixed rate λ_0 specified in equation (2), even though each neighbor of

individual j is assumed to receive the same influence $d_j\eta_j$ from her¹, each individual is allowed to influence her neighbors at her own rate η_j .

Similar to equation (2), equation (4) can be derived from a bayesian Nash Equilibrium. Let (x_i, ϵ_i) denotes an individual's type, where x_i is publicly observed characteristics and ϵ_i is private characteristics only observable by *i*. Individual *i*'s utility depends on her own action and characteristics as well as her neighbors' actions. Individual *i* chooses action d_i to maximize the following utility:

$$U_{i}(d_{i}, d_{-i}) = (x_{i}\beta + \epsilon_{i})d_{i} - \frac{1}{2}d_{i}^{2} + \sum_{j \in N_{i}} d_{j}d_{i}\eta_{j}$$

The first order condition yields equation (4)

2.3 Examples

To help readers conceptualize the heterogeneous endogenous effects model, here I apply the model to two specific contexts one invloving labor productivity and the other online opinion leaders.

• Peer Effects in Labor Productivity

Understanding the mechanism and magnitude of the dependence of labor productivity on coworkers is an important question for economists and policy makers. As found in Mas and Moretti (2009), workers respond more to the presence of coworkers with whom they frequently interact. In this case, the influence level of each individual to hers coworkers is not necessarily the same. Equation (4) can be used to incorporate such differences.

$$y_i = \sum_{j \in N_i} y_j \eta_j + x_i \beta + \epsilon_i,$$

where y_i is individual *i*'s productivity, x_i represents individual *i*'s characteristics (education levels, ages, etc) and N_i is the set of coworkers that works directly with *i*. η_j represents the size of influence of coworker j – all else being equal, the additional effect on individual *i*'s productivity if individual *j* becomes her coworker

Note that if we restrict the parameters η_j to be the same across different workers, then we are back to the classical SAR setting as laid out in equation (2). Thus, $\lambda = \frac{1}{n} \sum_{j=1}^{n} \eta_j$ can be interpreted as the averaged spillover effects in the canonical sense.

¹Further relaxation of the model considering different individual j's influence on each of her neighbors requires panel data.

Define

$$\lambda^* = \frac{1}{\sum \mathbf{1}_{\eta_j \neq 0}} \sum_{j=1}^n \eta_j \mathbf{1}_{\eta_j \neq 0}$$

as the averaged endogenous effects for influential workers. λ^* does not include non-influential individuals in the calculation. It is a more precise measure of endogenous effects compared with λ from equation (2).

• Online Opinion Leaders

A decision can represent whether to "tweet" a news story seen online. When individuals make such decisions, they are often influenced by several online opinion leaders – whether those people "tweet" the news or not. There are also many types of online opinion leaders, including political figures and some are celebrities. For certain types of news, some opinion leaders may be very influential while the rest may have no influence on the public. Opinion leaders may also influence each other when deciding whether to "tweet" the news or not. Assume a binary decision (0, 1) is made from a bayesian Nash Equilibrium, such that

$$d_i^* = \sum_{j \in N_i} d_j^* \eta_j + x_i \beta + \epsilon_i,$$

where d_i^* is the probability of individual *i* playing action 1, and $\sum_{j \in N_j} d_j^* \eta_j$ is the expected endogenous effects from *i*'s neighbors N_i . X_i is the individual *i*'s characteristics such as political views, age, career, etc.

Similarly, we can define $\lambda = \frac{1}{n} \sum_{j=1}^{n} \eta_j$ as the averaged endogenous effects. Since the number of opinion leaders is very small compared with total online users, λ can be very close to 0. A more precise measure would be

$$\lambda^* = \frac{1}{\sum \mathbf{1}_{\eta_j \neq 0}} \sum_{j=1}^n \eta_j \mathbf{1}_{\eta_j \neq 0}$$

 λ^* will be the average endogenous effects for online opinion leaders. On the other hand, it is also important to identify the set:

$$S = \{j : \eta_j \neq 0\}$$

as truly influential opinion leaders. If a similar type of news story needs to be spread the next time, contacting those leaders and obtaining their endorsement would be a good starting strategy.

2.4 Heterogeneous Endogenous Effects Model with Cliques

I propose an extension to my heterogeneous endogenous effects model which could address such challenges. Consider a network composed of many cliques (small groups of connected individuals).



Figure 1: Local Leader

Each clique has its local leader who only influences individuals within her own clique. Figure 1 provides an example of such a network structure. Note that in Figure 1, node S_2 , S_3 and S_4 represent local leaders who only influence individuals within their own cliques. On the contrary, node S_1 represents a global leader who can influence individuals across different cliques. For example, one can think about the local leaders S_2 , S_3 and S_4 as local news channels while S_1 is the national news channel. I assume that all local news will influence the public at a small but similar rate while different national channels can have different effects on their audience.

In the above network structure, if the number of local leaders is increasing with the number of cliques but the number of individuals in each clique stays fixed, it is impossible to identify the individualspecific influence of all those local leaders. To address this problem, I assume a homogeneous effect γ_0 among all individuals. This rate will capture all influence from local leaders. However, I allow global leaders to heterogeneously influence their neighbors at rates that differ from γ_0 and show that γ_0 and the heterogeneous effects can be consistently estimated.

More specifically, I consider the following model:

$$d_i = \sum_{j \in N_i} d_j \eta_j + \gamma_0 \sum_{j \in N_i} d_j + x_i \beta_0 + \epsilon_i,$$
(5)

which be represented in matrix form as:

$$D_n = \left(M_n \circ D_n\right)\eta_0 + M_n D_n \gamma_0 + X_n \beta_0 + \epsilon_n,\tag{6}$$

where $\eta'_0 = (\eta_1, \eta_2, \dots, \eta_n)'$. The new term $\gamma_0 \sum_{j \in N_i} d_j$ captures influence from the local level. Note that this is the same term as the spatial lag in the benchmark SAR model. The vector η_0 captures the heterogeneous endogenous effects of global leaders.

If no global leader exists, i.e. $\eta_j = 0, \forall j$, the model collapses back to the classical SAR model as

in equation (2). If there is no local level influence, i.e. $\gamma_0 = 0$, then the model coincides with the heterogeneous endogenous effects model in section 2.2.

2.5 Heterogeneous Endogenous Effects Model with Multiple Networks

In reality, individuals are often connected with each other through more than one type of network. For example, ones colleague (connection in an employment network) could also be her friend (connection in a friendship network), and ones uncle (connection in a relative network) could also be the person she lends money to (connection in a borrowing/lending network). In such scenarios, an individuals outcome could potentially be influenced by the outcomes of her neighbors from more than one type of network.

To capture different types of connections among the same set of individuals, we can incorporate multiple networks in my heterogeneous endogenous model. More specifically, a separate adjacency matrix can be constructed for each type of network. For instance, the (i, j)-th entry of the adjacency matrix representing friendship takes value 1 if individual i and individual j are friends and takes value 0 otherwise; that representing the borrowing/lending network takes value 1 if individual i and individual j lend money to each other and takes value 0 otherwise.

Let q be the total number of different types of network. Define M_n^l as the adjacency matrix for the *l*th network. The heterogeneous endogenous effects model with multiple networks is defined as

$$d_i = \sum_{l=1}^{q} \sum_{k \in N_i} d_k^l \eta_k^l + x_i \beta_0 + \epsilon_i \tag{7}$$

Note that in this model, different network could potentially bear different endogenous effects for the same individual. In equation (7), coefficient η_k^l represents the rate of endogenous effect of individual k through network l. As a result, we have nq + k coefficients for endogenous effects. In addition, I assume endogenous effects from different types of networks are linearly additive. The model can also be rewritten in matrix form as:

$$D_n = \sum_{l=1}^q \left(M_n^l \circ D_n \right) \eta_0^l + X_n \beta_0 + \epsilon_n, \tag{8}$$

where M_n^l is the adjacency matrix for network l. $\eta^l = (\eta_1^l, \eta_2^l, \dots, \eta_n^l)'$ is an n by 1 vector for $l = 1, 2, \dots, q$. Define a network l as efficient network if $\eta_i^l \neq 0$ for at least one individual $i = 1, 2, \dots, n$.

3 Identification

In this section, I discuss the conditions under which the heterogeneous endogenous effects model is identified and the extensions of this model. My assumptions combine both standard SARs type assumptions and LASSO type assumptions. SARs type assumptions ensure the existence of valid instruments to identify the model. LASSO type assumptions guarantee consistent model selection and estimation using a LASSO estimator. In what follows, I will first present the assumptions needed for a standard heterogeneous endogenous effects model to be identified. Then I will discuss identification assumptions for two model extensions laid out in the previous section – one heterogeneous endogenous effects model for networks consisting of multiple cliques and one with multiple types of networks.

Before discussing identification assumptions for the heterogeneous endogenous effects model, lets first recall the benchmark SAR model:

$$D_n = \lambda_0 M_n D_n + X_n \beta_0 + \epsilon_n, \tag{9}$$

Note that by rearranging the above equation, we can express endogenous variable $M_n D_n$ solely as a function of X_n and M_n , since:

$$D_n = J_n^{-1} X_n \beta_0 + J_n^{-1} \epsilon_n$$

where I_n is the *n* by *n* identity matrix and $J_n = I_n - \lambda_0 M_n$. It is straightforward that $J_n^{-1}X_n$ can serve as valid instruments for $M_n D_n$. As a result, the identification and estimation of equation (9) can be achieved through either 2SLS or GMM as proposed in papers such as Kelejian and Prucha (1995), Kelejian and Prucha (1998) Lee (2002), Lee (2003), and Lee (2004).

As will be explained in detail in subsequent sections, to estimate the individual specific effects in the heterogeneous endogenous effects model, I derive a set of instruments in a similar way by solving D_n as a function of exogenous variables and an adjacency matrix. The assumptions listed below essentially guarantee the existence and consistency of the 2SLS estimates.

3.1 Identification Assumptions for the Heterogeneous Endogenous Effects Model

Recall that the heterogeneous endogenous effects model is specified as:

$$D_n = \left(M_n \circ D_n\right)\eta_0 + X_n\beta_0 + \epsilon_n,$$

First note that without additional restrictions, this model could not be point identified through canonical method as the number of parameters n + k is greater than the number of observations n. To achieve identification, the key assumption that I maintain is that only a small number of individuals in the network are influential (i.e. $\eta_j \neq 0$).

Assumption 1. Let $S_n \subset \{1, 2, \dots, n\}$ denote the set of influential individuals (i.e. $\eta_j \neq 0$). Let $s_n = |S_n|$ be the number of elements in S_n .

$$s_n = o\left(\frac{\sqrt{n}}{\log n}\right), \quad \text{as } n \to \infty$$

Assumption 1 is usually referred to as "sparsity" assumption. The assumption that most individuals in a network are not influential is plausible under many circumstances. For example, opinion leaders on social media only constitute a very small fraction of internet users; there are only a couple of "cool" kids at school that might influence their friends' smoking decisions; passionate workers that can boost the productivity of their coworkers are also relatively rare. When many local leaders exist within a network, the sparsity assumption could be violated. I will address this issue in section 3.2 and show that identification can still be achieved with additional assumptions.

Assumption 2.

- There exists an $\eta_{\max} < 1$ such that $\|\eta_0\|_{\infty} \leq \eta_{\max}$
- The ϵ_i are *i.i.d* with 0 mean and variance σ^2
- The regressors x_i in X_n are non-stochastic and uniformly bounded for all n. $\lim_{n\to\infty} X'_n X_n/n$ exists and is nonsingular

Assumption 2 guarantees the invertibility of $(I_n - M_n \circ \eta_0)$. The restriction on η_0 excludes the unit root process and ensures the uniqueness of equilibrium. The assumptions on the error term and the assumption that X_n is a fixed design matrix are the same as those imposed in the mixed regression model² (see Lee, 2002). I focus on the case where X_n is an n by 1 vector and study identification as in Bramoullé et al. (2009). It is straightforward to generalize the algebra when X_n is n by k. More instruments can be constructed in this scenario.

To proceed, recall the definition of the operator " \circ " as $M_n \circ D_n = M_n \cdot \operatorname{diag}(D_n)$, where $\operatorname{diag}(\cdot)$ is the diagonalization operator. Note the following property of the " \circ ":

$$(M_n \circ D_n)\eta_0 = (M_n \circ \eta_0)D_n$$

 $^{^{2}}$ The assumption on error terms exclude exogenous effects and correlated effects from my model. An identification problem similar to the "reflection problem" arises when including exogenous effects. More instruments need to be constructed, which requires better data. These are interesting directions for future research.



Figure 2: Examples of networks which violate assumption 3

If the invertibility of $(I_n - M_n \circ \eta_0)$ is guaranteed, then

$$D_n = \left(M_n \circ D_n\right)\eta_0 + X_n\beta_0 + \epsilon_n \Leftrightarrow D_n = \sum_{i=0}^{\infty} \left(M_n \circ \eta_0\right)^i (X_n\beta_0 + \epsilon_n) \tag{10}$$

This is formally shown in Appendix B.

Since $(M_n \circ D_n)\eta_0$ is correlated with ϵ_n and η_0 is sparse (i.e. having at most s_n non-zero elements), we need at least s_n instruments to deal with the endogeneity in the model. Using equation (10), we can express the expectation of D_n as follows:

$$E(D_n) = X_n \beta_0 + \left(M_n \circ X_n \right) (\beta \eta_0) + \sum_{i=2}^{\infty} \left(M_n \circ \eta_0 \right)^i \beta_0 X_n, \tag{11}$$

Let $(\cdot)_S$ denote the operator such that $(M_n)_S$ is a sub matrix of M_n with its columns restricted to columns corresponding to the elements of S. The first and second terms of equation (11) suggest that X_n and $(M_n \circ X_n)_S$ can serve as valid instruments to point identify β_0 and η_0 .

Assumption 3. $[X_n, (M_n \circ X_n)_S]$ is full rank.

Assumption 3 is the key assumption that leads to identification. The linear independence among $(M_n \circ X_n)_S$ requires the assumption that any two influential individuals may not necessarily connect with identical neighbors. Moreover, assumption 3 also requires that neighbors of an influential individual cannot be a linear combination of neighbors of several other influential individuals, which rules out network structures as depicted in Figure 2:



Figure 3: Fixed Effects

In other words, as long as each influential individual has a neighbor that is not connected with any other influential individuals, assumption 3 is satisfied. One can think of the identification here as estimating fixed effects from influential individuals. Collinearity arises when the fixed effects of two influential individuals are imposed on exactly the same observations. As shown in Figure 3, the influence of S_1 can be identified by comparing red and yellow groups, while the influence of S_2 can be identified by comparing blue and black groups. Or the influence of S_1 can be identified by comparing green and blue groups, while the influence of S_2 can be identified by comparing red and green groups.

Further, as shown in Appendix B, one can rewrite equation (11) as:

$$E(D_n) = X_n \beta_0 + \left(M_n \circ X_n \right) \tilde{\eta}, \tag{12}$$

where $\tilde{\eta}_j = \eta_j f(\beta_0, X_n, M_n)$ for some function f depends on β_0, X_n , and M_n . Note that $\tilde{\eta}_j = 0$ as long as $\eta_j = 0$. As a result, the sparsity assumption is also satisfied in equation (12), and I can thus estimate equation (12) as the first stage in using a LASSO type estimator.

At this point, if the truly influential individuals set S_n were available to us, we would be able to estimate the model using 2SLS method or GMM. However, in most cases, S_n is not known beforehand. I propose to use a LASSO type estimator to both recover the set of influential individuals and estimate the model. For LASSO to achieve correct recovery, I need the following assumptions:

Assumption 4.

(Irrepresentable Condition) There exists $N \in \mathbb{N}$: $\forall n \geq N$, there is a $\vartheta \in (0,1)$ such that

$$P\left(\left\|diag((\hat{D}_n)_{S^c})\Sigma_n diag((\hat{D}_n)_S)^{-1}sign(\eta_0)\right\|_{\infty} \le \vartheta\right) = 1;$$

where

$$\Sigma_n = (M_n)'_{S^c}(M_n)_S \left((M_n)'_S(M_n)_S \right)^{-1},$$

(Beta Min Condition) There exists $N \in \mathbb{N}$: $\forall n \geq N$, there is a m > 0 such that

$$\min(|\eta_0|)_S \ge m/\sqrt{n},$$

Here $(M_n)_S$ represents the sub-matrix of M_n given by the columns corresponding to influential individuals. Similarly, $(M_n)_{S^c}$ represents the sub-matrix of M_n given by the columns corresponding to non-influential individuals.

Assumption 4 is required for the LASSO estimator to achieve a consistent selection for the set S_n in the second stage. The Irrepresentable Condition imposes restrictions on non-influential individuals such that the neighbors of a non-influential individual will not be exactly the same as those of any influential individual. This is because when two individuals connect with exactly the same neighbors, we cannot distinguish which individual is the true source of influence. This assumption rules out identification in complete networks (i.e. all individuals are connected). The Beta Min Condition requires the magnitude of the endogenous effects to be sufficiently strong in order to be detected by LASSO. As shown in Zhao and Yu (2006), the Irrepresentable Condition together with the Beta Min Condition are necessary and sufficient conditions for LASSO to achieve consistent model selection. If consistent selection is not required, these two conditions can also be relaxed to weaker conditions (such as the compatibility condition as in Bühlmann and van de Geer (2011)). As shown in van de Geer et al. (2014), with the compatibility condition, inference on the de-sparse coefficients as discussed in the next section is still valid.

Assumption 5.

(Maximum Neighbors Condition) There exists $N \in \mathbb{N}$: $\forall n \geq N$,

$$\|M'_n \mathbf{1}_n\|_{\infty} \le O(\log n),$$

holds with probability equal to 1

(Variance Condition)

$$\frac{1}{n}M'_nW_n(I-M_n\circ\eta_0)^{-1}(I-M_n\circ\eta_0)^{-1'}W_nM_n\to\Omega,$$

where $W_n=\left(I-X_n(X'_nX_n)^-X'_n\right)$

The Maximum Neighbors Condition requires the network structure (edges) to be sparse. More specifically, it requires that the number of direct neighbors not increase faster than $O(\log n)$ when the number of influential individuals increases at speed $o\left(\frac{\sqrt{n}}{\log n}\right)$. This rate can be improved when

the number of influential individuals is fixed. The Maximum Neighbors Condition is an asymptotic

bound on the number of neighbors for each individual as the network increases. This condition is required to prevent shrinkage bias carried from the first stage LASSO estimation from growing faster than $o(1/\sqrt{n})$ in the second stage.

The Variance Condition requires the variance-covariance matrix to converge to a limit. In classical SARs, the spatial weight matrix is assumed to be uniformly row sum bounded. This assumption implies the Variance Condition but imposes restrictions on the network structure. Each individual may only connect with a finite number of neighbors. In my case, the identification of an influential individual comes from the difference in responses between neighbors that solely connect with her and individuals who connect with no influential individuals. For example, consider two groups of individuals that have the same characteristic X where one group all connects with individual j and the other does not. If the mean response of the two groups is significantly different, we can conclude that j is influential. To identify the influence of individual j as a fixed effect, the number of individuals affected by individual j must grow as the sample size increases. As a result, the row sum for influential individuals cannot be bounded by a fixed number.

The heterogeneous endogenous effects model is identified under assumptions 1-5 as a linear system with a unique solution. I discuss the identification of my model with cliques and with multiple networks in the following two sections.

3.2 Identification Assumptions with Cliques

Recall the heterogeneous endogenous effects model with cliques, represented as follows:

$$D_n = \left(M_n \circ D_n\right)\eta_0 + M_n D_n \gamma_0 + X_n \beta_0 + \epsilon_n$$

Define global leaders as those influential individuals who influence multiple cliques and whose neighborhoods increase as n increases. Define local leaders as influential individuals who are not global leaders.

Assumption 1'. Among n individuals in the network, let $S_n \subset \{1, 2, \dots, n\}$ be the set of global leaders. Let $s_n = |S_n|$ be the number of elements in S_n . Assume:

$$s_n = o\left(\frac{\sqrt{n}}{\log n}\right), \quad as \ n \to \infty$$

Assumption 1' only requires the number of global leaders to be sparse. My model does not impose any restriction on the number of local leaders. As a result, it does not rule out situations where everyone is (locally) influential. Local leaders' influence will be captured by the γ_0 , coefficient of classical spatial lag.

To ensure invertibility of the matrix $(I_n - M_n \circ \eta_0 - M_n \gamma_0)$, I modify the first part of assumption 2 as:

Assumption 2'. There exists an $\eta_{\max} < 1$ such that $\|\eta_0 + \gamma_0\|_{\infty} \leq \eta_{\max}$

Similar to assumption 2, this assumption excludes unit root processes. Since there exists a local level influence γ_0 in the network, global level influence η_0 needs to be further bounded above by 1. As a result, equation (6) can be transformed into the following:

$$E(D_n) = X_n \beta_0 + \left(M_n \circ X_n\right)(\beta_0 \eta_0) + M_n X_n(\beta_0 \gamma) + \sum_{i=2}^{\infty} \left(M_n \circ \eta_0 + \gamma M_n\right)^i \beta_0 X_n$$

Equation (6) introduces one more coefficient γ_0 compared with equation (4). As a result, assumption 3 is modified to include an extra instrument $M_n X_n$, which is also the classic instrument used in equation (2):

Assumption 3'. $[X_n, (M_n \circ X_n)_S, M_n X_n]$ is full rank.

Assumption 3' is similar to assumption 3 and requires the additional instrument $M_n X_n$ to be linearly independent with $[X_n, (M_n \circ X_n)_S]$. The remaining assumptions 4 and 5 are unchanged.

3.3 Identification Assumptions with Multiple Networks

Recall the heterogeneous endogenous effects model with multiple networks, represented as follows:

$$D_n = \sum_{j=1}^q \left(M_n^j \circ D_n \right) \eta_0^j + X_n \beta_0 + \epsilon_n$$

First notice that the number of coefficients in this model becomes nq + k. The number of observed networks q is also allowed to increase as the number of observations increases. As a result, the sparsity assumption will be imposed on both the influential individuals and the effective networks. I assume that some of the networks are completely irrelevant (i.e. $\eta_0^j = 0$) and that relevant networks are not necessarily passing influence for everyone (i.e. $\eta_0^j \neq 0$ but $\eta_{0,i}^j = 0$ for some i).

Second, to ensure invertibility, for any matrix norm $\|.\|$:

$$\left\|\sum_{j=1}^{q} \left(M_n^j \circ \eta_0^j\right)\right\| \le \sum_{j=1}^{q} \left\| \left(M_n^j \circ \eta_0^j\right) \right\| \le \sum_{j=1}^{q} \|\eta_0^j\|_{\infty} \left\| \left(M_n^j\right) \right\|$$

Because M_n^j is the adjacency matrix such that each entry is 0 or 1, $\sum_{j=1}^q \|\eta_0^j\|_{\infty} < 1$ guarantees the invertibility of $I - \sum_{j=1}^q \left(M_n^j \circ \eta_0^j\right)$.

Third, I require $\left[X_n, \left(M_n^1 \circ X_n\right)_S, \left(M_n^2 \circ X_n\right)_S, \cdots, \left(M_n^q \circ X_n\right)_S\right]$ to be full rank. Compared with the standard model, this assumption requires the independence condition to hold across different networks. Again, we cannot identify the source of influence if two influential individuals connect to the same neighbors. Fourth, I assume conditions that guarantee a consistent selection of square-root sparse group LASSO. And, finally, the Maximum Neighbor Condition needs to be satisfied in all q adjacency matrices. Since the five conditions for multiple networks are very similar to assumption 1-5, I list them formally in the appendix as assumption 1*-5*.

4 Estimation

I propose an estimator similar to the two-stage least square method but use LASSO in both stages. The estimator proposed here is differs from the "double selection" estimator proposed in Belloni et al. (2014a) as I plugin the fitting from the first stage directly to the second stage. It is in the same framework as that proposed in Zhu (2016). I call this estimator a two-stage LASSO (2SLSS) estimator. In this section, I define this 2SLSS procedure and propose a bias corrected version of the estimator. I show how this procedure can be extended to estimate my model for networks consisting of multiple cliques and my model for multiple networks.

4.1 Two-Stage LASSO Estimator

I propose to estimate equation (4) using the following estimator:

Two-Stage LASSO Estimator:

• First Stage:

$$(\tilde{\beta}, \tilde{\eta}) = \arg\min_{\beta, \eta} \|D_n - X_n\beta - (M_n \circ X_n)\eta\|_2 + \lambda |\eta|_1$$

Obtain a LASSO fitting \hat{D}_n

$$\hat{D}_n = X_n \tilde{\beta} + \left(M_n \circ X_n \right) \tilde{\eta}$$

• Second Stage:

$$(\hat{\beta}, \hat{\eta}) = \arg\min_{\beta, \eta} \|D_n - (M_n \circ \hat{D}_n)\eta - X_n\beta\|_2 + \lambda |\eta|_1$$

As shown in section 3, $(M_n \circ D_n)$ is correlated with ϵ_n . Thus equation (4), equation (6) and equation (8) cannot be estimated directly using LASSO or sparse group LASSO. The instruments proposed in section 3 are $[X_n, (M_n \circ X_n)_S]$. We do not observe the set S but note that $[X_n, (M_n \circ X_n)]$ is a set of regressors that contains the valid instruments.

The two-stage least square method can be used to address endogeneity in SARs as in Lee (2003). In the first stage, $M_n X_n$ are used as instruments to estimate D_n . In the second stage, $M_n \hat{D}_n$ is used to replace $M_n D_n$ to avoid endogeneity.

Following the same idea, I estimate a first stage using $[X_n, (M_n \circ X_n)]$. Since there are n + k regressors, I use the square-root LASSO to select those instruments in set S. I choose the square-root LASSO over standard LASSO to avoid a pre-estimation of the unknown variance of the error term σ^2 . I construct a synthetic \hat{D}_n variable using square-root LASSO estimates. In the second stage, I replace D_n with \hat{D}_n in the regressors and estimate the coefficients $\hat{\eta}$ using the square-root LASSO again.

The statistical properties of two-stage estimators using LASSO have been studied in Zhu (2016), where she derives bounds for the estimator and proves consistency of model selection in a general setting. Zhu (2016) studied the over identified case where the number of endogenous regressors goes to infinity while the number of instruments for each regressor also goes to infinity. I studied the just identified case using the instruments proposed in section 3, where the number of endogenous regressors is the same as the number of instruments and both go to infinity.

4.2 De-sparse 2SLSS Estimator

The estimator $\hat{\beta}$ and $\hat{\eta}$ suffers from LASSO shrinkage bias. Moreover, post model selection inference conditioning on the selected model $\hat{S}_n = \{i | \hat{\eta} \neq 0\}$ suffers from the omitted variable bias and thus is not uniformly valid. (see Leeb and Potscher, 2005, 2008, 2009). I construct a "de-sparse" estimator under my setting and derive the asymptotic distribution for it. I propose the following de-sparse LASSO estimator:

De-sparse 2SLSS Estimator:

• Define

$$\hat{e} = \hat{\eta} + \hat{\Theta}(M_n \circ \hat{D}_n)'(D_n - X_n\hat{\beta} - (M_n \circ \hat{D}_n)\hat{\eta})/n$$

• Define

$$\hat{b} = \hat{\beta} - (X'_n X_n)^- X'_n (M_n \circ \hat{D}_n)' (D_n - X_n \hat{\beta} - (M_n \circ \hat{D}_n) \hat{\eta}) / n$$

 $\hat{\beta}$ and $\hat{\eta}$ are estimators from the 2SLSS. $\hat{\Theta}$ is defined by the nodewise regression as in Meinshausen and Bühlmann (2006). Nodewise regression explores the correlation between the columns of the design matrix $W_n(M_n \circ \hat{D}_n)$ by regressing each column on all the rest of the columns while penalizing the coefficients. An approximation of the inverse of the matrix $\frac{1}{n}(M_n \circ \hat{D}_n)'W_n(M_n \circ \hat{D}_n)$ can be constructed based on nodewise regression. Further, define $\hat{S}_n = \{i | \hat{\eta} \neq 0\}$, which represents the LASSO selected active set. The estimators (\hat{e}, \hat{b}) are adjusted for the LASSO shrinkage bias and are a consistent estimator for β and η . They are similar to the estimators proposed in van de Geer et al. (2014), but are constructed through a two-stage process as well as using square-root LASSO.

The de-sparse LASSO estimator does not depend on the selected active set. Thus, it does not suffer from the non-uniformity problem. Notice that the double selection method proposed in Belloni et al. (2014a) could also be applied to conduct inference on $\hat{\beta}$. Belloni et al. (2014b) shows the first order equivalence of the double selection method and the de-sparse method. On the other hand, the main interest of this paper is the coefficients $\hat{\eta}$. The double selection method does not provide a way to conduct inference for all the coefficients in the model, while the de-sparse LASSO estimator does.

My de-sparse LASSO estimator differs from the one proposed in Zhu (2016). Since the instruments are known in my case, I can derive the asymptotics for my estimator explicitly. By considering a sparse network structure (e.g. Maximum Neighbors Condition), I can show that the shrinkage bias from the first stage is negligible $(o(1/\sqrt{n}))$. The estimator proposed in Zhu (2016) adjusts shrinkage bias from both the first and second stages. In order to show consistency, she assumes the convergence of the product between the residual of nodewise regression and the endogenous regressors.

I will defer the proof of consistency for the LASSO selected set \hat{S}_n and consistency and asymptotic distribution for my estimator (\hat{e}, \hat{b}) to section 5. In the remainder of this subsection, I will define

the estimators for the two extended models.

4.3 2SLSS with Cliques

To estimate equation (6), I propose the following 2SLSS:

Two-Stage LASSO Estimator with Homogenous Effects:

• First Stage:

$$(\tilde{\beta}, \tilde{\gamma}, \tilde{\eta}) = \arg\min_{\beta, \gamma, \eta} \|D_n - X_n\beta - M_nX_n\gamma - (M_n \circ X_n)\eta\|_2 + \lambda(|\eta|_1 + |\gamma|)$$

Obtain a LASSO fitting \hat{D}_n

$$\hat{D}_n = X_n \tilde{\beta} + M_n X_n \tilde{\gamma} + \left(M_n \circ X_n \right) \tilde{\eta}$$

• Second Stage:

$$(\hat{\beta}, \hat{\gamma}, \hat{\eta}) = \arg\min_{\beta, \gamma, \eta} \|D_n - M_n \hat{D}_n \gamma - (M_n \circ \hat{D}_n) \eta - X_n \beta\|_2 + \lambda(|\eta|_1 + |\gamma|)$$

The estimator is similar to that for the previous model except that the classical spatial lag $M_n X_n$ is now included in the estimation. In the above estimator, I penalize η s and γ at the same rate because I have no prior knowledge of these two effects. One can penalize them at a different rate or not penalize γ if one believes that influence from local leaders is more likely than that from global leaders or vice versa. Since γ and η s are both penalized coefficients, a similar de-sparse LASSO estimator can be constructed for γ :

De-sparse 2SLSS Estimator with Cliques:

• Define

$$\hat{r} = \hat{\gamma} + \hat{\Theta}(M_n \hat{D}_n)' (D_n - X_n \hat{\beta} - M_n \hat{D}_n \tilde{\gamma} - (M_n \circ \hat{D}_n) \hat{\eta}) / n$$

Note that $\hat{\Theta}$ should be an approximation for the inverse of the matrix $\frac{1}{n}[M_n\hat{D}_n, (M_n\circ\hat{D}_n)]'W_n[M_n\hat{D}_n, (M_n\circ\hat{D}_n)]$ in this case.

4.4 Multiple Networks

When multiple networks exist, each individual will have network-specific endogenous effects. The number of unknown coefficients increases from n + k to nq + k compared with the standard case. These coefficients can also be classified into q different groups based on networks. By applying the sparsity assumption to the relevant networks, we can estimate the model using the square-root sparse group LASSO instead of the square-root LASSO and propose the following estimator. The square-root sparse group LASSO penalizes both the l_1 and l_2 norm in each group. It can identify all the relevant groups under weaker assumptions compared with the square-root LASSO estimator.

Two-Stage LASSO Estimator with Multiple Networks:

• First Stage:

$$(\tilde{\beta}, \tilde{\eta}) = \arg\min_{\beta, \eta} \left\{ \left\| D_n - X_n \beta - \sum_{j=1}^q (M_n^j \circ X_n) \eta^j \right\|_2 + \left(\sum_{j=1}^q \left(\lambda_1 \|\eta^j\|_2 + \lambda_2 \|\eta^j\|_1 \right) \right) \right\}$$

Obtain a LASSO fitting \hat{D}_n

$$\hat{D}_n = X_n \tilde{\beta} + \sum_{j=1}^q (M_n^j \circ X_n) \tilde{\eta}^j$$

• Second Stage:

$$(\hat{\beta}, \hat{\eta}) = \arg\min_{\beta, \eta} \left\{ \left\| D_n - X_n \beta - \sum_{j=1}^q (M_n^j \circ \hat{D}_n) \eta^j \right\|_2 + \left(\sum_{j=1}^q \left(\lambda_1 \| \eta^j \|_2 + \lambda_2 \| \eta^j \|_1 \right) \right) \right\}$$

The square-root sparse group LASSO introduces two tuning parameters, λ_1 and λ_2 , to penalize both the l_1 and the l_2 norm in each network. Similar to the LASSO estimator, the geometric shape of the penalties allows the square-root sparse group LASSO to identify sparsity not only within each network (group) but also among networks (groups). In other words, some networks could be completely irrelevant (i.e. $\eta^j = 0$) and within relevant networks, some individuals can have no influence on their neighbors (i.e. $\eta^j \neq 0$ but $\eta_i^j = 0$ for some *i*). The sparse group Lasso was first proposed by Simon et al. (2013). They provide an algorithm to solve this problem without deriving any statistical properties. I modify the estimator by taking the square-root of the mean square error term in the minimization problem. Similar to the square-root LASSO proposed in Belloni et al. (2011), the method becomes pivotal since it does not require a pre-estimation of the standard deviation σ . I will prove the statistical properties of square-root sparse group LASSO in section 5.

The de-sparse LASSO estimator for square-root sparse group LASSO is proposed as follows:

De-sparse 2SLSS Estimator for Square-root Sparse Group LASSO:

• Define

$$\hat{e}_m = \hat{\eta} + \hat{\Theta}_Z \hat{Z}'_n (D_n - \hat{Z}_n \hat{\eta} - X_n \hat{\beta})/n$$
$$\hat{b}_m = \hat{\beta} - (X'_n X_n)^{-1} X'_n \hat{Z}_n \hat{\Theta}_Z X'_n (D_n - \hat{Z}_n \hat{\eta} - X_n \hat{\beta})/n$$

where $\hat{Z}_n = \left[(M_n^1 \circ \hat{D}_n), (M_n^2 \circ \hat{D}_n), \cdots, (M_n^q \circ \hat{D}_n) \right]$ and $\hat{\Theta}_Z$ is the approximation of the inverse of the matrix $\frac{1}{n} \hat{Z}'_n W_n \hat{Z}_n$.

5 Statistical Properties

In this section, I consider the statistical properties for the de-sparse 2SLSS estimators $(\hat{e}, \hat{b}, \hat{S}_n)$ proposed in section 4. I show consistency and derive asymptotic normality for my de-sparse estimators. In order to show consistency and asymptotic normality for the de-sparse 2SLSS estimator with multiple networks, I derive the statistical properties for square-root sparse group LASSO, which have not been previously defined in statistics literature.

5.1 Consistency

The proof of consistency has two parts. 1) I show that the selected active set converges to the true non-zero parameter set. 2) I show that the de-sparse estimators converge to the true parameters.

Theorem 1. In heterogeneous endogenous effects model and with assumption 1-5, if $\lambda \propto \sqrt{\frac{\log n}{n}}$

• $\lim_{n\to\infty} \mathbb{P}(\hat{S}_n = S) = 1$

- $\hat{e} \rightarrow \eta_0$
- $\hat{b} \rightarrow \beta_0$

The consistency of the LASSO active set \hat{S}_n follows from assumption 4 as is shown in Zhao and Yu (2006). The consistency of \hat{e} and \hat{b} can be shown by taking the Karush-Kuhn-Tucker conditions of the LASSO minimization problem in the second stage. The shrinkage bias carried from the first stage: $\frac{1}{n}(M_n \circ \hat{D}_n)' \left(M_n \circ (\hat{D}_n - D_n)\right) \eta_0$ can be shown of order $o(1/\sqrt{n})$. The details of this proof are provided in the appendix.

In the presence of cliques, if γ is penalized, it can be treated as one of the components in η . On the other hand, if it is not penalized, it can be treated as one of the components in β . The consistency follows directly from Theorem 1:

Corollary 1. In the heterogeneous endogenous effects model with cliques and under assumptions 1'-3', assumptions 4-5, if $\lambda \propto \sqrt{\frac{\log n}{n}}$

- $\lim_{n\to\infty} \mathbb{P}(\hat{S}_n = S) = 1$
- $\hat{e} \rightarrow \eta_0$
- $\hat{r} \rightarrow \gamma_0$
- $\hat{b} \rightarrow \beta_0$

In the presence of multiple networks, theorem 2 summarizes the consistency results.

Theorem 2. In the heterogeneous endogenous effects model with multiple networks and under assumptions 1^*-5^* , if $\lambda_1 \propto \sqrt{\frac{\log n}{n}}$ and $\lambda_2 \propto \sqrt{\frac{\log n}{n}}$

- $\lim_{n\to\infty} \mathbb{P}(\hat{S}_n = S) = 1$
- $\hat{e}_m^j \to \eta_0^j \text{ for } j = 1, \cdots, q$
- $\hat{b}_m \to \beta_0$

The derivation of theorem 2 is similar to that of theorem 1 expect that the square-root LASSO is replaced with the square-root sparse group LASSO. Theorem 1, corollary 1 and theorem 2 establish the consistency for my de-sparse 2SLSS estimators.

5.2 Asymptotics

Post inference or inference after model selection are not uniformly valid. Define the set:

$$B(s) = \{\eta \in \mathbb{R}^n | \{j, \eta_j \neq 0\} \le s\}$$

As shown in Leeb and Potscher (2005), Leeb and Potscher (2008), and Kasy (2015)

$$\sup_{\eta_0 \in B(s)} \left| P\left(\frac{\sqrt{n}(\hat{\eta}_j - \eta_0)}{\hat{V}_j} < t\right) - \Phi(t) \right| \nrightarrow 0$$
(13)

where $\hat{\eta}_j$ can be any estimator based on a selected model, \hat{V}_j is the associated standard deviation and $\Phi(t)$ is the normal CDF function. When η_j is of order $O(1/\sqrt{n})$, the probability that LASSO fails to select this regressor into the active set can be non-zero. The resulting post model selection estimator will carry the omitted variable bias because of the exclusion of regressor j from the model. Thus, post inference conditioning on the selected model cannot converge to the true parameters uniformly over the models defined by sparsity.

On the other hand, the de-sparse LASSO estimator is uniformly valid since the inference is not conditioned on the selected model (see van de Geer et al., 2014). I follow the same idea and show that my de-sparse 2SLSS estimators achieve asymptotic normality with square-root LASSO and square-root sparse group LASSO.

Theorem 3. In the heterogeneous endogenous effects model and under assumption 1-5, if $\lambda \propto \sqrt{\log n/n}$

$$\sqrt{n}(\hat{e} - \eta_0) = E_1 + \Delta_1,$$
$$\sqrt{n}(\hat{b} - \beta_0) = E_2 + \Delta_2,$$

where

$$E_1 \sim N(0, \sigma^2 \Theta_1 diag(\Gamma) \Omega diag(\Gamma) \Theta'_1),$$

$$E_2 \sim N(0, \sigma^2 \Theta_2 diag(\Gamma) \Omega diag(\Gamma) \Theta'_2),$$

and

$$\begin{split} \|\Delta_1\|_{\infty} &= o_p(1), \quad \|\Delta_2\|_{\infty} = o_p(1), \\ \Gamma &= \lim_{n \to \infty} (I - M_n \circ \eta_0)^- X_n \beta_0, \\ \Theta_1 &= \lim_{n \to \infty} \hat{\Theta}, \quad Z_n = (M_n \circ \hat{D}_n), \quad \tilde{Z}_n = X_n (X'_n X_n)^{-1} X'_n Z, \\ \Theta_2 &= \lim_{n \to \infty} \frac{1}{n} \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big) \end{split}$$

Theorem 3 shows that the 2SLSS estimator achieves normality at the standard rate \sqrt{n} . The shifts

 Δ_1 and Δ_2 represent the bias from using nodewise regression and they are shown to be $o_p(1)$ with the proper choice of tuning parameters.

Corollary 2. In the heterogeneous endogenous effects model with cliques and under assumption 1'-3', and assumptions 4-5, if $\lambda \propto \sqrt{\log n/n}$

$$\sqrt{n} \begin{pmatrix} (\hat{e} - \eta_0) \\ (\hat{r} - \gamma_0) \end{pmatrix} = E_1 + \Delta_1,$$
$$\sqrt{n} (\hat{b} - \beta_0) = E_2 + \Delta_2,$$

where

$$E_1 \sim N(0, \sigma^2 \Theta_1 diag(\Gamma) \Omega diag(\Gamma) \Theta'_1),$$

$$E_2 \sim N(0, \sigma^2 \Theta_2 diag(\Gamma) \Omega diag(\Gamma) \Theta'_2),$$

and

$$\begin{split} \|\Delta_{1}\|_{\infty} &= o_{p}(1), \quad \|\Delta_{2}\|_{\infty} = o_{p}(1), \\ \Gamma &= \lim_{n \to \infty} (I - M_{n} \circ \eta_{0})^{-} X_{n} \beta_{0}, \\ \Theta_{1} &= \lim_{n \to \infty} \hat{\Theta}, \quad Z_{n} = [(M_{n} \circ \hat{D}_{n}), M_{n} \hat{D}_{n}], \quad \tilde{Z}_{n} = X_{n} (X_{n}' X_{n})^{-1} X_{n}' Z, \\ \Theta_{2} &= \lim_{n \to \infty} \frac{1}{n} \Big(I - Z_{n} \hat{\Theta} \tilde{Z}_{n}' / n \Big)' X_{n} (X_{n}' X_{n})^{-1} X_{n}' \Big(I - Z_{n} \hat{\Theta} \tilde{Z}_{n}' / n \Big) \end{split}$$

For my setting with multiple networks, I derive the following results:

Theorem 4. In the heterogeneous endogenous effects model with multiple networks and under assumptions 1^*-5^* , if $\lambda_1 \propto \sqrt{\frac{\log n}{n}}$ and $\lambda_2 \propto \sqrt{\frac{\log n}{n}}$

$$\sqrt{n}(\hat{e}_m - \eta_0) = E_{m1} + \Delta_{m1},$$
$$\sqrt{n}(\hat{b}_m - \beta_0) = E_{m2} + \Delta_{m2},$$

where

$$E_{m1} \sim N(0, \sigma^2 \Theta_{Z1} diag(\Gamma) \Omega_m diag(\Gamma) \Theta'_{Z2}),$$

$$E_{m2} \sim N(0, \sigma^2 \Theta_{Z2} diag(\Gamma) \Omega_m diag(\Gamma) \Theta'_{Z2}),$$

and

$$\begin{split} \|\Delta_{m1}\|_{\infty} &= o_p(1), \quad \|\Delta_{m2}\|_{\infty} = o_p(1), \\ \Theta_{Z1} &= \lim_{n \to \infty} \hat{\Theta}_Z, \quad Z_n = (M_n \circ \hat{D}_n), \quad \tilde{Z}_n = X_n (X'_n X_n)^{-1} X'_n Z, \\ \Theta_{Z2} &= \lim_{n \to \infty} \frac{1}{n} \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - Z_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big)$$

The proof of Theorem 2 and Theorem 4 requires the following results from the square-root sparse group LASSO: 1) Bounds on the prediction, i.e. $\left\|\sum_{j=1}^{q} (M^{j} \circ X_{n})(\hat{\eta}^{j} - \eta_{0}^{j}) + X_{n}(\hat{\beta} - \beta_{0})\right\|_{2} \lesssim \lambda$. and 2) Consistency of selection i.e. $\hat{S}_{n} = S$. I prove these two statistical properties in the appendix.

6 Simulations

In this section, I report Monte Carlo simulation results for my heterogeneous endogenous effects model and its extension with cliques and with multiple networks. My results are robust when applied to networks generated by different algorithms and to networks of different sizes.

6.1 Heterogeneous Endogenous Effects Model

To assess the finite sample performance of my estimator for the heterogeneous endogenous effects model, I use the Erdos-Renyi algorithm to simulate a network of size n. Individuals are added into the graph one at a time. When one individual is added to the network, she has probability p of generating a link with all existing individuals independently. I choose p = 0.1 and p = 0.2 in the simulation. I avoid a large p because collinearity among regressors may arise when links become very dense, violating assumption 5.

I set the first 5 individuals to be influential by letting their coefficients η_j be non-zero. To guarantee the existence of endogenous effects, I arbitrarily specify the connections among these five individuals. The adjacency matrix M_n for the five influential individuals is given in the appendix. If the connections among these five individuals are not fixed, there is a possibility that no connections are formed among these five and thus there is no endogeneity in the network. In this case, the results will be too good in such a case.

The true parameters are fixed as $\beta_0 = 3$, $\eta_{0,1} = \eta_{0,2} = \eta_{0,3} = \eta_{0,4} = \eta_{0,5} = 0.5$, and $\eta_{0,j} = 0$ for j > 5. Individual characteristics X_n are generated from a standard normal distribution.

Individual outcomes Y_n are then generated as $Y_n = (I - M_n \circ \eta_0)^{-1} (X_n \beta_0 + \epsilon_n)$ where ϵ_n is drawn independently from a standard normal distribution.

I use (M_n, X_n, Y_n) as observations and apply my two-stage LASSO estimator. I construct the desparse 2SLSS estimator and repeat the above process 200 times in a manner similar to van de Geer et al. (2014).

I report the average coverage probability (Avgcov) and average length (Avglength) of confidence intervals for the coefficients for influential individuals, $\{\eta_1, \dots, \eta_5\}$, the coefficient for individual characteristics, β_0 , and the coefficients for non-influential individuals, the $\eta_j s$ (j > 5). For example:

Avgcov
$$S_0 = s_0^{-1} \sum_{j \in S_0} \mathbb{P}[\eta_{0,j} \in CI_j]$$
 (14)

Avglength
$$S_0 = s_0^{-1} \sum_{j \in S_0} length(CI_j)$$
 (15)

I separately report the average coverage and average length for each of the five influential individuals. As shown in table 8, the coverage is around the nominal 95% level and the length of the confidence intervals decreases as the sample size grows.

Since we can construct confidence intervals for all n coefficients, joint inference can be performed under the control of False Discover Rate (FDR). As shown in equation (16), the power reported in table 8 represents the average percentage in the active set (i.e. $\{1, 2, 3, 4, 5\}$) that is significant after controlling for the False Discover Rate (FDR) at 5% using the Benjamini-Hochberg method. The FDR reported in table 8 represents the average percentage of the non-active set (i.e. $\{6, 7, \dots, n\}$) that is significant after controlling the FDR at 5% using the Benjamini-Hochberg method. The exact definition is as in equation (17).

Power
$$= s_0^{-1} \sum_{j \in S_0} \mathbb{P}[H_{0,j} \text{ is rejected}]$$
 (16)

$$FDR = \sum_{j \in S_0^c} \mathbb{P}[H_{0,j} \text{ is rejected}] / \sum_{j=1}^n \mathbb{P}[H_{0,j} \text{ is rejected}]$$
(17)

The power varies because the networks change when the sample size increases. It is strictly increasing when the network is sparse (i.e. p = 0.1). The power decreases in the p = 0.2 case as the problem of endogeneity increases when the network is dense. The empirical FDR is controlled well, which all under the 5% rate. Notice that the confidence interval's length is large when the sample size equals 50. This is because when the number of individuals is small, some individuals might only connect to 1 or 2 other individuals. This means that the regressors that represent this individual are all 0s except for a small numbers of non-zero terms, which leads to a large standard error.

The two-stage LASSO estimator requires the choice of two tuning parameters (i.e. the two λ s from both stages as in section 4.1). Moreover, when calculating $\hat{\Theta}$ in the De-sparse 2SLSS estimator (section 4.2) and using the nodewise regression, one also need to choose a tuning parameter. Following the suggestion in Belloni et al. (2011), I use a benchmark choice of λ for the first stage and nodewise regression (i.e. $\lambda \propto \Phi^{-1}(1 - \alpha/(2n))/\sqrt{n}$), where $\Phi^{-1}(.)$ is the inverse of normal cdf function. For the second stage, I use cross-validation to pick λ to enhance finite sample performance.

I further increase the number of influential individuals to 10 and report the results in table 9. Again, to guarantee the existence of endogeneity, the adjacency matrix for these ten individuals is set as shown in the appendix. All average coverages and average confidence interval lengths are separately reported for these ten individuals.

The choice of the tuning parameters is similar to those used to generate table 8 for networks with 50 and 200 individuals. For networks with 500 individuals, I use benchmark λ to replace cross validation in the second stage. The idea is to show the converge of the process, such that valid coverage can still be generated under theory guide tuning parameters (see Belloni et al., 2011).

As shown in table 9, all coverages are very close to the nominal levels. The average lengths of confidence intervals is slightly larger compared with table 8. This is due to the increase in influential individuals; it is more difficult to differentiate them from those irrelevant individuals.

Table 10 presents the result when a network is generated using the Watts-Strogatz mechanism or the "small world" network. Define the pN (even number) as the mean degree for each node and a special parameter $\omega = 0.4$. The WattsStrogatz mechanism works as follows:

- construct a graph with N nodes each connected to pN neighbors, which $\frac{pN}{2}$ on each side.
- For each node n_i , take every edge (n_i, n_j) with i < j and rewrite it with probability ω . Rewrite means replace (n_i, n_j) with (n_i, n_k) where k is choosing uniformly among all nodes that are not currently connected with n_i

The influential individuals are chosen as the 1st, 5th, 15th, 40th and 50th individuals in the network. As shown in table 10, my estimator is robust under a "small world" algorithm. Nominal level is reached as the size of the network grows and the length of confidence intervals is slightly smaller than in the standard case.

6.2 Heterogeneous Endogenous Effects Model with Cliques

Table 12 presents results for the heterogeneous endogenous effects model with cliques. The outcome variable Y_n is now generated as $Y_n = (I - M_n \circ \eta_0 - M_n \gamma_0)^{-1} (X_n \beta_0 + \epsilon_n)$. The coefficient of the homogeneity effects γ_0 is set at 0.05.

The choice of the tuning parameters is similar to that used to generate table 8 for networks with 50 and 200 individuals. For networks with 500 individuals, I use benchmark λ (i.e. $\lambda \propto \Phi^{-1}(1 - \alpha/(2n))/\sqrt{n}$) to replace cross validation in the second stage.

The coverage is above the 95% nominal level in all cases. I also report the mean coverage and average length of the confidence interval for the coefficient of the homogeneous effects. My model gives above 95% coverage in all cases. I also report the empirical probability of rejecting a null hypothesis of zeros effects at 95% nominal level. The probability of rejecting the test converges to 1 when the sample size grows to 500.

6.3 Heterogeneous Endogenous Effects Model with Multiple Networks

In this Monte Carlo exercise, I include two different networks generated by the Erdos-Renyi algorithm, where one is influential and the other is not. I use the two-stage LASSO estimator with multiple networks to estimate the parameters. The square-root sparse group LASSO requires two tuning parameters, one for the l_2 norm and the other for the l_1 norm. I set the two parameters to be equal to each other as the correlations among the columns of the adjacency matrices are very small. The choice of tuning parameters is similar to that used to generate table 1 for networks with 50 and 200 individuals. For networks with 500 individuals, I use a rule of thumb to choose λ instead of cross-validation in the second stage. Table 11 summarizes the results. As in previous results, all coverages exceed the nominal 95% level.

I report the empirical probabilities such that at least one individual is detected in a given network controlling for the FDR at 5% using the Benjamini-Hochberg method. I also report the average number of detections conditioning on at least one individual who is detected in a given network. Tables 11 shows that network 1, which is the relevant network, is more likely to be detected in all cases than network 2, the irrelevant network. The average number of identified individuals for network 1 is also more than that of network 2.

7 Empirical Application

I use the proposed estimator to study the importance of different networks in spreading the participation in a micro finance program within rural Indian villages. I show that different kinds of networks have different effects on individuals decisions. I identify the influential individuals in each village. My analysis shows that leaders among agricultural laborers, Anganavadi teachers, construction workers, small business owners and mechanics are very likely to be influential in the villages.

7.1 Background

A non-profit organization named Bharatha Swamukti Samsthe (BSS) has been running micro finance programs in rural southern Karnataka, India since 2007. It provides small loan products to poor women and, through them, to their families. The villages covered by the program are geographically isolated and heterogeneous in terms of caste.

When BSS initially introduces a micro finance program to a village, the credit officers of BSS first approached a number of "predefined leaders", such as teachers, shopkeepers and village elders. BSS held a private meeting with these leaders and explained the program. Then these predefined leaders passed the information onto other villagers. Those who were interested in the program and contacted BSS were trained and assigned to groups to receive credit. Each group consisted of 5 borrowers and group members were jointly liable for loans. Loans were around 10,000 rupees (approximately \$200) at an annualized rate of approximately 28%. Note that 74.5 percent of the households in rural area said the monthly income of their highest earning member is less than 5,000 rupees (source: Socio-Economic Caste Census-2011). This loan had to be repaid within 50 weeks.

In 2006, 75 villages in Karnataka were surveyed 6 months before the initiation of the BSS micro finance program. This survey consisted of a village questionnaire and a detailed follow-up survey conducted among a subsample of villagers. The village questionnaire gathered demographic information on all households in a village including GPS coordinates, age, gender, number of rooms, whether the house had electricity, and whether the house had a latrine. The data set also contains information on the "pre-defined leaders" set who helped spread the information to the entire village. The follow-up survey collected data from a villager sample stratified according to age, education level, caste, occupancy, etc. It also asked questions about social network structures along 12 dimensions, including:

- Friends: Name the 4 non-relatives whom you speak to the most.
- Visit-go: In your free time, whose house do you visit?
- Visit-come: Who visits your house in his or her free time?
- Borrow-kerorice: If you needed to borrow kerosene or rice, to whom would you go?

- Lend-kerorice: Who would come to you if he/she needed to borrow kerosene or rice?
- Borrow-money: If you suddenly needed to borrow Rs. 50 for a day, whom would you ask?
- Lend-money: Who do you trust enough that if he/she needed to borrow Rs. 50 for a day you would lend it to him/her?
- Advice-come: Who comes to you for advice?
- Advice-go: If you had to make a difficult personal decision, whom would you ask for advice?
- Medical-help: If you had a medical emergency and were alone at home whom would you ask for help in getting to a hospital?
- Relatives: Name any close relatives, aside from those in this household, who also live in this village.
- Temple-company: Do you visit a temple/mosque/church? Do you go with anyone else? What are the names of these people?

For the 43 villages where micro finance was introduced by the time of 2011, BSS also collects information on which villagers have joined the program. These survey questions reveal the underlying structures for connections among any two individuals in the network. Figure 4 presents all those connections at the household-level in a graph. Each node in the graph represents a household. A green node indicates that the household joined the micro finance program, while a blue node indicates that it did not. Bigger nodes represent those households in which at least one family member has been chosen as being among the "pre-defined leaders". An edge between two nodes signifies that the two nodes are connected in at least one of the 12 networks. The darker the color of the edge, the more connections it represents.

This dataset provides an ideal framework for application of the heterogeneous endogenous effects model. First, it allows me to model endogenous effects. An individual may decide to join the micro finance program if her neighbors or friends plan to join. Second, the endogenous effects are individual specific. Given the diversity of the villagers, it is possible that some villagers are more influential than others. Third, it allows me to implement the heterogeneous endogenous effects model with multiple networks. The questions asked regarding multiple dimensions of the network structure allow me to explore which network is most influential.



Figure 4: Network in Village 1
7.2 Data

In this empirical study, I focus on the 38 villages that have been introduced to the micro finance programs by BSS and have data publicly available ³. For each village, I can observe both its social network structure and the villagers' decisions about joining the program. I drop the data for one village (Village 46) that contains incorrect entries on the index of households. Table 13 summarizes the descriptive statistics for each village.

Among the 12 questions about the social network structure, 4 pairs essentially capture the same connections among the villagers ⁴. Therefore, I consolidate each pair of questions into one dimension:

- Visit-go-come
 - In your free time, whose house do you visit?
 - Who visits your house in his or her free time?
- Borrow-Lend-kerorice
 - If you needed to borrow kerosene or rice, to whom would you go?
 - Who would come to you if he/she needed to borrow kerosene or rice?
- Borrow-Lend-money
 - If you suddenly needed to borrow Rs. 50 for a day, whom would you ask?
 - Who do you trust enough that if he/she needed to borrow Rs. 50 for a day you would lend it to him/her?
- Help decision
 - Who comes to you for advice?
 - If you had to make a difficult personal decision, whom would you ask for advice?

I restructure all the data at the household level as only women are allowed to apply for the micro finance program because the goal of BSS is to support families through the women in them. As a

³The dataset can be downloaded from http://web.stanford.edu/ jacksonm/Data.html

⁴Assuming every villager truthfully answers a pair of questions, the adjacency matrices associated with each question are the same. It is also plausible to treat villagers' answers to each question as a separate directed graph. However, these questions do not allow for clear determination of the directions. For example, if villager A visits villager B's house, it is not clear whether villager A influences villager B or vice versa

result, a womans decision to join or not join the micro finance program becomes her familys decision. A connection between two villagers becomes a connection between two families. A "predefined leader" is a villager selected by BSS to help spread information about the micro finance program to the other villagers. At the household level, I use the term "predefined leader" for a household that contains at least one such villager.

7.3 Sparsity and Equilibrium

To demonstrate how my method identifies influential households, I model families' decisions regarding joining the micro finance program as a network game with Bayesian Nash Equilibrium. For household *i*, let d_i^* be the expected probability that *i* chooses to join the micro finance program. The decision of household *i* depends on its neighbors' decisions as well as the types of connections between them. The decision also depends on its characteristics X_i and on unobserved information ϵ_i . Formally, it can be written as:

$$d_i^* = \sum_{l \in N_i} d_l^* (\sum_{j=1}^q \eta_l^j) + x_i \beta + \epsilon_i$$

Rewritten in matrix form:

$$D_n^* = \sum_{j=1}^q \left(M_n^j \circ D_n^* \right) \eta^j + X_n \beta + \epsilon_n$$

By Assumption 1-3, there is a unique equilibrium that determines D_n .

I assume that only a small number of households are influential over their neighbors. Leaders and followers are usually observed in those rural villages. Big decisions are often made by the village elders or by the more educated among the villagers. BSS recognized the importance of leaders and gathered a group of predefined leaders, asking them to inform the rest of the villagers about their program. I do not consider the local level influence in these villages given the size and how complicated the network structures are. Households are closely connected by these 8 networks as shown in Figure 4 and there is no form of clique visible.

Because the villages are considered geographically isolated, I apply my estimator separately to each of the 38 villages. I use the number of rooms per person in a houshold as the independent variable X_n . Number of rooms per person is a proxy for the wealth in the family. As shown in table 1, it is negatively correlated with the decision to join the micro finance program. The richer the family, the less likely the family is to participate in the micro-finance program. I further check the robustness of my independent variable by including additional controls. The adjacency matrix M_n^j is constructed from the questions in the survey. Households i and k are connected in network j if either i or k reported the other in question j. Finally, d_i^* is replaced with the household's choice.

	(1)	(2)	(3)
	Decision	Decision	Decision
Ave # room	1798***	1602***	0922***
	(.0180)	(.0201)	(.0215)
number of person		$.0084^{*}$	$.0085^{*}$
		(.0038)	(.0040)
Electricity			.0022
			(.0303)
latrine			0908***
			(.0183)
Ave #workers			.0242
			(.0347)
Ave age			0052***
			(.0007)
Constant	5949***	6457***	5018***
	(.0132)	(.0266)	(.0443)
N	8,375	8,375	8,375
R^2	0.0118	0.0124	0.0212

 Table 1: Independent Variable

Standard errors in parentheses * p < 0.05, ** p < 0.01, *** p < 0.001

Dependent variable is households' decision on whether to join the micro finance program or not. All design control village fixed effects.

The instruments are constructed as $(M_n^j \circ X_n)$ for $j = 1, 2, \dots, 8$. I use the heterogeneous endogenous effects model with multiple networks to: 1) Identify the effective networks affecting a household's decision and 2) Identify that households that are leaders in the village and study the association between observable characteristics and leader status. If a new program is going to try to recruit these households, the organizers can target those influential households and try to persuade them to join first.

7.4 Results

7.4.1 Identifying Effective Networks

First, I study how LASSO selects networks. I define a coefficient for a household's endogenous effect in a network as significant according to two different criteria. The first criterion, "Cross-Validation", determines a coefficient to be significant if LASSO predicts the coefficient to be non-zero after cross-validation. The second criterion, "De-sparse", first constructs a bias-adjusted coefficient and calculates its standard error. It then determines a coefficient to be significant if the Benjamini-Hochberg method rejects the null hypothesis of zero effect at the 5% false discovery rate. A network is defined as significant if at least one coefficient for heterogeneous endogenous effects in this network is significant.

Table 2 presents the empirical probability of the 8 networks being significant among the 37 villages. Note that certain types of networks (such as visit go-come) are more likely to pass influence then others (such as temple company). For example, by cross-validation criterion, the visit go-come network is detected as significant in 19 out of the 37 villages (i.e. 51%) whereas temple company is detected as significant in only 5 out of the 37 villages (i.e. 14%). I also present the average number of households associated with significant endogenous effects in each significant network. For example, according to the cross-validation criterion, 342 households in 19 villages have significant coefficients associated with the visit go-come network, which averages to 18 households per detection. On the other hand, 32 households in 5 villages have significant coefficients associated with the temple company network, which averages to 6 households per detection.

In terms of variable selection, if Assumption 4 holds, the cross-validation criterion may consistently select the truly influential households with high probability even in a finite sample. On the other hand, the de-sparse criterion is likely to be conservative because of its use of the false discovery control process. In terms of coefficients estimated, de-sparse estimators are asymptotically consistent. On the other hand, estimates based on the LASSO estimator suffer from shrinkage bias and are not consistent.

Table 3 reports the average absolute heterogeneous endogenous effects within significant networks using the de-sparse estimators. For example, if all else is equal, an additional influential neighbor in the visit go-come network will, on average, increase the probability of joining the micro-finance program by 16%; moreover, an additional influential neighbor in both the visit go-come network and the friendship network will increase the probability of joining the micro-finance program by 16% on average. The magnitude of those coefficients should not be over interpreted as exogenous effects and correlated effects are not considered in the model. Similar to Table 2, certain types of networks (such as visit go-come) pass stronger influence than others (such as temple company). Note that, in most of the cases, networks that are more likely to pass influence also pass stronger influence. The relative network is an exception. Even though the relative network is less likely to pass influence compared to the friendship network, the borrow-lend-money network and the help decision network, it passes stronger influence once it is significant. Table 3 also presents the percentage of positive effects detected among different networks. For networks such as visit-go-come and friendship, more than 70% of influential villagers are "true leaders" – if they decide to join the micro-finance program, their neighbors will follow them and join the program. On the contrary, for the temple company network, it is almost equally likely for neighbors of influential households to either follow the same decision or choose the opposite.

[insert table 2]

Table	2:	Second	Stage:	network	usage
Table	4.	Decona	bluge.	nouvoin	abage

		visit	borrow-lend	borrow-lend	friendship	medical	help	relatives	temple
		go-come	keroric	money		help	decision		company
Cross 1	probability 3	51%	43%	41%	41%	30%	32%	30%	14%
Validation	identified $^{\rm 4}$	18	12	13	14	9	14	9	6
De sparse ²	probability 3	51%	46%	51%	43%	32%	41%	43%	19%
De-sparse	identified 4	3	3	3	2	3	3	3	2

0. Reported are the probability of detection among the 38 villages.

1. Cross Validation represents those networks identified from lasso using cross validation.

2. De-sparse represents those networks identified from De-sparse criterion using FDR control.

3. Probability reports the empirical probability that at least one regressor in the group is significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

4. Identified reports the averaged number of significant regressors in the group conditioning on the network being significant. False discover rate is controlled at 5% using Benjamini-Hochberg method.

[insert table 3]

The results in Table 2 and Table 3 suggest villagers are more likely to discuss the micro-finance program when they visit each other, chat with friends, and meet with people to whom they are economically connected. Villagers are not likely to talk about the micro finance program when they go to the temple.

	visit	borrow-lend	borrow-lend	friendship	medical	help	relatives	temple
	go-come	keroric	money		help	decision		company
absolute magnitudes $^{\rm 1}$	0.1543	0.1443	0.1245	0.1194	0.1214	0.1217	0.1404	0.0555
percentage of positive effect								
	77%	67%	69%	70%	68%	77%	67%	55%

Table 3: Second Stage: average endogenous effect

1. The magnitudes are reported based on the de-sparse estimator.

Reported numbers are conditioning being significant using De-sparse method.

To verify my findings above, I provide exogenous evidence using centrality measures. Intuitively, the more a villager is exposed to a network, the more likely she is to be connected to influential villagers, and hence she is more likely to join the program. Following Banerjee et al. (2013), I measure the centrality of each villager in each network through "degree", "closeness", "betweenness" and "eigenvector". (See Appendix for definitions) Then I regress households' decisions on whether to join the micro finance program on each centrality measure separately while controlling for village fixed effects:

$$d_j = C_j^q \beta + \gamma_j + \epsilon_j \tag{18}$$

where d_j is household j's decision; C_j^q is household j's centrality in the network q; γ_j is the village fixed effect; and ϵ_j is the error term.

Table 7 presents the regression results for equation (18). Visit go-come and borrow-lend kerorice are positively correlated with degree, closeness and eigenvector centrality. Friendship, borrowlend money and medical help are positively correlated with degree and closeness centrality. This is consistent with my findings that these four networks are more effective in passing influence. Meanwhile, neither help decision, relative nor temple company are found to be correlated with any of the centrality measures defined above. This is also consistent with the lower probability of passing influence as found in table 2. Note that none of the networks are found to correlate with betweenness centrality. This is because betweenness centrality is based on the shortest paths in a network, which is not a direct measure of the exposure of an individual to a network.

[insert table 6]

	visit go-come	borrow-lend kerorice	borrow-lend money	friendship	medical help	help decision	relatives	temple company
degree	0.0025^{**}	0.0032^{**}	0.0020^{**}	0.0022^{**}	0.0032^{**}	0.0013	0.0035	0.0061
	(0.0009)	(0.0011)	(0.0010)	(0.0010)	(0.0014)	(0.0011)	(0.0019)	(0.0032)
closeness	$\begin{array}{c} 32.9116^{***} \\ (9.5639) \end{array}$	40.2695^{***} (10.7603)	$29.7981^{**} \\ (9.3901)$	31.0882^{***} (9.0446)	32.5602^{**} (11.1383)	18.1944 (9.8242)	7.9509 (16.5557)	$231.0770 \\ (134.3147)$
betweenness	1.3565	0.1751	0.2940	1.6713	1.1662	-0.5728	0.3055	-0.2093
	(1.0101)	(0.8240)	(0.9504)	(1.0207)	(0.8634)	(0.8283)	(0.7736)	(0.2178)
eigenvector	3.6201^{***}	1.5239^{**}	0.1161	1.3927	-0.7338	-0.2387	0.7753	3.3015
	(0.8888)	(0.6250)	(0.8245)	(0.8271)	(0.7709)	(0.7603)	(0.5672)	(3.5585)

 Table 4: Centrality Measure

Standard errors in parentheses * p < 0.05, ** p < 0.01, *** p < 0.001

Definition for centrality measures are in appendix.

7.4.2 Identifying Influential Households

Second, I focus on how LASSO selects households. I compare the LASSO selected influential households with the BSS selected "predefined leaders". It is important to point out that these "predefined leaders" are *not* necessarily influential villagers in a network. Recall that predefined leaders are a set of villagers that BSS select to help spread the information about the micro finance program. The fact that a villager is selected as a "predefined leader" to *pass information* about the micro finance program does not a priori guarantee her or her family's *influence* – her decision to join the micro finance program may not lead to her neighbors' decisions to join. In the analyses below, I will examine how influential villagers are associated with "predefined leaders" and explore their potential differences.

1. Influential Predefined Households

In table 7, I report results indicating that influential households selected by LASSO partly overlap with "predefined leaders". This is intuitive because some "predefined leaders" such as school headmasters and village elders are highly respected figures in a village. Therefore, their decisions are likely to be followed by others in the village. On average, BSS selected 27 villagers as "predefined leaders" in each village. In comparison, Cross-Validation criterion selects around 22 villagers and De-sparse criterion selects around 6. Furthermore, on average, 4 out of 22 influential villagers (i.e. 19%) selected by Cross-Validation criterion are also BSS "predefined leaders"; 1 out of 6 influential villagers (i.e. 13%) selected by De-sparse criterion are also BSS "predefined leaders". In Table 14 below, I show that small business owners are more likely to be both influential and selected as "predefined leaders".

[insert table 7]

Table 5: Second Stage: coverage of predefined leaders

	coverage 2	total number of discovery ³
Cross Validation 4	19%	22
De-sparse 5	13%	6

1. predefined leaders are a set of villagers defined by BSS, who helped spread the information about the micro-finance program.

2. Coverage reports the percentage of individuals detected by LASSO and also selected as "predefined leaders" in total detection.

3. Total number of discovery reports the total number of individuals discovered by lasso using each method.

4. Cross Validation represents those individuals identified from lasso using cross validation .

5. De-sparse represents those individuals identified from De-sparse criterion controlling FDR.

6. The average number of predefined leaders in one village is 27

2. Influential Non-Predefined Households

In this and the following section, I focus on understanding the differences between the influential households selected by LASSO and the "predefined households" selected by BSS. I investigate the likelihood that a household being selected by LASSO or by BSS, as associated with the careers of its family members. More specifically, I regress whether a household is selected as "predefined leader" (Design (1)), whether a household is selected by LASSO as influential (Design (2)), and whether a household joins the micro finance program (Design (3)), separately on dummy variables based on the full set of careers as reported in the survey data controlling for other household characteristics and village fixed effects. The full results of these regressions are reported in Table 14 in the Appendix.

Table 6 summarizes all careers that have a significant impact on the likelihood of a household being selected by LASSO as influential. Note that except for small business owners, all the other careers in this table are not significantly associated with the likelihood of a household being selected by BSS as being among the "predefined leaders". Over 67% of the villagers are agricultural laborers and 75% of the LASSO selected influential households have agricultural laborers in the family. Anganwadi Teacher is a set of groups that provides pre-school education to the children. They are part of the government's health care system in the rural areas. There are 31 Anganwadi Teachers in all villages, and LASSO detects 7 of their families to be influential. BSS also selects 7 of them as "predefined leaders" but only 2 of the 7 are selected by LASSO as influential. Other careers that are

correlated with LASSO selection include police officer, mechanic, and skilled laborers. These are more educated individuals and it seems compelling that they are selected as influential individuals.

Table 7 summarizes all careers that have a significant impact on the likelihood of a household being selected by BSS as being among the "predefined leaders". Poojari are Indian priests in those villages and they are very likely to be included as "predefined leaders". However, they are not likely to influence people to join the micro finance program. Other careers as tailor, hotel workers, veteran, and barber are included as "predefined leaders" because individuals doing these jobs can spread information quickly in the village. However, LASSO does not find these individuals to be influence.

[insert table 9]

8 Conclusions

In this paper, I propose a novel SAR model which allows for *heterogeneous* endogenous effects. Specifically, each individual has an individual-specific endogenous effect on her neighbors. My approach is useful for modeling a network with leaders and followers. For example, it can model how online opinion leaders influence the public or how experienced workers boost coworkers' productivity.

I propose a set of instruments as well as a two stage LASSO (2SLSS) method to estimate my model. The instruments are constructed as a function of the independent variables and an adjacency matrix. I use a LASSO type estimator to select the valid instruments in the first stage and the influential individuals in the second stage. I propose a bias correction for my two-stage estimator following van de Geer et al. (2014). I derive the asymptotic normality for my "de-sparse" two-stage LASSO estimator and conduct robust inference including confidence intervals.

My model can be extended to allow for more flexible structures. To apply LASSO, I assume that the number of influential individuals is sparse. I propose heterogeneous endogenous effects model with cliques to incorporate locally influential individuals, where the sparsity assumption is only applied to globally influential individuals. My model can also be extended to situations where there are multiple networks. I propose the use of the square-root sparse group LASSO in my 2SLSS process. I derive the convergence rate and prove the consistency of selection for the square-root sparse group LASSO estimator.

	(1)	(2)	(3)
Agriculture labour	-0.0141	0.0476^{*}	0.0672***
	(0.0136)	(0.0286)	(0.0134)
Anganwadi Teacher	0.0386	0.0664	0.1248^{**}
	(0.0602)	(0.1269)	(0.0593)
Blacksmith	-0.0752	-0.2279	0.1606^{*}
	(0.0927)	(0.1954)	(0.0913)
Construction/mud work	0.0050	0.2199^{***}	0.0562^{**}
	(0.0258)	(0.0544)	(0.0254)
Small business	0.2006^{***}	0.1287^{***}	0.0606^{***}
	(0.0227)	(0.0479)	(0.0224)
Police officer	-0.1459	-0.0374	0.3282^{*}
	(0.1917)	(0.4044)	(0.1890)
Mechanic	0.0106	-0.1237	0.1274^{**}
	(0.0634)	(0.1337)	(0.0625)
Skilled labour/work for company	0.0469	0.0252	0.0809^{*}
	(0.0491)	(0.1036)	(0.0484)
Control other careers	Y	Y	Y
Control village fix effect	Y	Υ	Υ

Table 6: Second Stage: who are they

Standard errors in parentheses * p < 0.1, ** p < 0.05, *** p < 0.01

design (1) uses whether one is predefined leaders as response variable

design (2) uses whether one joins the micro-finance program as response variable

design (3) uses whether one is selected by lasso as response variable

	(1)	(2)	(3)
Small business	0.2006***	0.1287***	0.0606***
	(0.0227)	(0.0479)	(0.0224)
Tailor Garment worker	0.0903^{***}	0.1169^{*}	0.0309
	(0.0304)	(0.0642)	(0.0300)
Hotel worker	0.3299^{***}	0.4257^{***}	0.0759
	(0.0750)	(0.1581)	(0.0739)
Poojari	0.3697^{***}	-0.1542	0.1501
	(0.1369)	(0.2887)	(0.1349)
Veterinary clinic	0.8649^{***}	1.9114^{***}	0.0377
	(0.3314)	(0.6990)	(0.3266)
Barber/saloon	0.4883^{***}	-0.0036	0.0443
	(0.1005)	(0.2119)	(0.0990)
Doctor/Health assistant	0.2691^{**}	0.2703	0.0874
	(0.1053)	(0.2222)	(0.1038)
Control other careers	Y	Y	Y
Control village fix effect	Υ	Y	Y

Table 7: Second Stage: who are they

Standard errors in parentheses * p < 0.1, ** p < 0.05, *** p < 0.01

design (1) uses whether one is predefined leaders as response variable

design (2) uses whether one joins the micro-finance program as response variable

design (3) uses whether one is selected by lasso as response variable

I apply my method to study villagers' decisions to participate in micro-finance programs in rural areas of Indian. I show that leaders in those villages have significant influence over their neighbors' decision to join the micro-finance program, and I provide rankings for the different social and economic networks among villagers. Based on how effectively each network spreads the impact of influential individuals' decisions, my method shows that some networks such as "visit go-come" and "borrow money" are much more effective in influencing villagers' decisions than other networks such as "temple company" and "medical help". I further show that individuals from certain careers such as agricultural workers, Anganwadi teachers and small business owners are more likely to influence other villagers.

There are two interesting directions for future research. First, it is possible to include heterogeneous exogenous effects in the model. These effects aim to capture how an individual's outcome varies with the exogenous characteristics of her neighbors. However, when both exogenous and endogenous effects are included in standard SARs, an identification problem known as the "reflection problem" may arise (see Manski, 1993). A similar problem arises also in my model if heterogeneous exogenous effects are included. Bramoullé et al. (2009) show that under additional assumptions on the adjacency matrix, this problem can be solved in SARs. With similar restrictions on the adjacency matrix, it is possible to construct a new set of instruments to include heterogeneous exogenous effects in my model.

Second, it might be possible to use penalized GMM type estimator to estimate my model. 2SLS and GMM are the two most commonly used estimators to deal with endogeneity in SARs. My 2SLSS can be rewritten as a penalized GMM problem. The current progress on penalized GMM estimators include Fan and Liao (2014) and Luo and Chernozhukov (2016). But no uniformly valid inference method currently exists for penalized GMM.

A Proofs

Lemma 1. Write $D_n = (I - M_n \circ \eta)^- X_n \beta + (I - M_n \circ \eta)^- \epsilon = f(M_n \circ X_n) + \epsilon_1$. In the first stage problem (19), under assumption 5 and square root lasso regularization parameter $\lambda \ge c\Lambda/n$,

$$\frac{\|\hat{D}_n - f\|_2}{\sqrt{n}} \le Cs^{1/2}\lambda$$

where,

$$\Lambda = n \left\| \nabla \sqrt{\hat{Q}(\beta^0)} \right\|_{\infty} = \max_{1 \le i \le p} \left\{ \frac{\sqrt{n}((X^i)'\epsilon)}{\|\epsilon\|_2} \right\}$$

Lemma one is the same as Theorem 1 in Belloni et al. (2011)

A.1 Theorem 1

Proof. In the second stage, \hat{D}_n is used to replace D_n

$$(\hat{\beta}, \hat{\eta}) = \arg\min_{\beta, \eta} (\|D_n - X_n\beta - (M_n \circ \hat{D}_n)\eta\|_2 / \sqrt{n} + 2\lambda \|\eta\|_1 / \sqrt{n})$$

take the derivative for the second equation:

$$-(M_n \circ \hat{D}_n)'(D_n - X_n\hat{\beta} - (M_n \circ \hat{D}_n)\hat{\eta})/n + \hat{Q}\lambda\hat{\kappa} = 0 \quad (A)$$
$$-X'_n(D_n - X_n\hat{\beta} - (M_n \circ \hat{D}_n)\hat{\eta})/n = 0 \quad (B)$$

where $\hat{Q} = \|D_n - X_n\hat{\beta} - (M_n \circ \hat{D}_n)\hat{\eta}\|_2$

Substitute $D_n = X_n \beta_0 + (M_n \circ D_n) \eta_0 + \epsilon_n$. Equation (A) can be transformed as:

$$\frac{1}{n}(M_n \circ \hat{D}_n)' X_n(\hat{\beta} - \beta_0) + \frac{1}{n}(M_n \circ \hat{D}_n)' \Big((M_n \circ \hat{D}_n)\hat{\eta} - (M_n \circ D_n)\eta_0 \Big) \\ + \hat{Q}\lambda\hat{\kappa} = \frac{(M_n \circ \hat{D}_n)'\epsilon}{n}$$

and further that:

$$\frac{1}{n}(M_n \circ \hat{D}_n)' X_n(\hat{\beta} - \beta_0) + \frac{1}{n}(M_n \circ \hat{D}_n)'(M_n \circ \hat{D}_n)(\hat{\eta} - \eta_0) + \underbrace{\frac{1}{n}(M_n \circ \hat{D}_n)' (M_n \circ (\hat{D}_n - D_n))\eta_0}_{(C)} + Q\lambda\hat{\kappa} = \frac{(M_n \circ \hat{D}_n)'\epsilon}{n}$$

Equation (C) can be written as:

$$\frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ (\hat{D}_n - D_n) \Big) \eta_0 = \frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ (\hat{D}_n - f) \Big) \eta_0 \\ + \frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ (f - D_n) \Big) \eta_0$$

And

$$\frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ (f - D_n) \Big) \eta_0 = -\frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ \epsilon_1 \Big) \eta_0 \\ = -\frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ \eta^0 \Big) (I - M_n \circ \eta_0)^{-1} \epsilon \\ = -\frac{1}{n}(M_n \circ \hat{D}_n)' \left((I - M_n \circ \eta_0)^{-1} - I \right) \epsilon$$

Thus (A) is equivlent to:

$$\frac{1}{n}(M_n \circ \hat{D}_n)' X_n(\hat{\beta} - \beta_0) + \frac{1}{n}(M_n \circ \hat{D}_n)'(M_n \circ \hat{D}_n)(\hat{\eta} - \eta_0) + \frac{1}{n}(M_n \circ \hat{D}_n)' \Big(M_n \circ (\hat{D}_n - f)\Big)\eta_0 + \hat{Q}\lambda\hat{\kappa} = \frac{1}{n}(M_n \circ \hat{D}_n)'(I - M_n \circ \eta_0)^{-1}\epsilon$$

Notice that:

$$\begin{aligned} \left\| \frac{1}{n} (M_n \circ \hat{D}_n)' \Big(M_n \circ (\hat{D}_n - f) \Big) \eta_0 \right\|_{\infty} &\leq \left\| \frac{1}{n} M_n' (M_n \circ \eta_0) (\hat{D}_n - f) \right\|_{\infty} \| \hat{D}_n \|_{\infty} \\ &\leq \frac{1}{n} \left\| M_n' (M_n \circ \eta_0) \right\|_{op2} \left\| (\hat{D}_n - f) \right\|_2 \| \hat{D}_n \|_{\infty} \end{aligned}$$

where $\|.\|_{op2}$ is the operation norm of the matrix in $l_2 \to l_{\infty}$ space, which is the maximum l_2 norm of the row.

From lemma 1,

$$\|(\hat{D}_n - f)\|_2 = O\left(\sqrt{s\log n}\right) = o(n^{1/4})$$

Since each entry of M_n is either 1 or 0, and η_0 has $o(\frac{\sqrt{n}}{\log n})$ non-zero entries. $\|M_n \circ \eta_0\|_{op2} \le o(\frac{n^{1/4}}{\sqrt{\log n}})$. By assumption 5, $\|M_n\|_{op2} = O(\sqrt{\log n})$

$$\left\| M'_n(M_n \circ \eta_0) \right\|_{op2} \le \left\| M_n \right\|_{op2} \left\| M_n \circ \eta_0 \right\|_{op2} = o(n^{1/4})$$

And $\|\eta_0\|_{\infty} < 1$:

$$\left\|\frac{1}{n}(M_n \circ \hat{D}_n)' \left(M_n \circ (\hat{D}_n - f)\right) \eta_0\right\|_{\infty} = o(1/\sqrt{n})$$

Similarly (B) can be transformed in the same way, so (A) and (B) are:

$$\frac{1}{n}(M_n \circ \hat{D}_n)' X_n(\hat{\beta} - \beta_0) + \frac{1}{n}(M_n \circ \hat{D}_n)'(M_n \circ \hat{D}_n)(\hat{\eta} - \eta_0) + \hat{Q}\lambda\hat{\kappa} + o(1/\sqrt{n}) = \frac{(M_n \circ \hat{D}_n)'\epsilon_1}{n} \quad (A') \frac{1}{n}X'_n X_n(\hat{\beta} - \beta_0) + \frac{1}{n}X'_n (M_n \circ \hat{D}_n)(\hat{\eta} - \eta_0) + o(1/\sqrt{n}) = \frac{X'_n \epsilon_1}{n} \quad (B')$$

From (B')

$$(\hat{\beta} - \beta_0) = (X'_n X_n)^{-} X'_n \epsilon_1 - (X'_n X_n)^{-} X'_n (M_n \circ \hat{D}_n) (\hat{\eta} - \eta_0)$$

And substitute this into (A')

$$\frac{1}{n}(M_n \circ \hat{D}_n)' \Big(I - X_n (X'_n X_n)^- X'_n \Big) (M_n \circ \hat{D}_n)(\hat{\eta} - \eta_0) + \hat{Q}\lambda\hat{\kappa} \\ = \frac{1}{n}(M_n \circ \hat{D}_n)' \Big(I - X_n (X'_n X_n)^- X'_n \Big) \epsilon_1$$

Define $W_n = \left(I - X_n (X'_n X_n)^- X'_n\right),$ 1 \tilde{x}'

$$\frac{1}{n}\tilde{X}_1'\tilde{X}_1(\hat{\eta}-\eta_0) + \hat{Q}\lambda\hat{\kappa} = \frac{1}{n}\tilde{X}_1'\epsilon_1$$

where $\tilde{X}_1 = W_n(M_n \circ \hat{D}_n)$.

Define $\hat{\Theta}$ generated from the nodewise regression on \tilde{X}_1 as in Meinshausen and Bühlmann (2006). $\hat{\Theta}$ is a reason able approximation to the inverse of $\tilde{X}'_1 \tilde{X}_1/n$. Thus,

$$\hat{\eta} - \eta_0 + \hat{\Theta}\hat{Q}\lambda\hat{\kappa} = \frac{1}{n}\hat{\Theta}\tilde{X}_1'\epsilon_1 - \Delta/\sqrt{n}$$

where

$$\Delta := \sqrt{n} (\hat{\Theta} \tilde{X}_1' \tilde{X}_1 / n - I) (\hat{\eta} - \eta_0)$$

van de Geer et al. (2014) show that $\|\Delta\|_{\infty} = o_p(1)$ when λ for the nodewise regression is chosen at rate $\sqrt{\log n/n}$

Notice that from (A)

$$\hat{Q}\lambda\hat{\kappa} = (M_n \circ \hat{D}_n)'(D_n - X_n\hat{\beta} - (M_n \circ \hat{D}_n)\hat{\eta})/n$$

Thus

$$\hat{e} = \hat{\eta} + \hat{\Theta}(M_n \circ \hat{D}_n)'(D_n - X_n \hat{\beta} - (M_n \circ \hat{D}_n)\hat{\eta})/n$$
$$= \eta_0 + \frac{1}{n} \hat{\Theta} \tilde{X}'_1 \epsilon_1 - \Delta/\sqrt{n}$$
$$\to \eta_0 \quad \text{as } n \to \infty$$

Similarly

$$\begin{aligned} (\hat{\beta} - \beta_0) &= (X'_n X_n)^{-} X'_n \epsilon_1 - (X'_n X_n)^{-} X'_n (M_n \circ \hat{D}_n) (\hat{\eta} - \eta_0) \\ &= (X'_n X_n)^{-} X'_n \Big(I - (M_n \circ \hat{D}_n) \hat{\Theta} \tilde{X}'_1 / n \Big) \epsilon_1 + (X'_n X_n)^{-} X'_n (M_n \circ \hat{D}_n) \Delta / \sqrt{n} \\ &+ (X'_n X_n)^{-} X'_n (M_n \circ \hat{D}_n) \hat{\Theta} \lambda \hat{\kappa} \end{aligned}$$

 So

$$\begin{aligned} \hat{b} &= \hat{\beta} - (X'_n X_n)^- X'_n (M_n \circ \hat{D}_n)' (D_n - (M_n \circ \hat{D}_n) \hat{\eta} - X_n \hat{\beta}) / n \\ &= \beta_0 + (X'_n X_n)^- X'_n \Big(I - (M_n \circ \hat{D}_n) \hat{\Theta} \tilde{X}'_1 / n \Big) \epsilon_1 + (X'_n X_n)^- X'_n (M_n \circ \hat{D}_n) \Delta / \sqrt{n} \\ &\to \beta_0 \quad \text{as } n \to \infty \end{aligned}$$

Notice that the estimator \hat{b} is a special case in Chernozhukov et al. (2015)

Now consider the design matrix in the second stage, $M_n \circ \hat{D}_n$. Let $(.)_S$ be the operator that restricts a matrix to its columns indexed in S.

Define $\Sigma_{1,1,n}^x = \frac{1}{n} \left(M_n \circ \hat{D}_n \right)'_S \left(M_n \circ \hat{D}_n \right)_S.$ Define $\Sigma_{2,1,n}^x = \frac{1}{n} \left(M_n \circ \hat{D}_n \right)'_{S^c} \left(M_n \circ \hat{D}_n \right)_S.$

$$\Sigma_{1,1,n}^{x} = \frac{1}{n} diag\left((\hat{D}_{n})_{S}\right) \left(M_{n}\right)_{S}' \left(M_{n}\right)_{S} diag\left((\hat{D}_{n})_{S}\right)$$
$$= \frac{1}{n} diag\left((\hat{D}_{n})_{S}\right) \Sigma_{1,1,n} diag\left((\hat{D}_{n})_{S}\right)$$

 So

$$(\Sigma_{1,1,n}^x)^{-1} = n \cdot diag \left((\hat{D}_n)_S \right)^{-1} \Sigma_{1,1,n}^{-1} diag \left((\hat{D}_n)_S \right)^{-1}$$

And,

$$\begin{split} \Sigma_{2,1,n}^{x} &= \frac{1}{n} diag \left((\hat{D}_{n})_{S^{c}} \right) \left(M_{n} \right)_{S^{c}}^{\prime} \left(M_{n} \right)_{S} diag \left((\hat{D}_{n})_{S} \right) \\ &= \frac{1}{n} diag \left((\hat{D}_{n})_{S^{c}} \right) \Sigma_{2,1,n} diag \left((\hat{D}_{n})_{S} \right) \end{split}$$

Thus

$$\left\| \Sigma_{2,1,n}^{x} \left(\Sigma_{1,1,n}^{x} \right)^{-1} \operatorname{sign}(\eta_{0}) \right\|_{\infty} = \left\| \operatorname{diag} \left((\hat{D}_{n})_{S^{c}} \right) \Sigma_{2,1,n} \Sigma_{1,1,n}^{-1} \operatorname{diag} \left((\hat{D}_{n})_{S} \right)^{-1} \operatorname{sign}(\eta_{0}) \right\|_{\infty}$$

Assume $\hat{D}_n \to \Gamma$, $\Sigma_{2,1,n} \Sigma_{1,1,n}^{-1} \to \Sigma$, then I require

$$\|diag(\Gamma_S)\Sigma diag(\Gamma_{S^c})sign(\eta_0)\|_{\infty} < 1$$

The consistency of the active set $\lim_{n\to\infty} \mathbb{P}(\hat{S}_n = S) = 1$ follows from Zhao and Yu (2006) under Assumption 5.

A.2 Theorem 2

Proof. In the presence of multiple networks, we use sparse square root lasso:

$$\min_{\eta} \left\{ \frac{1}{\sqrt{n}} \left\| D_n - \sum_{j=1}^q (M_n^j \circ \hat{D}_n) \eta^j - X_n \beta \right\|_2 + \left(\sum_{j=1}^q \sqrt{T_j} \left(\lambda_1 \| \eta^j \|_2 + \lambda_2 \| \eta^j \|_1 \right) \right) \right\}$$

The KKT condition with respect to the jth group can be written as:

$$\frac{-(M_n^j \circ \hat{D}_n)'(D_n - \sum_{j=1}^q (M_n^j \circ \hat{D}_n)\hat{\eta}^j - X_n\hat{\beta})}{\sqrt{n}\|D_n - \sum_{j=1}^q (M_n^j \circ \hat{D})\hat{\eta}^j - X_n\hat{\beta}\|_2} + \lambda_1\tau^j + \lambda_2\nu^j = 0$$
(A1)

For any $\hat{\beta}_i^j \neq 0$ in group j,

$$\tau_i^j = \frac{\sqrt{T_j}\hat{\eta}_i^j}{\|\hat{\eta}^j\|_2}, \text{ and } \nu_i^j = \sqrt{T_j}sign(\hat{\eta}_i^j)$$

Let $\tau = (\tau_1, \tau_2, \cdots, \tau_p)', \nu = (\nu_1, \nu_2, \cdots, \nu_p)'$. Let $\hat{Z}_n = \left[(M_n^1 \circ \hat{D}_n), (M_n^2 \circ \hat{D}_n), \cdots, (M_n^q \circ \hat{D}_n) \right],$ $Z_n = \left[(M_n^1 \circ D_n), (M_n^2 \circ D_n), \cdots, (M_n^q \circ D_n) \right].$ $\hat{Q} := \|D_n - \sum_{j=1}^q (M_n^j \circ \hat{D}_n) \hat{\eta}^j - X_n \hat{\beta} \|_2.$ Let $\eta = (\eta^{1'}, \eta^{2'}, \cdots, \eta^{q'})'.$ Plug in $D_n = Z_n \eta_0 + X_n \beta_0 + \epsilon_n.$ (A1) can be transformed as:

$$-\frac{\hat{Z}_{n}'\epsilon}{n} + \frac{\hat{Z}_{n}'\hat{Z}_{n}}{n}(\hat{\eta} - \eta_{0}) + \frac{\hat{Z}_{n}'X_{n}}{n}(\hat{\beta} - \beta_{0}) + \underbrace{\frac{1}{n}\hat{Z}_{n}'(\hat{Z}_{n} - Z_{n})\eta_{0}}_{(C*)} + \sqrt{n}\hat{Q}\lambda_{1}\tau + \sqrt{n}\hat{Q}\lambda_{2}\nu = 0 \quad (A2)$$

The derivative with respect to β is

$$-X'_{n}(D_{n} - X_{n}\hat{\beta} - \hat{Z}_{n}\hat{\eta})/n = 0$$

$$\Leftrightarrow \frac{1}{n}X'_{n}X_{n}(\hat{\beta} - \beta_{0}) + \frac{1}{n}X'_{n}\hat{Z}(\hat{\eta} - \eta_{0}) + \frac{1}{n}X'_{n}\Big(\hat{Z}_{n} - Z_{n}\Big)\eta_{0} = \frac{X'_{n}\epsilon}{n}$$
(A3)

Notice that $D_n = (I - \sum_{j=1}^q M_n^j \circ \eta^j)^- X_n \beta + (I - \sum_{j=1}^q M_n^j \circ \eta^j)^- \epsilon = f(\sum_{j=1}^q M_n^j \circ X_n) + \epsilon_1.$

Equation (C^*) can be written as:

$$\frac{1}{n}\hat{Z}'\Big([M_n^1\circ(\hat{D}_n-D_n),M_n^2\circ(\hat{D}_n-D_n),\cdots,M_n^q\circ(\hat{D}_n-D_n)]\Big)\eta_0 \\
= \frac{1}{n}\hat{Z}'\Big([M_n^1\circ(\hat{D}_n-f),M_n^2\circ(\hat{D}_n-f),\cdots,M_n^q\circ(\hat{D}_n-f)]\Big)\eta_0 \\
+ \frac{1}{n}\hat{Z}'\Big([M_n^1\circ(f-D_n),M_n^2\circ(f-D_n),\cdots,M_n^q\circ(f-D_n)]\Big)\eta_0$$

By Theorem 3, $\|\hat{D}_n - f\|_2/\sqrt{n} < M\lambda\sqrt{s_n}$. When $\lambda \asymp \sqrt{\frac{\log n}{n}}$, $\|\hat{D}_n - f\|_2 = o(n^{1/4})$

By assumption 1^\ast and 5^\ast

$$\begin{split} \left\| \frac{1}{n} \hat{Z}' \sum_{j=1}^{q} \left(M_{n}^{j} \circ (\hat{D}_{n} - f) \right) \eta_{0}^{j} \right\|_{\infty} &\leq \sum_{j=1}^{q} \left\| \frac{1}{n} \hat{Z}' \left(M_{n}^{j} \circ (\hat{D}_{n} - f) \right) \eta_{0}^{j} \right\|_{\infty} \\ &\leq \sum_{j=1}^{q} \max_{i=1,\cdots,q} \left\| \frac{1}{n} M_{n}^{i} (M_{n}^{j} \circ \eta_{0}) (\hat{D}_{n} - f) \right\|_{\infty} \| \hat{D}_{n} \|_{\infty} \\ &\leq \frac{1}{n} \sum_{j=1}^{q} \max_{i=1,\cdots,q} \| M_{n}^{i} (M_{n}^{j} \circ \eta_{0}) \|_{op2} \left\| (\hat{D}_{n} - f) \right\|_{2} \| \hat{D}_{n} \|_{\infty} = o(1/\sqrt{n}) \end{split}$$

Since

$$\max_{i=1,\cdots,q} \|M_n^i(M_n^j \circ \eta_0)\|_{op2} \le \max_{i=1,\cdots,q} \|M_n^i\|_{op2} \|(M_n^j \circ \eta_0)\|_{op2} = o(n^{1/4})$$

where $\|.\|_{op2}$ is the operator norm from $l_2 \to l_{\infty}$

And

$$\begin{aligned} \frac{1}{n}\hat{Z}'\Big([M_n^1\circ(f-D_n),M_n^2\circ(f-D_n),\cdots,M_n^q\circ(f-D_n)]\Big)\eta_0\\ &=-\frac{1}{n}\hat{Z}'\Big([M_n^1\circ\epsilon_1,M_n^2\circ\epsilon_1,\cdots,M_n^q\circ\epsilon_1]\Big)\eta_0\\ &=-\frac{1}{n}\hat{Z}'\sum_{j=1}^q\Big(M_n^j\circ\eta_0^j\Big)(I-\sum_{j=1}^qM_n^j\circ\eta_0^j)^-\epsilon\\ &=-\frac{1}{n}\hat{Z}'\left(\Big(I-\sum_{j=1}^qM_n^j\circ\eta_0^j\Big)^--I\right)\epsilon\end{aligned}$$

Thus (A2) and (A3) can be written as:

$$\frac{1}{n}\hat{Z}'_{n}X_{n}(\hat{\beta}-\beta_{0}) + \frac{1}{n}\hat{Z}'_{n}\hat{Z}_{n}(\hat{\eta}-\eta_{0}) + \sqrt{n}\hat{Q}\lambda_{1}\tau + \sqrt{n}\hat{Q}\lambda_{2}\nu + o(1/\sqrt{n}) = \frac{\hat{Z}'_{n}\epsilon_{1}}{n}$$
$$\frac{1}{n}X'_{n}X_{n}(\hat{\beta}-\beta_{0}) + \frac{1}{n}X'_{n}\hat{Z}_{n}(\hat{\eta}-\eta_{0}) + o(1/\sqrt{n}) = \frac{X'_{n}\epsilon_{1}}{n}$$

Define $\tilde{Z}_m = W_n \hat{Z}_n$. Find $\hat{\Theta}_Z$ as an approximation for the inverse of $\tilde{Z}'_m \tilde{Z}_m / n$

$$(\hat{\eta} - \eta_0) + \sqrt{n}\hat{\Theta}_Z Q(\lambda_1 \tau + \lambda_2 \nu) = \hat{\Theta}_Z \tilde{Z}'_m \epsilon_1 / n - \Delta_m / \sqrt{n}$$
$$(\hat{\beta} - \beta_0) - \sqrt{n} (X'_n X_n)^{-1} X'_n \hat{Z}_n \hat{\Theta}_Z Q(\lambda_1 \tau + \lambda_2 \nu) = (X'_n X_n)^{-1} X'_n \Big(I - \hat{Z}_n \hat{\Theta}_Z \tilde{Z}'_m / n \Big) \epsilon_1$$
$$+ (X'_n X_n)^{-1} X'_n \hat{Z}_n \Delta_m / \sqrt{n}$$

where $\Delta_m = \sqrt{n} (\hat{\Theta}_Z \tilde{Z}'_m \tilde{Z}_m / n - I) (\hat{\eta} - \eta_0)$

This suggest the following estimator:

$$\hat{e}_m = \hat{\eta} + \hat{\Theta}_Z \hat{Z}'_n (D_n - \hat{Z}_n \hat{\eta} - X_n \hat{\beta}) / n \to \eta_0$$
$$\hat{b}_m = \hat{\beta} - (X'_n X_n)^{-1} X'_n \hat{Z}_n \hat{\Theta}_Z X'_n (D_n - \hat{Z}_n \hat{\eta} - X_n \hat{\beta}) / n \to \beta_0$$

Now consider the design matrix in the second stage, $\left[\left(M_n^1 \circ \hat{D}_n\right), \cdots, \left(M_n^q \circ \hat{D}_n\right)\right]$. Let $(.)_S$ be the operator that restricts a matrix to its columns indexed in S.

Define
$$\Sigma_{1,1,n}^x = \frac{1}{n} \left[\left(M_n^1 \circ \hat{D}_n \right)_{S_1}, \cdots, \left(M_n^q \circ \hat{D}_n \right)_{S_q} \right]' \left[\left(M_n^1 \circ \hat{D}_n \right)_{S_1}, \cdots, \left(M_n^q \circ \hat{D}_n \right)_{S_q} \right].$$

Define $\tilde{\Sigma}_{2,1,n}^x = \frac{1}{n} \left[\left(\tilde{M}_{S_1}^1 \circ \hat{D}_n \right), \cdots, \left(\tilde{M}_{S_q}^q \circ \hat{D}_n \right) \right]' \left[\left(\tilde{M}_{S_1}^1 \hat{D}_n \right), \cdots, \left(\tilde{M}_{S_q}^q \circ \hat{D}_n \right) \right].$

where $M_{S^c}^{j}$ is defined as M_n^{j} with all non-influential individuals columns being replaced with 0s

Notice that

$$\begin{split} \Sigma_{1,1,n}^{x} &= \frac{1}{n} diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right) \left[\left(M_{n}^{1} \right)_{S_{1}}, \cdots, \left(M_{n}^{q} \right)_{S_{q}} \right]' \left[\left(M_{n}^{1} \right)_{S_{1}}, \cdots, \left(M_{n}^{q} \right)_{S_{q}} \right] \\ & \quad \cdot diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right) \\ &= \frac{1}{n} diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right) \Sigma_{1,1,n} diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right) \end{split}$$

 So

$$\left(\Sigma_{1,1,n}^{x}\right)^{-1} = n \cdot diag\left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}}\right]\right)^{-1} \Sigma_{1,1,n}^{-1} diag\left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}}\right]\right)^{-1} \sum_{n=1}^{\infty} diag\left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}}\right)\right)^{-1} \sum_{n=1}^{\infty} diag\left(\left[(\hat{D}_{n})_{S_{q$$

And,

$$\begin{split} \Sigma_{2,1,n}^{x} &= \frac{1}{n} diag \left(\left[\hat{D}_{n}, \cdots, \hat{D}_{n} \right] \right) \left[\left(\tilde{M}_{n}^{1} \right)_{S_{1}^{c}}, \cdots, \left(\tilde{M}_{n}^{q} \right)_{S_{q}^{c}} \right]' \left[\left(M_{n}^{1} \right)_{S_{1}}, \cdots, \left(M_{n}^{q} \right)_{S_{q}} \right] \\ & \cdot diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right) \\ &= \frac{1}{n} diag \left(\left[\hat{D}_{n}, \cdots, \hat{D}_{n} \right] \right) \tilde{\Sigma}_{2,1,n} diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right) \end{split}$$

Thus

$$\Sigma_{2,1,n}^{x} \left(\Sigma_{1,1,n}^{x} \right)^{-1} = diag \left(\left[\hat{D}_{n}, \cdots, \hat{D}_{n} \right] \right) \tilde{\Sigma}_{2,1,n} \Sigma_{1,1,n}^{-1} diag \left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}} \right] \right)^{-1}$$

Notice that the jth group in the vector

$$\left(\tilde{\Sigma}_{2,1,n}^{x}\left(\Sigma_{1,1,n}^{x}\right)^{-1}u\right)^{j} = diag(\hat{D}_{n})\left(\tilde{\Sigma}_{2,1,n}\Sigma_{1,1,n}^{-1}diag\left(\left[(\hat{D}_{n})_{S_{1}},\cdots,(\hat{D}_{n})_{S_{q}}\right]\right)^{-1}u\right)^{j}$$

$$\max_{\substack{u: \|u\|_{2} \leq \sqrt{n} \text{ } 1 \leq j \leq q}} \max_{1 \leq j \leq q} \frac{\|\left(\tilde{\Sigma}_{2,1,n}^{x} \left(\Sigma_{1,1,n}^{x}\right)^{-1} u\right)^{j}\|_{2}}{\sqrt{n}}$$

=
$$\max_{\substack{u: \|u\|_{2} \leq 1}} \max_{1 \leq j \leq q} \left\| diag(\hat{D}_{n}) \left(\tilde{\Sigma}_{2,1,n} \Sigma_{1,1,n}^{-1} diag\left(\left[(\hat{D}_{n})_{S_{1}}, \cdots, (\hat{D}_{n})_{S_{q}}\right]\right)^{-1} u\right)^{j} \right\|_{2}$$

The consistency of the active set $\lim_{n\to\infty} \mathbb{P}(\hat{S}_n = S) = 1$ follows Theorem 4.

A.3 Theorem 3

Consider the error term in Theorem 1:

$$\frac{(M_n \circ \hat{D}_n)' W_n (I - M_n \circ \eta_0)^{-1} \epsilon}{n} = diag(\hat{D}_n) \frac{M'_n W_n (I - M_n \circ \eta_0)^{-1} \epsilon}{n}$$

And by assumption,

$$\frac{1}{n}M'_{n}W_{n}(I - M_{n} \circ \eta_{0})^{-1}(I - M_{n} \circ \eta_{0})^{-1'}W_{n}M_{n} \to \Omega$$

Thus,

$$\frac{M'_n W_n (I - M_n \circ \eta_0)^- \epsilon}{\sqrt{n}} \to N(0, \Omega)$$

Notice that the limit exist as

$$\left\|\frac{M'_{n}W_{n}(I-M_{n}\circ\eta_{0})^{-}\epsilon}{\sqrt{n}}\right\|_{\infty} \leq \left\|\frac{M'_{n}(I-M_{n}\circ\eta_{0})^{-}\epsilon}{\sqrt{n}}\right\|_{\infty} \|W_{n}\|_{\infty}$$
$$\leq \frac{1}{\sqrt{n}}\|M'_{n}\|_{op1}\|(I-M_{n}\circ\eta_{0})^{-}\|_{op1}\|\epsilon\|_{1}\|W_{n}\|_{\infty}$$

where $\|.\|_{op1}$ norm is the operation norm from $l_1 \to l_{\infty}$, which is the maximum entry in the matrix. Notice that

$$\|(I - M_n \circ \eta_0)^-\|_{op1} = \|\sum_{k=0}^{\infty} (M_n \circ \eta_0)^k\|_{op1}$$
$$\leq \sum_{k=0}^{\infty} \|(M_n \circ \eta_0)\|_{op1}^k$$
$$\leq \frac{1}{1 - \eta_{max}}$$

Also, $\|M'_n\|_{op1} = 1$ and $\|\epsilon\|_1/\sqrt{n} = O(1)$

As a result

$$\left\|\frac{M'_n W_n (I - M_n \circ \eta_0)^- \epsilon}{\sqrt{n}}\right\|_{\infty} \le O(1)$$

And the limit exists.

Let $\Gamma = \lim_{n \to \infty} \hat{D}_n$, $\Theta_1 = \lim_{n \to \infty} \hat{\Theta}$, $\hat{Z}_n = (M_n \circ \hat{D}_n)$, $\tilde{Z}_n = X_n (X'_n X_n)^{-1} X'_n \hat{Z}_n$ and $\Theta_2 = \lim_{n \to \infty} \frac{1}{n} \Big(I - \hat{Z}_n \hat{\Theta} \tilde{Z}'_n / n \Big)' X_n (X'_n X_n)^{-1} X'_n \Big(I - \hat{Z}_n \hat{\Theta} \tilde{Z}'_n / n \Big)$ We have:

$$\sqrt{n}(\hat{e} - \eta_0) = E_1 + \Delta_1,
\sqrt{n}(\hat{b} - \beta_0) = E_2 + \Delta_2,
E_1 \sim N(0, \sigma^2 \Theta_1 diag(\Gamma) \Omega diag(\Gamma) \Theta'_1),
E_2 \sim N(0, \sigma^2 \Theta_2 diag(\Gamma) \Omega diag(\Gamma) \Theta'_2),$$
(19)

where $\|\Delta_1\|_{\infty} = \sqrt{n} (\hat{\Theta} \tilde{X}'_n \tilde{X}_n / n - I) (\hat{\eta} - \eta_0) = o_p(1),$ and $\|\Delta_2\|_{\infty} = (X'_n X_n)^{-1} X'_n (M_n \circ \hat{D}_n) \sqrt{n} (\hat{\Theta} \tilde{X}'_n \tilde{X}_n / n - I) (\hat{\eta} - \eta_0) = o_p(1)$

A.4 Theorem 4

The error term in Theorem 2:

$$\tilde{Z}'_n \epsilon_1 / n = \left[(M_n^1 \circ \hat{D}_n), (M_n^2 \circ \hat{D}_n), \cdots (M_n^q \circ \hat{D}_n) \right]' W_n (I - M_n \circ \eta_0)^{-1} \epsilon / n$$
$$= diag(\hat{D}_n) \left[M_n^1, M_n^2, \cdots M_n^q \right]' W_n (I - M_n \circ \eta_0)^{-1} \epsilon / n$$

and by assumption,

$$\left[M_n^1, M_n^2, \cdots, M_n^q\right]' W_n (I - M_n \circ \eta_0)^{-1} \epsilon / \sqrt{n} \to N(0, \Omega_m)$$

Let
$$\Gamma = \lim_{n \to \infty} \hat{D}_n$$
, $\Theta_{Z1} = \lim_{n \to \infty} \hat{\Theta}_Z$, $\hat{Z}_n = \left[(M_n^1 \circ \hat{D}_n), (M_n^2 \circ \hat{D}_n), \cdots, (M_n^q \circ \hat{D}_n) \right]$,
 $\tilde{Z}_n = X_n (X'_n X_n)^{-1} X'_n \hat{Z}_n$ and $\Theta_{Z2} = \lim_{n \to \infty} \frac{1}{n} \left(I - \hat{Z}_n \hat{\Theta}_Z \tilde{Z}'_n / n \right)' X_n (X'_n X_n)^{-1} X'_n \left(I - \hat{Z}_n \hat{\Theta}_Z \tilde{Z}'_n / n \right)$

We have:

$$\begin{aligned}
\sqrt{n}(\hat{e}_m - \eta_0) &= E_{m1} + \Delta_{m1}, \\
\sqrt{n}(\hat{b}_m - \beta_0) &= E_{m2} + \Delta_{m2}, \\
E_{m1} &\sim N(0, \sigma^2 \Theta_{Z1} diag(\Gamma) \Omega_2 diag(\Gamma) \Theta'_{Z1}), \\
E_{m2} &\sim N(0, \sigma^2 \Theta_{Z2} diag(\Gamma) \Omega_2 diag(\Gamma) \Theta'_{Z2}),
\end{aligned}$$
(20)

where $\|\Delta_{m1}\|_{\infty} = \sqrt{n} (\hat{\Theta}_Z \tilde{Z}'_n \tilde{Z}_n / n - I) (\hat{\eta} - \eta_0) = o_p(1),$ and $\|\Delta_{m2}\|_{\infty} = (X'_n X_n)^- X'_n \hat{Z}_n \sqrt{n} (\hat{\Theta}_Z \tilde{Z}'_n \tilde{Z}_n / n - I) (\hat{\eta} - \eta_0) = o_p(1)$

A.5 Square-root Sparse Group LASSO

To prove Theorem 2 and Theorem 4, we need the following results from square-root sparse group LASSO: 1) Bounds on the prediction, i.e. $\left\|\sum_{j=1}^{q} (M^{j} \circ X_{n})(\hat{\eta}^{j} - \eta_{0}^{j}) + X_{n}(\hat{\beta} - \beta_{0})\right\|_{2} \lesssim \lambda$. And 2) Consistency of selection i.e. $\hat{S}_{n} = S$.

First, define the effective networks as:

$$SG := \{1 \le j \le q : |\eta_0^j|_1 \ne 0\}$$

Define the influential individuals in network j as:

$$S^{j} := \{ 1 \le i \le n : \eta_{0i}^{j} \ne 0 \}$$

Define the number of influential individuals in network j as $|S^j| = s_n^j$. Define the number of the all influential individuals as $|S| = \sum_{j=1}^q s_n^j = s_n$. Define the number of non-zero groups as $|SG| = s_g$. Let $\eta_S \in \mathbb{R}^{s_n}$ be the coefficients for the influential individuals and $\eta_{S^c} \in \mathbb{R}^{nq-s_n}$ be the coefficients for the non-influential individuals.

Theorem 5. Assume $\kappa > 0$, $\gamma > 1$ and $\alpha \in (0, 1)$. Assume $\max_j s_n^j \leq \frac{n}{\log n}$ and $s_g \leq \frac{n}{\log q}$. Let $\lambda = \lambda_1 + \lambda_2$. Under assumptions 1*-5*, the following holds with probability greater than $1 - \alpha$:

$$\left\|\sum_{j=1}^{q} (M^{j} \circ X_{n})(\hat{\eta}^{j} - \eta_{0}^{j}) + X_{n}(\hat{\beta} - \beta_{0})\right\|_{2} \lesssim \frac{\sigma\lambda\sqrt{ns_{n}}}{\kappa},\tag{21}$$

and

$$\sum_{j=1}^{q} \sqrt{T_j} \| (\hat{\eta} - \eta_0)^j \|_2 \lesssim \frac{\sigma \lambda s_n}{\kappa}$$
(22)

Theorem 5 establishes bounds on the LASSO prediction. Notice that Theorem 5 is still valid under a weaker assumption (compatibility assumption) than assumption 4^{*}. The details are shown in the proofs.

The advantage of using sparse group LASSO compare to standard LASSO is that consistent model selection can be achieved under a weaker condition. Standard LASSO requires the l_2 norm of the correlations between all irrelevant regressors and relevant regressors to be small. On the other hand, sparse group LASSO only requires that the l_2 norm of correlations between irrelevant regressors in each group and relevant regressors to be small.

Theorem 6. Define $c = \lambda_1/\lambda$. For constant $\vartheta < 1$, $\alpha \in (0,1)$ and D > 0, if $c < \frac{1-\vartheta}{2}$ and under assumptions 1^*-5^* , with probability greater than $1 - \alpha$:

- 1. $\hat{\eta}_{S^c} = 0$,
- 2. for all $1 \leq j \leq q$,

$$\|(\hat{\eta} - \eta_0)^j\|_{\infty} \le D(c\sqrt{n} + 1 - c)\sigma\lambda_j$$

3. if $\min\{\eta_0^j\} \ge D\sqrt{T_j}\sigma\lambda$, then

$$S = S$$

Theorem 6 shows that consistent selection can be achieved if the design matrix satisfies the Irrepresentable condition together with a Beta-min condition. The ratio between λ_2 and λ_1 is $\frac{1}{c} - 1 > \frac{1+\vartheta}{1-\vartheta}$. Thus when the correlation among the irrelevant regressors in each group and relevant regressors is low, we can penalize more on the l_2 norm and *vice versa*.

A.6 Theorem 5

In the proof of theorem 3 and 4, I consider the following standard lasso problem:

$$\min_{\beta} \left\{ \frac{\|Y - X\beta\|_2}{\sqrt{n}} + \left(\sum_{j=1}^q \sqrt{T_j} \left(\lambda_1 \|\eta_j\|_2 + \lambda_2 \|\eta\|_1 \right) \right) \right\}$$
(A4)

where the true data generating process is $Y = X\beta + \sigma\epsilon$, where ϵ s a mean 0 process with variance 1.

I use β^0 , σ to represent the true parameter values. Let p be the total number of regressors. Let $\{G_1, \dots, G_q\}$ be a partition of $\{1, \dots, p\}$ and $T_i = |G_i|, i = 1, \dots, q$ be the number of regressors in each group. Denote $\beta^j \in \mathbb{R}^{T_j}$ as the coefficients for regressors in group j. Both p, q and T_i s can go to ∞ as $n \to \infty$. Define the active group as:

$$SG := \{1 \le j \le q : |\beta^{0j}|_1 \ne 0\}$$

Define the active set among all regressors as:

$$S := \{1 \le i \le p : \beta_i^0 \ne 0\}$$

Define the size of true support of β^0 as |S| = s; define the number of non-zeros groups as $|SG| = s_g$. Let $\beta_S \in \mathbb{R}^{s_1}$ be the set of coefficients on the true support and $\beta_{S^c} \in \mathbb{R}^{p-s_1}$ be the coefficients for those irrelevant regressors. Define $\hat{Q}(\beta) := \frac{\|Y - X\beta\|_2^2}{n}$. Define $\hat{\delta} := \hat{\beta} - \beta^0$.

The advantage of using square root type lasso is the tuning parameter λ_1 and λ_2 can be chosen independently from σ . In sparse group lasso, the noise component can be viewed in two different ways:

1. I want λ_1 to be sufficiently large to overrule the noise component in grouped lasso, defined as:

$$V = \max_{1 \le j \le q} \left\{ \frac{\sqrt{n} \| (X'\epsilon)^j \|_2}{\sqrt{T_j} \| \epsilon \|_2} \right\}$$

2. I want λ_2 to be sufficiently large to overrule the noise component in standard lasso within each group, defined as:

$$\Lambda = n \left\| \nabla \sqrt{\hat{Q}(\beta^0)} \right\|_{\infty} = \max_{1 \le j \le q} \left\{ \frac{\sqrt{n} \| (X'\epsilon)^j \|_{\infty}}{\|\epsilon\|_2} \right\}$$

Lemma 2. Assume the noise terms ϵ_i are *i.i.d* standard normal random variables. Let $\alpha \in (0, 1)$ be given such that $p/\alpha > 8$ and $n > \log(1/\alpha)$. If

$$\lambda \ge \sqrt{2\log(4p/\alpha)/n}$$

Then

$$\mathbb{P}(\Lambda \ge n\lambda) \le \alpha/2$$

Lemma 3 is a direct result as case (ii) of Lemma 1 in Belloni et al. (2011) Notice that a direct inequality:

$$\mathbb{P}(V \ge n\lambda) \le \mathbb{P}(\Lambda \ge n\lambda)$$

as $||(X'\epsilon)^j||_2/\sqrt{T_j} \le ||(X'\epsilon)^j||_{\infty}$.

Define the event $\mathcal{A}_1 := \{ V \le n\lambda/\bar{\gamma} \}$, the set $\mathcal{A}_2 := \{ \Lambda \le n\lambda/\bar{\gamma} \}$. We can choose

$$\min\{\lambda_1, \lambda_2\} \ge \sqrt{2\log(4p/\alpha)/n}$$

So that:

$$\mathbb{P}(\mathcal{A} := \mathcal{A}_1 \cap \mathcal{A}_2) \ge \mathbb{P}(\mathcal{A}_1) + \mathbb{P}(\mathcal{A}_2) - 1 \ge 1 - \alpha$$

Here, $\bar{\gamma} = \frac{\gamma+1}{\gamma-1}$, where γ is defined as below.

Define

$$\Delta_{\gamma}^{1} := \{ \delta \in \mathbb{R}^{p} : \sum_{j \in SG^{c}} \sqrt{T_{j}} \|\delta^{j}\|_{2} \leq \gamma \sum_{j \in SG} \sqrt{T_{j}} \|\delta^{j}\|_{2} \}$$
$$\Delta_{\gamma}^{2} := \{ \delta \in \mathbb{R}^{p} : \sum_{j \in SG^{c}} \sqrt{T_{j}} \|\delta_{S^{c}}^{j}\|_{1} \leq \gamma \sum_{j \in SG} \sqrt{T_{j}} \|\delta_{S}^{j}\|_{1} \}$$

Compatibility Condition (CC). We say that the Compatibility Condition is met for $\kappa > 0$ and $\|X\hat{\delta}\|_2$ $\gamma>1$ if :

$$\sum_{j \in SG} \sqrt{T_j} \|\hat{\delta}_S^j\|_1 \le \frac{\sqrt{s_n} \|X\delta\|_2}{\sqrt{n\kappa}}$$

for all $\delta \in \Delta^1_{\gamma} \cap \Delta^2_{\gamma}$

Proof. • First, by definition of (A4)

$$\sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^{0})} \leq \underbrace{\lambda_{1} \sum_{j=1}^{q} \sqrt{T_{j}} (\|\beta^{0j}\|_{2} - \|\hat{\beta}^{j}\|_{2})}_{(1)} + \underbrace{\lambda_{2} \sum_{j=1}^{q} \sqrt{T_{j}} (\|\beta^{0j}\|_{1} - \|\hat{\beta}^{j}\|_{1})}_{(1)}$$
(A5)

$$(1) = \lambda_1 \sum_{j \in SG} \sqrt{T_j} (\|\beta^{0j}\|_2 - \|\hat{\beta}^j\|_2) - \lambda_1 \sum_{j \in SG^c} \sqrt{T_j} (\|\hat{\beta}^j\|_2)$$

$$\leq \lambda_1 \sum_{j \in SG} \sqrt{T_j} (\|\|\beta^{0j}\|_2 - \|\hat{\beta}^j\|_2|) - \lambda_1 \sum_{j \in SG^c} \sqrt{T_j} (\|\hat{\beta}^j\|_2)$$

$$\leq \lambda_1 \sum_{j \in SG} \sqrt{T_j} (\|\|\hat{\delta}^j\|_2) - \lambda_1 \sum_{j \in SG^c} \sqrt{T_j} (\|\hat{\delta}^j\|_2)$$

$$(2) = \lambda_2 \sum_{j=1}^{q} \sqrt{T_j} (\|\beta_S^{0j}\|_1 - \|\hat{\beta}_S^j\|_1 - \|\hat{\beta}_{S^c}^j\|_1)$$

$$\leq \lambda_2 \sum_{j=1}^{q} \sqrt{T_j} (\|\hat{\delta}_S^j\|_1 - \|\hat{\delta}_{S^c}^j\|_1)$$

• Second, by convexity,

$$\sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^0)} \ge \nabla \sqrt{\hat{Q}(\beta^0)}(\hat{\delta}) \ge -\frac{|\epsilon' X \hat{\delta}|}{\sqrt{n} \|\epsilon\|_2}$$

$$\begin{split} |\epsilon' X \hat{\delta}| &= \Big| \sum_{j=1}^{q} \epsilon' X^{j} \hat{\delta}^{j} \Big| \leq \sum_{j=1}^{q} \| (\epsilon' X^{j})' \|_{2} \| \hat{\delta}^{j} \|_{2} \\ &\leq \max_{1 \leq j \leq q} \left\{ \frac{\sqrt{n} \| (\epsilon' X^{j})' \|_{2}}{\sqrt{T_{j}} \| \epsilon \|_{2}} \right\} \frac{\| \epsilon \|_{2}}{\sqrt{n}} \sum_{j=1}^{q} \sqrt{T_{j}} \| \hat{\delta}^{j} \|_{2} \\ &= V \frac{\| \epsilon \|_{2}}{\sqrt{n}} \sum_{j=1}^{q} \sqrt{T_{j}} \| \hat{\delta}^{j} \|_{2} \end{split}$$

Also

$$\begin{split} |\epsilon' X \hat{\delta}| &= \big| \sum_{j=1}^{q} \epsilon' X^{j} \hat{\delta}^{j} \big| \leq \sum_{j=1}^{q} \| (\epsilon' X^{j})' \|_{\infty} \| \hat{\delta}^{j} \|_{1} \\ &\leq \left\| \frac{\sqrt{n} \| (\epsilon' X^{j})' \|_{\infty}}{\sqrt{T_{j}} \| \epsilon \|_{2}} \right\|_{\infty} \frac{\| \epsilon \|_{2}}{\sqrt{n}} \sum_{j=1}^{q} \sqrt{T_{j}} \| \hat{\delta}^{j} \|_{1} \\ &= \frac{\Lambda}{\sqrt{T_{j}}} \frac{\| \epsilon \|_{2}}{\sqrt{n}} \sum_{j=1}^{q} \sqrt{T_{j}} \| \hat{\delta}^{j} \|_{1} \end{split}$$

On set \mathcal{A} , we have $\lambda/\bar{\gamma} \geq V$, Thus

$$\sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^0)} \ge -\frac{\lambda}{\bar{\gamma}} \sum_{j=1}^q \sqrt{T_j} \|\hat{\delta}^j\|_2 \tag{A6}$$

Again, on set \mathcal{A} , we also have $\lambda/\bar{\gamma} \geq \Lambda/\sqrt{T_{min}}$, Thus

$$\sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^0)} \ge -\frac{\lambda}{\bar{\gamma}} \sum_{j=1}^q \sqrt{T_j} (\|\hat{\delta}_S^j\|_1 + \|\hat{\delta}_{S^c}^j\|_1)$$
(A7)

• Third, Combine (A6) and (A7), for any $c \in [0, 1]$

$$\sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^{0})} \ge -\frac{c\lambda}{\bar{\gamma}} \sum_{j=1}^{q} \sqrt{T_{j}} \|\hat{\delta}^{j}\|_{2} - \frac{(1-c)\lambda}{\bar{\gamma}} \sum_{j=1}^{q} \sqrt{T_{j}} (\|\hat{\delta}_{S}^{j}\|_{1} + \|\hat{\delta}_{S^{c}}^{j}\|_{1})$$
(A8)

Set $\lambda_1 = c\lambda$ and $\lambda_2 = (1 - c)\lambda$, we can combine (A8) and (A5) to get

$$c\lambda \sum_{j \in SG} \sqrt{T_j} (|\|\hat{\delta}^j\|_2) - c\lambda \sum_{j \in SG^c} \sqrt{T_j} (\|\hat{\delta}^j\|_2) + (1-c)\lambda \sum_{j=1}^q \sqrt{T_j} (\|\hat{\delta}^j_S\|_1 - \|\hat{\delta}^j_{S^c}\|_1)$$

$$\geq -\frac{c\lambda}{\bar{\gamma}} \sum_{j=1}^q \sqrt{T_j} \|\hat{\delta}^j\|_2 - \frac{(1-c)\lambda}{\bar{\gamma}} \sum_{j=1}^q \sqrt{T_j} (\|\hat{\delta}^j_S\|_1 + \|\hat{\delta}^j_{S^c}\|_1)$$

Thus,

$$\left(1 + \frac{1}{\bar{\gamma}}\right) c\lambda \sum_{j \in SG} \sqrt{T_j} \|\hat{\delta}^j\|_2 + \left(1 + \frac{1}{\bar{\gamma}}\right) (1 - c)\lambda \sum_{j=1}^q \sqrt{T_j} \|\hat{\delta}^j_S\|_1$$

$$\geq \left(1 - \frac{1}{\bar{\gamma}}\right) c\lambda \sum_{j \in SG^c} \sqrt{T_j} \|\hat{\delta}^j\|_2 + \left(1 - \frac{1}{\bar{\gamma}}\right) (1 - c)\lambda \sum_{j=1}^q \sqrt{T_j} \|\hat{\delta}^j_{S^c}\|_1$$

which implies:

$$c\gamma \sum_{j \in SG} \sqrt{T_j} \|\hat{\delta}^j\|_2 + (1-c)\gamma \sum_{j \in SG^c} \sqrt{T_j} \|\hat{\delta}^j_S\|_1$$

$$\geq c \sum_{j \in SG^c} \sqrt{T_j} \|\hat{\delta}^j\|_2 + (1-c) \sum_{j \in SG^c} \sqrt{T_j} \|\hat{\delta}^j_{S^c}\|_1$$
(A9)

(A9) $\Rightarrow \hat{\delta} \in \Delta^1_{\gamma} \cap \Delta^2_{\gamma}$. Thus,

$$\sum_{j \in SG} \sqrt{T_j} \|\hat{\delta}^j\|_2 \le \sum_{j \in SG} \sqrt{T_j} \|\hat{\delta}_S^j\|_1 \le \frac{\sqrt{s_n} \|X\hat{\delta}\|_2}{\sqrt{n\kappa}}$$
(A10)

 $\bullet\,$ Forth, from (A10),

$$\begin{split} \sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^0)} &\leq \lambda_1 \sum_{j \in SG} \sqrt{T_j} (\|\hat{\delta}^j\|_2) - \lambda_1 \sum_{j \in SG^c} \sqrt{T_j} (\|\hat{\delta}^j\|_2) \\ &+ \lambda_2 \sum_{j=1}^q \sqrt{T_j} (\|\hat{\delta}^j_S\|_1 - \|\hat{\delta}^j_{S^c}\|_1) \\ &\leq c\lambda \sum_{j \in SG} \sqrt{T_j} (\|\hat{\delta}^j\|_2) + (1-c)\lambda \sum_{j=1}^q \sqrt{T_j} \|\hat{\delta}^j_S\|_1 \\ &\leq \lambda \frac{\sqrt{s_n} \|X\delta\|_2}{\sqrt{n\kappa}} \end{split}$$
(A11)

• Fifth,

$$\begin{split} \hat{Q}(\hat{\beta}) - \hat{Q}(\beta_0) &= \frac{1}{n} \left((Y - X\hat{\beta})'(Y - X\hat{\beta}) - (Y - X\beta_0)'(Y - X\beta_0) \right) \\ &= \frac{1}{n} \{ (Y - X\hat{\beta})'(Y - X\hat{\beta}) \\ &- (Y - X\hat{\beta} + X\hat{\beta} - X\beta_0)'(Y - X\beta_0) \} \\ &= \frac{1}{n} \{ (Y - X\hat{\beta})'(Y - X\hat{\beta}) \\ &- (Y - X\hat{\beta})'(Y - X\beta_0) - X(\hat{\beta} - \beta_0)'(Y - X\beta_0) \} \\ &= \frac{1}{n} \{ -(Y - X\hat{\beta})'X(\hat{\beta} - \beta_0) - X(\hat{\beta} - \beta_0)'(Y - X\beta_0) \} \\ &= \frac{1}{n} \{ (\hat{\beta} - \beta_0)'X'X(\hat{\beta} - \beta_0) - 2X(\hat{\beta} - \beta_0)'(Y - X\beta_0) \} \\ &= \frac{\|X\hat{\delta}\|_2^2}{n} - \frac{2\sigma\epsilon' X\hat{\delta}}{n} \end{split}$$

• Sixth, from (A11),

$$\begin{split} \frac{\|X\hat{\delta}\|_2^2}{n} &= \hat{Q}(\hat{\beta}) - \hat{Q}(\beta_0) + \frac{2\sigma\epsilon' X\hat{\delta}}{n} \\ &= \left(\sqrt{\hat{Q}(\hat{\beta})} - \sqrt{\hat{Q}(\beta^0)}\right) \left(\sqrt{\hat{Q}(\hat{\beta})} + \sqrt{\hat{Q}(\beta^0)}\right) + \frac{2\sigma\epsilon' X\hat{\delta}}{n} \\ &\leq \lambda \frac{\sqrt{s_n} \|X\delta\|_2}{\sqrt{n}\kappa} \left(2\sqrt{\hat{Q}(\beta^0)} + \lambda \frac{\sqrt{s_n} \|X\delta\|_2}{\sqrt{n}\kappa}\right) + 2V \frac{\|\sigma\epsilon\|_2}{n^{3/2}} \sum_{j=1}^q \sqrt{T_j} \|\hat{\delta}^j\|_2 \end{split}$$

From (A10):

$$\leq \frac{s_n\lambda}{\kappa^2 n} \frac{\|X\delta\|_2^2}{n} + 2\lambda \frac{\|\sigma\epsilon\|_2}{\sqrt{n}} \frac{\sqrt{s_n} \|X\delta\|_2}{\sqrt{n\kappa}} + 2V \frac{\|\sigma\epsilon\|_2}{n^{3/2}} \frac{\sqrt{s_n} \|X\hat{\delta}\|_2}{\sqrt{n\kappa}}$$
$$n\lambda/\bar{\gamma} \geq V:$$
$$\leq \frac{s_n\lambda}{\kappa^2 n} \frac{\|X\delta\|_2^2}{n} + 2\left(1 + \frac{1}{\bar{\gamma}}\right) \frac{\|\sigma\epsilon\|_2}{\sqrt{n}} \lambda \frac{\sqrt{s_n} \|X\hat{\delta}\|_2}{\sqrt{n\kappa}}$$

As a result:

$$\left(1 - \left(\frac{\lambda\sqrt{s_1}}{\kappa}\right)^2\right) \frac{\|X\hat{\delta}\|_2^2}{n} \le 2\left(1 + \frac{1}{\bar{\gamma}}\right) \frac{\|\sigma\epsilon\|_2}{\sqrt{n}} \lambda \frac{\sqrt{s_1}\|X\hat{\delta}\|_2}{\sqrt{n}\kappa}$$
(A12)

(A12) concludes the first statement in Theorem 3.

For the second claim, use the fact that $\delta \in \Delta^1_{\gamma}$ and the Compatibility Condition:

$$\sum_{j=1}^{q} \sqrt{T_j} \|\delta^j\|_2 \le (\gamma+1) \sum_{j \in SG^c} \sqrt{T_j} \|\delta^j\|_2 \le \frac{(\gamma+1)\sqrt{s_n} \|X\delta\|_2}{\sqrt{n\kappa}}$$
$$\lesssim \frac{(\gamma+1)\sqrt{s_n}}{\sqrt{n\kappa}} \frac{\sqrt{n\sigma\lambda}\sqrt{s_n}}{\kappa} \lesssim \frac{\sigma\lambda s_n}{\kappa}$$

-	_

Lemma 3. Assumption^{*} (4) implies Irrepresentable Condition: for $0 < \vartheta < 1$ if $\Sigma_{1,1}$ is invertible and

$$\max_{u:\|u\|_{\infty} \le \sqrt{T_k}} \max_{1 \le j \le q} \frac{\|\left(\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} u\right)^j\|_{\infty}}{\sqrt{T_j}} \le \vartheta$$

for all k.

Proof.

$$\|v\|_2 \le \sqrt{T_k} \Rightarrow \|v\|_\infty \le \sqrt{T_k}$$

Thus, Assumption* (4) \Rightarrow

$$\max_{u:\|u\|_{\infty} \le \sqrt{T_k}} \max_{1 \le j \le q} \frac{\|\left(\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} u\right)^j\|_2}{\sqrt{T_j}} \le \vartheta$$

Since

$$\|\left(\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}u\right)^{j}\|_{2} \ge \|\left(\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}u\right)^{j}\|_{\infty}$$

Thus

$$\max_{u:\|u\|_{\infty} \le \sqrt{T_k}} \max_{1 \le j \le q} \frac{\|\left(\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} u\right)^j\|_{\infty}}{\sqrt{T_j}} \le \vartheta$$

A.7 Theorem 6

From Lemma 2 and Lemma 3, Define $\mathcal{B}_1 = \{V \leq n\lambda/(\bar{\gamma} \vee 2\bar{\vartheta})\}$ and $\mathcal{B}_2 = \{\Lambda \leq \sqrt{T_{min}}\lambda/(\bar{\gamma} \vee 2\bar{\vartheta})\}$. I can choose

$$\min\{\lambda_1, \lambda_2\} \ge \max\left\{ (\bar{\gamma} \lor 2\vartheta) \sqrt{2\log(4p/\alpha)/n}, \frac{\sqrt{n}(\bar{\gamma} \lor 2\vartheta)}{\sqrt{T_{min}}} \sqrt{2\log(4p/\alpha)/n} \right\}$$

So that:

$$\mathbb{P}(\mathcal{B} := \mathcal{B}_1 \cap \mathcal{B}_2) \ge \mathbb{P}(\mathcal{B}_1) + \mathbb{P}(\mathcal{B}_2) - 1 \ge 1 - \alpha$$
(A13)

Here, $\bar{\vartheta} = \frac{1+\vartheta}{1-\vartheta}$, where ϑ is defined assumption 4^* .

Proof. Choose λ big enough so that (A13) holds.

• First take the derivative of (A4) with respect to each column *i*:

$$\frac{-(X_i^j)'(Y-X\hat{\beta})}{\|Y-X\hat{\beta}\|_2} = \sqrt{n\lambda_1\tau_i^j} + \sqrt{n\lambda_2\nu_i^j}$$
(A14)

Let $\tau = (\tau_1, \tau_2, \cdots, \tau_p)', \nu = (\nu_1, \nu_2, \cdots, \nu_p)' \hat{Q} := ||Y - X\hat{\beta}||_2$ For any $\hat{\beta}_i^j \neq 0$ in group j,

$$\tau_i^j = \frac{\sqrt{T_j}\hat{\beta}_i^j}{\|\hat{\beta}^j\|_2}, \text{ and } \nu_i^j = \sqrt{T_j}sign(\hat{\beta}_i^j)$$

By KKT condition, $\|\tau^j\|_2 \leq \sqrt{T_j}$ and $|\nu_i^j| \leq \sqrt{T_j}$.

 $Y = X\beta + \epsilon$. Let $d\hat{elta} = \hat{\beta} - \beta^0$. We can rewrite (A14) in matrix form:

$$\sigma X'\epsilon - X'X\hat{\delta} = \sqrt{n}\hat{Q}\lambda_1\tau + \sqrt{n}\hat{Q}\lambda_2\nu \tag{A15}$$

or

$$\begin{pmatrix} \sigma(X'\epsilon)_S \\ \sigma(X'\epsilon)_{S^c} \end{pmatrix} - n \begin{pmatrix} \Sigma_{1,1} & \Sigma_{1,2} \\ \Sigma_{2,1} & \Sigma_{2,2} \end{pmatrix} \begin{pmatrix} \hat{\delta}_S \\ \hat{\delta}_{S^c} \end{pmatrix} = \begin{pmatrix} (\sqrt{n}\hat{Q}\lambda_1\tau + \sqrt{n}\hat{Q}\lambda_2\nu)_S \\ (\sqrt{n}\hat{Q}\lambda_1\tau + \sqrt{n}\hat{Q}\lambda_2\nu)_{S^c} \end{pmatrix}$$

• Second, the upper part of (A15) can be transform to

$$-n\Sigma_{1,1}\hat{\delta}_S - n\Sigma_{1,2}\hat{\delta}_{S^c} = \sqrt{n}\hat{Q}(\lambda_1\tau_S + \lambda_2\nu_S) - \sigma(X'\epsilon)_S$$

or, equivalently,

$$-n\hat{\delta}'_{S^{c}}\Sigma_{2,1}\hat{\delta}_{S} - n\hat{\delta}'_{S^{c}}\Sigma_{2,1}\Sigma_{1,1}^{-1}\Sigma_{1,2}\hat{\delta}_{S^{c}} = \sqrt{n}\hat{Q}\hat{\delta}'_{S^{c}}\Sigma_{2,1}\Sigma_{1,1}^{-1}(\lambda_{1}\tau_{S} + \lambda_{2}\nu_{S}) - \sigma\hat{\delta}'_{S^{c}}\Sigma_{2,1}\Sigma_{1,1}^{-1}(X'\epsilon)_{S}$$
(A16)

Notice that for all $\hat{\delta}_i^j \neq 0$ but $i \in S^c$, either $j \in SG^c$ or $j \in SG$.

Define

$$SG_1 \subset SG := \{1 \le j \le q : \exists \beta_i^{0j} = 0\}$$

Define

$$S^{jc} := \{ 1 \le i \le T_j : \beta_i^{0j} = 0 \}$$

Let $l_j = |S^{jc}|$ denotes the size of the sparsity in group j.

The right hand side of (A16) can be broken into two parts. The first part consider all sparse term in nonzero groups while the second term consider all zero groups:

$$(A16) = \underbrace{\sqrt{n}\hat{Q}\lambda}_{j\in SG_{1}}\sum_{i\in S^{jc}}\hat{\delta}_{i}^{j}\left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}\left(c\tau_{S}+(1-c)\nu_{S}-\frac{\sigma}{\sqrt{n}\lambda\hat{Q}}(X'\epsilon)_{S}\right)\right]_{i}^{j}}_{(1)} + \underbrace{\sqrt{n}\hat{Q}\lambda}_{j\in SG^{c}}\sum_{i\in S^{jc}}\hat{\delta}_{i}^{j}\left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}\left(c\tau_{S}+(1-c)\nu_{S}-\frac{\sigma}{\sqrt{n}\lambda\hat{Q}}(X'\epsilon)_{S}\right)\right]_{i}^{j}}_{(2)}$$

$$(1) = \underbrace{\sqrt{n}\hat{Q}\lambda c}_{j\in SG_{1}} \sum_{i\in S^{jc}} \hat{\delta}_{i}^{j} \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} \left(\tau_{S} - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}} (X'\epsilon)_{S} \right) \right]_{i}^{j}}_{(3)} + \underbrace{\sqrt{n}\hat{Q}\lambda(1-c)}_{j\in SG_{1}} \sum_{i\in S^{jc}} \hat{\delta}_{i}^{j} \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} \left(\nu_{S} - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}} (X'\epsilon)_{S} \right) \right]_{i}^{j}}_{(4)}$$

$$(2) = \underbrace{\sqrt{n}\hat{Q}\lambda c}_{j\in SG^c} \sum_{i\in S^{jc}} \hat{\delta}_i^j \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} \left(\tau_S - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}} (X'\epsilon)_S \right) \right]_i^j }_{(5)} + \underbrace{\sqrt{n}\hat{Q}\lambda(1-c)}_{j\in SG^c} \sum_{i\in S^{jc}} \hat{\delta}_i^j \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} \left(\nu_S - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}} (X'\epsilon)_S \right) \right]_i^j }_{(6)}$$

By Holder:

$$(3) \leq \sqrt{n}\hat{Q}\lambda c \sum_{j \in SG_1} \left\{ \sum_{i \in S^{jc}} |\hat{\delta}_i^j| \left\| \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} \left(\tau_S - \frac{\sigma}{\sqrt{n}\lambda \hat{Q}} (X'\epsilon)_S \right) \right]_i^j \right\|_{\infty} \right\}$$

Observe again that if $n\lambda/2\bar{\vartheta} \ge \hat{\Lambda}/\sqrt{T_{min}} \Rightarrow \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}((X_i^j)'\epsilon) \le \frac{\sqrt{T_j}}{2\bar{\vartheta}}$ for any i and

$$\|\tau^j\|_{\infty} = \left\|\frac{\sqrt{T_j}\hat{\beta}_i^j}{\|\hat{\beta}^j\|_2}\right\|_{\infty} \le \sqrt{T_j}$$

By Lemma 4

$$(3) \leq \sqrt{n}\hat{Q}\lambda c \max_{u:\|u\|_{\infty} \leq \left(\sqrt{T_{j}} + \frac{\sqrt{T_{j}}}{2\vartheta}\right)} \sum_{j \in SG_{1}} \left\{ \left(\sum_{i \in S^{jc}} |\hat{\delta}_{i}^{j}|\right) \left\| \left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}u\right]^{j}\right\|_{\infty} \right\}$$
$$\leq \left(1 + \frac{1}{2\vartheta}\right)\sqrt{n}\hat{Q}\lambda c \max_{u:\|u\|_{\infty} \leq \sqrt{T_{j}}} \sum_{j \in SG_{1}} \sqrt{T_{j}} \left\{ \left(\sum_{i \in S^{jc}} |\hat{\delta}_{i}^{j}|\right) \frac{\left\| \left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}u\right]^{j}\right\|_{\infty}}{\sqrt{T_{j}}} \right\}$$
$$\leq \vartheta \left(1 + \frac{1}{2\vartheta}\right)\sqrt{n}\hat{Q}\lambda c \sum_{j \in SG_{1}} \sqrt{T_{j}} \left(\sum_{i \in S^{jc}} |\hat{\delta}_{i}^{j}|\right)$$

By Holder:

$$(4) \le \sqrt{n}\hat{Q}\lambda(1-c)\sum_{j\in SG_1}\left\{\left(\sum_{i\in S^{jc}}|\hat{\delta}_i^j|\right)\left\|\left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}\left(\nu_S - \frac{\sqrt{n}\sigma}{\lambda\hat{Q}}(X'\epsilon)_S\right)\right]^j\right\|_{\infty}\right\}$$

Observe that if $\lambda/2\bar{\vartheta} \ge \hat{\Lambda}/\sqrt{T_{min}} \Rightarrow \frac{\sqrt{n\sigma}}{\lambda\hat{Q}}((X^i)'\epsilon) \le \frac{\sqrt{T_j}}{2\bar{\vartheta}}$ for any i and $\|\nu_i^j\|_{\infty} \le \sqrt{T_j}$.

$$(4) \leq \sqrt{n}\hat{Q}\lambda(1-c) \max_{u:\|u\|_{\infty} \leq \left(1+\frac{1}{2\vartheta}\right)\sqrt{T_{k}}} \sum_{j \in SG_{1}} \left\{ \left(\sum_{i \in S^{jc}} |\hat{\delta}_{i}^{j}|\right) \left\| \left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}u\right]^{j}\right\|_{\infty} \right\}$$

$$\leq \left(1+\frac{1}{2\vartheta}\right)\sqrt{n}\hat{Q}\lambda(1-c) \max_{u:\|u\|_{\infty} \leq \sqrt{T_{k}}} \sum_{j \in SG_{1}} \sqrt{T_{j}} \left\{ \left(\sum_{i \in S^{jc}} |\hat{\delta}_{i}^{j}|\right) \frac{\left\| \left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}u\right]^{j}\right\|_{\infty}}{\sqrt{T_{j}}} \right\}$$

$$\leq \vartheta \left(1+\frac{1}{2\vartheta}\right)\sqrt{n}\hat{Q}\lambda(1-c) \sum_{j \in SG_{1}} \sqrt{T_{j}} \left\{ \left(\sum_{i \in S^{jc}} |\hat{\delta}_{i}^{j}|\right)\right\}$$

Since $\|\tau^j\|_2 \leq \sqrt{T_j}$ and $n\lambda/\bar{\vartheta} \geq \lambda/2\bar{\vartheta} \geq \hat{V} \Rightarrow \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \|(X'\epsilon)^j\|_2 \leq \frac{\sqrt{T_j}}{\bar{\vartheta}}$

$$(5) \leq \sqrt{n}\hat{Q}\lambda c \sum_{j\in SG^{c}} \sqrt{T_{j}} \left(\sum_{i\in S^{jc}} (\hat{\delta}_{i}^{j})^{2} \right)^{1/2} \frac{\left\| \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} \left(\tau_{S} - \frac{\sigma}{\sqrt{n}\lambda \hat{Q}} (X'\epsilon)_{S} \right) \right]^{j} \right\|_{2}}{\sqrt{T_{j}}} \\ \leq \sqrt{n}\hat{Q}\lambda c \max_{v: \|v^{k}\|_{2} \leq (1+\frac{1}{\vartheta})\sqrt{T_{k}}} \sum_{j\in SG^{c}} \sqrt{T_{j}} \|\hat{\delta}^{j}\|_{2} \frac{\left\| \left[\tilde{\Sigma}_{2,1} \Sigma_{1,1}^{-1} v \right]^{j} \right\|_{2}}{\sqrt{T_{j}}} \\ \leq \vartheta \left(1 + \frac{1}{\vartheta} \right) c \sqrt{n} \hat{Q}\lambda \sum_{j\in SG^{c}} \sqrt{T_{j}} \|\hat{\delta}^{j}\|_{2} \\ = \left(1 - \frac{1}{\vartheta} \right) c \sqrt{n} \hat{Q}\lambda \sum_{j\in SG^{c}} \sqrt{T_{j}} \|\hat{\delta}^{j}\|_{2}$$
For any
$$i \in S$$
, $\nu_i^j = \sqrt{T_j} sign(\beta_i^j)$. Thus $|\nu_i^j| \leq \sqrt{T_j}$. $n\lambda/\bar{\vartheta} \geq n\lambda/2\bar{\vartheta} \geq \hat{\Lambda}/\sqrt{T_{min}} \Rightarrow$
 $\frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \| (X'\epsilon)^j \|_{\infty} \leq \frac{\sqrt{T_j}}{\vartheta}$

$$(6) \leq \sqrt{n}\hat{Q}\lambda(1-c) \sum_{j\in SG^c} \|\hat{\delta}^j\|_1 \left\| \left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1} \left(\nu_S - \frac{\sigma}{\sqrt{n\lambda}\hat{Q}}(X'\epsilon)_S \right) \right]^j \right\|_{\infty}$$

$$\leq \sqrt{n}\hat{Q}\lambda c \max_{v:\|v^k\|_2 \leq (1+\frac{1}{\vartheta})\sqrt{T_k}} \sum_{j\in SG^c} \sqrt{T_j} \|\hat{\delta}^j\|_1 \frac{\left\| \left[\tilde{\Sigma}_{2,1}\Sigma_{1,1}^{-1}v \right]^j \right\|_{\infty}}{\sqrt{T_j}}$$

$$\leq \vartheta \left(1 + \frac{1}{\vartheta} \right) (1-c)\sqrt{n}\hat{Q}\lambda \sum_{j\in SG^c} \sqrt{T_j} \|\hat{\delta}^j\|_1$$

$$= \left(1 - \frac{1}{\vartheta} \right) (1-c)\sqrt{n}\hat{Q}\lambda \sum_{j\in SG^c} \sqrt{T_j} \|\hat{\delta}^j\|_1$$

• Third, the bottom part of (A15) can be transform to:

$$-n\Sigma_{2,1}\hat{\delta}_S - n\Sigma_{2,2}\hat{\delta}_{S^c} = \sqrt{n}\hat{Q}(\lambda_1\tau_{S^c} + \lambda_2\nu_{S^c}) - \sigma(X'\epsilon)_{S^c}$$

or, equivalently,

$$-n\hat{\delta}'_{S^c}\Sigma_{2,1}\hat{\delta}_S - n\hat{\delta}'_{S^c}\Sigma_{2,2}\hat{\delta}_{S^c} = \sqrt{n}\hat{Q}\hat{\delta}'_{S^c}(\lambda_1\tau_{S^c} + \lambda_2\nu_{S^c}) - \sigma\hat{\delta}'_{S^c}(X'\epsilon)_{S^c}$$
(A17)

For $j \in SG_1$ and $i \in S^{jc}$,

$$\hat{\beta}_{i}^{j} \neq 0 \Rightarrow \hat{\delta}_{i}^{j}(c\tau_{i}^{j} + (1-c)\nu_{i}^{j}) = c\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\beta}^{j}\|_{2}} + (1-c)\sqrt{T_{j}}|\delta_{i}^{j}|$$

$$\hat{\beta}_{i}^{j} = 0 \Rightarrow \hat{\delta}_{i}^{j}(c\tau_{i}^{j} + (1-c)\nu_{i}^{j}) = 0 = c\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\beta}^{j}\|_{2}} + (1-c)\sqrt{T_{j}}|\delta_{i}^{j}|$$

For $j \in SG^c$ and all i,

$$\hat{\beta}_{i}^{j} \neq 0 \Rightarrow \hat{\delta}_{i}^{j}(c\tau_{i}^{j} + (1-c)\nu_{i}^{j}) = c\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\delta}^{j}\|_{2}} + (1-c)\sqrt{T_{j}}|\delta_{i}^{j}|$$
$$\hat{\beta}_{i}^{j} = 0 \Rightarrow \hat{\delta}_{i}^{j}(c\tau_{i}^{j} + (1-c)\nu_{i}^{j}) = 0 = c\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\delta}^{j}\|_{2}} + (1-c)\sqrt{T_{j}}|\delta_{i}^{j}|$$

As in the previous section, the right hand side of (A17) can be broken into two parts:

$$(A18) = \underbrace{\sqrt{n}\hat{Q}\lambda}_{j\in SG_{1}}\sum_{i\in S^{jc}} \left(c\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\beta}^{j}\|_{2}} + (1-c)\sqrt{T_{j}}|\delta_{i}^{j}| - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_{i}^{j}(X'\epsilon)_{i}^{j}\right)}_{(7)} + \underbrace{\sqrt{n}\hat{Q}\lambda}_{j\in SG^{c}}\sum_{i\in S^{jc}} \left(c\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\delta}^{j}\|_{2}} + (1-c)\sqrt{T_{j}}|\delta_{i}^{j}| - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_{i}^{j}(X'\epsilon)_{i}^{j}\right)}_{(8)}$$

$$(7) = \underbrace{\sqrt{n}\hat{Q}\lambda c}_{j\in SG_{1}} \sum_{i\in S^{jc}} \left(\frac{\sqrt{T_{j}}(\hat{\delta}_{i}^{j})^{2}}{\|\hat{\beta}^{j}\|_{2}} - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_{i}^{j}(X'\epsilon)_{i}^{j} \right)}_{(9)} + \underbrace{\sqrt{n}\hat{Q}\lambda(1-c)}_{j\in SG_{1}} \sum_{i\in S^{jc}} \left(\sqrt{T_{j}}|\delta_{i}^{j}| - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_{i}^{j}(X'\epsilon)_{i}^{j} \right)}_{(10)}$$

$$(8) = \underbrace{\sqrt{n}\hat{Q}\lambda c}_{j\in SG^c} \sum_{i\in S^{jc}} \left(\frac{\sqrt{T_j}(\hat{\delta}_i^j)^2}{\|\hat{\delta}^j\|_2} - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_i^j(X'\epsilon)_i^j \right)}_{(11)} + \underbrace{\sqrt{n}\hat{Q}\lambda(1-c)}_{j\in SG^c} \sum_{i\in S^{jc}} \left(\sqrt{T_j}|\delta_i^j| - \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_i^j(X'\epsilon)_i^j \right)}_{(12)}$$

By Holder, we have

$$(9) = \sqrt{n}\hat{Q}\lambda c \sum_{j\in SG_1} \left(\sum_{i\in S^{jc}} \frac{\sqrt{T_j}(\hat{\delta}_i^j)^2}{\|\hat{\beta}^j\|_2} - \sum_{i\in S^{jc}} \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_i^j(X'\epsilon)_i^j \right)$$
$$\geq \sqrt{n}\hat{Q}\lambda c \sum_{j\in SG_1} \left\{ \frac{\sum_{i\in S^{jc}} \sqrt{T_j}(\hat{\delta}_i^j)^2}{\|\hat{\beta}^j\|_2} - \left(\sum_{i\in S^{jc}} |\hat{\delta}_i^j|\right) \left\| \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}(X'\epsilon)_i^j \right\|_{\infty} \right\}$$

Observe again that if $n\lambda/2\bar{\vartheta} \ge \hat{\Lambda}/\sqrt{T_{min}} \Rightarrow \frac{\sigma}{\sqrt{n}\lambda\hat{Q}}((X^i)'\epsilon) \le \frac{\sqrt{T_j}}{2\bar{\vartheta}}$ for any i

$$(9) \geq \sqrt{n}\hat{Q}\lambda c \sum_{j\in SG_1} \left\{ \frac{\sqrt{T_j}\sum_{i\in S^{jc}} (\hat{\delta}_i^j)^2}{\|\hat{\beta}^j\|_2} - \frac{\sqrt{T_j}}{2\bar{\vartheta}} \Big(\sum_{i\in S^{jc}} |\hat{\delta}_i^j|\Big) \right\}$$
$$\geq -\frac{1}{2\bar{\vartheta}}\sqrt{n}\hat{Q}\lambda c \sum_{j\in SG_1} \sqrt{T_j} \Big(\sum_{i\in S^{jc}} |\hat{\delta}_i^j|\Big)$$

$$(10) = \sqrt{n}\hat{Q}\lambda(1-c)\left\{\sqrt{T_j}\sum_{j\in SG_1}\left(\sum_{i\in S^{jc}}|\delta_i^j|\right) - \left(\sum_{i\in S^{jc}}\frac{\sigma}{\sqrt{n}\lambda\hat{Q}}\hat{\delta}_i^j(X'\epsilon)_i^j\right)\right\}$$
$$\geq \sqrt{n}\hat{Q}\lambda(1-c)\sum_{j\in SG_1}\left\{\sqrt{T_j}\left(\sum_{i\in S^{jc}}|\delta_i^j|\right) - \left(\sum_{i\in S^{jc}}|\hat{\delta}_i^j|\right)\left\|\frac{\sigma}{\sqrt{n}\lambda\hat{Q}}(X'\epsilon)_i^j\right\|_{\infty}\right\}$$
$$\geq \left(1-\frac{1}{2\overline{\vartheta}}\right)\sqrt{n}\hat{Q}\lambda(1-c)\sum_{j\in SG_1}\sqrt{T_j}\left(\sum_{i\in S^{jc}}|\hat{\delta}_i^j|\right)$$

Since $n\lambda/\bar{\vartheta} \ge \lambda/2\bar{\vartheta} \ge \hat{V} \Rightarrow \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} ||(X'\epsilon)^j||_2 \le \frac{\sqrt{T_j}}{\vartheta}$ for any j: $(11) = \sqrt{n}\hat{Q}\lambda c \sum_{j\in SG^c} \left(\sqrt{T_j} ||\hat{\delta}^j||_2 - \sum_{i\in S^{jc}} \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \hat{\delta}_i^j (X'\epsilon)_i^j\right)$ $\ge \sqrt{n}\hat{Q}\lambda c \sum_{j\in SG^c} \left(\sqrt{T_j} ||\hat{\delta}^j||_2 - \sqrt{T_j} ||\hat{\delta}^j||_2 \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \frac{||(X'\epsilon)^j||_2}{\sqrt{T_j}}\right)$ $\ge \left(1 - \frac{1}{\vartheta}\right) c\sqrt{n}\hat{Q}\lambda \sum_{j\in SG^c} \sqrt{T_j} ||\hat{\delta}^j||_2$

Since
$$n\lambda/\bar{\vartheta} \ge \hat{\Lambda}/\sqrt{T_{min}} \Rightarrow \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \| (X'\epsilon)^j \|_{\infty} \le \frac{\sqrt{T_j}}{\vartheta}$$
 for any j:

$$(12) = \sqrt{n}\hat{Q}\lambda(1-c) \sum_{j\in SG^c} \left(\sqrt{T_j} \|\hat{\delta}^j\|_1 - \sum_{i\in S^{jc}} \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \hat{\delta}_i^j (X'\epsilon)_i^j\right)$$

$$\ge \sqrt{n}\hat{Q}\lambda(1-c) \sum_{j\in SG^c} \left(\sqrt{T_j} \|\hat{\delta}^j\|_1 - \|\hat{\delta}^j\|_1 \frac{\sigma}{\sqrt{n\lambda}\hat{Q}} \| (X'\epsilon)^j\|_{\infty}\right)$$

$$\ge \left(1 - \frac{1}{\bar{\vartheta}}\right) (1-c)\sqrt{n}\hat{Q}\lambda \sum_{j\in SG^c} \sqrt{T_j} \|\hat{\delta}^j\|_1$$

Subtract (A17) from (A16) and notice (11) and (12) canceled with (5) and (6), we have:

$$\begin{split} n^{2} \hat{\delta}_{S^{c}}^{\prime}(\Sigma_{2,2} - \Sigma_{2,1} \Sigma_{1,1}^{-1} \Sigma_{1,2}) \hat{\delta}_{S^{c}} \\ &\leq \left\{ -\left(1 - \frac{1}{2\bar{\vartheta}}\right) (1 - c) + \frac{1}{2\bar{\vartheta}} c + \vartheta \left(1 + \frac{1}{2\bar{\vartheta}}\right) c + \vartheta \left(1 + \frac{1}{2\bar{\vartheta}}\right) (1 - c) \right\} \\ &\cdot \sqrt{n} \hat{Q} \lambda \sum_{j \in SG_{1}} \left(\sum_{i \in S^{j_{c}}} |\hat{\delta}_{i}^{j}|\right) \\ &= \left\{ -1 + \frac{1}{2\bar{\vartheta}} + \vartheta \left(1 + \frac{1}{2\bar{\vartheta}}\right) + c + \vartheta \left(1 + \frac{1}{2\bar{\vartheta}}\right) c - \vartheta \left(1 + \frac{1}{2\bar{\vartheta}}\right) c \right\} \\ &\cdot \sqrt{n} \hat{Q} \lambda \sum_{j \in SG_{1}} \left(\sum_{i \in S^{j_{c}}} |\hat{\delta}_{i}^{j}|\right) \\ &= \left\{ -\frac{1}{2} (1 - \vartheta) + c \right\} \sqrt{n} \hat{Q} \lambda \sum_{j \in SG_{1}} \left(\sum_{i \in S^{j_{c}}} |\hat{\delta}_{i}^{j}|\right) \\ &\leq 0 \end{split}$$

The last inequality is due to Substitution Condition.

However, since $\Sigma_{2,2} - \Sigma_{2,1} \Sigma_{1,1}^{-1} \Sigma_{1,2} \ge 0$, this implies $\hat{\delta}_{S^c} = 0$, which establish the first claim.

For the second claim, substitute $\hat{\delta}_{S^c} = 0$ into (A16) we have:

$$-n\hat{\delta}_S = \sqrt{n}\hat{Q}\lambda\Sigma_{1,1}^{-1}\left(c\tau_S + (1-c)\nu_S - \frac{\sigma(X'\epsilon)_S}{\sqrt{n}\hat{Q}\lambda}\right)$$

Recall again $n\lambda/\bar{\vartheta} \ge \hat{\Lambda}/\sqrt{T_{min}}$

$$\begin{split} \|\hat{\delta}^{j}\|_{\infty} &\leq \frac{\hat{Q}\lambda c}{n^{1/2}} \left\| \left[\tilde{\Sigma}_{1,1}^{-1} \left(\tau_{S} - \frac{\sigma(X'\epsilon)_{S}}{\sqrt{n}\hat{Q}\lambda} \right) \right]^{j} \right\|_{\infty} \\ &+ \frac{\hat{Q}\lambda(1-c)}{n^{1/2}} \left\| \left[\tilde{\Sigma}_{1,1}^{-1} \left(\nu_{S} - \frac{\sigma(X'\epsilon)_{S}}{\sqrt{n}\hat{Q}\lambda} \right) \right]^{j} \right\|_{\infty} \\ &\leq \frac{\hat{Q}\lambda c}{n^{1/2}} \max_{v: \|v^{k}\|_{2} \leq (1+\frac{1}{\vartheta})\sqrt{T_{j}}} \left\| \left[\tilde{\Sigma}_{1,1}^{-1}v \right]^{j} \right\|_{\infty} \\ &+ \frac{\hat{Q}\lambda(1-c)}{n^{1/2}} \max_{u: \|u\|_{\infty} \leq (1+\frac{1}{\vartheta})\sqrt{T_{j}}} \left\| \left[\tilde{\Sigma}_{1,1}^{-1}u \right]^{j} \right\|_{\infty} \\ &\leq (1+\frac{1}{\vartheta}) \frac{\hat{Q}\lambda}{n^{1/2}} \sqrt{T_{j}} \max_{u: \|u\|_{\infty} \leq \sqrt{T_{j}}} \frac{\left\| \left[\tilde{\Sigma}_{1,1}^{-1}u \right]^{j} \right\|_{\infty}}{\sqrt{T_{j}}} \\ &\leq D\sqrt{T_{j}}\sigma\lambda \end{split}$$

The third claim follows with Beta Min Condition.

B Useful Algebra Transformation

B.1 (4)

Since $(M \circ D)\eta = (M \circ \eta)D$,

$$D_{n} = (M_{n} \circ D_{n})\eta_{0} + X_{n}\beta_{0} + \epsilon_{n}$$

$$\Leftrightarrow D_{n} = (M_{n} \circ \eta_{0})D_{n} + X_{n}\beta_{0} + \epsilon_{n}$$

$$\Leftrightarrow (I_{n} - (M_{n} \circ \eta_{0}))D_{n} = X_{n}\beta_{0} + \epsilon_{n}$$

$$\Leftrightarrow D_{n} = (I_{n} - (M_{n} \circ \eta_{0}))^{-}(X_{n}\beta_{0} + \epsilon_{n})$$

$$\Leftrightarrow D_{n} = \sum_{i=0}^{\infty} (M_{n} \circ \eta_{0})^{i}(X_{n}\beta_{0} + \epsilon_{n})$$

B.2 (5)

$$E(D_n) = \sum_{i=0}^{\infty} \left(M_n \circ \eta_0 \right)^i \beta_0 X_n$$

= $\beta_0 X_n + \beta_0 \left(M_n \circ \eta_0 \right) X_n + \sum_{i=2}^{\infty} \left(M_n \circ \eta_0 \right)^i \beta_0 X_n$
= $X_n \beta_0 + \left(M_n \circ X_n \right) (\beta_0 \eta_0) + \sum_{i=2}^{\infty} \left(M_n \circ \eta_0 \right)^i \beta_0 X_n$

$$\mathbf{B.3} \quad (6)$$

Let $M_n = (m_1, m_2, \cdots, m_n)$, where m^j is the *j*th column of M_n . $\eta_0 = (\eta_1, \eta_2, \cdots, \eta_n)'$ Then $\left(M_n \circ \eta_0\right)^2 = \left(M_n \circ \eta_0\right) \left(m_1\eta_1, m_2\eta_2, \cdots, m_n\eta_n\right)$ $= \left[\left(M_n \circ \eta_0\right)m_1\eta_1, \left(M_n \circ \eta_0\right)m_2\eta_2, \cdots, \left(M_n \circ \eta_0\right)m_n\eta_n\right]$ $= \left[\left(M_n \circ m_1\right)\eta_0\eta_1, \left(M_n \circ m_2\right)\eta_0\eta_2, \cdots, \left(M_n \circ m_n\right)\eta_0\eta_n\right]$

Thus

$$\begin{pmatrix} M_n \circ \eta_0 \end{pmatrix}^2 \beta_0 X_n = \begin{pmatrix} M_n \circ m_1 \end{pmatrix} \eta_0 \eta_1 x_{n1} \beta_0 + \begin{pmatrix} M_n \circ m_2 \end{pmatrix} \eta_0 \eta_2 x_{n2} \beta_0 + \dots + \begin{pmatrix} M_n \circ m_n \end{pmatrix} \eta_0 \eta_n x_{nn} \beta_0$$
$$= \begin{pmatrix} M_n \circ m_1 \end{pmatrix} \eta_0 \delta_1^1 + \begin{pmatrix} M_n \circ m_2 \end{pmatrix} \eta_0 \delta_2^1 + \dots + \begin{pmatrix} M_n \circ m_n \end{pmatrix} \eta_0 \delta_n^1$$

$$\begin{pmatrix} M_n \circ \eta_0 \end{pmatrix}^3 \beta_0 X_n = \sum_{i=1}^n \left(M_n \circ \eta_0 \right) \left(M_n \circ m_i \right) \eta_0 \delta_i^1$$

= $\sum_{i=1}^n \left(M_n \circ \eta_0 \right) \left(m_1 m_{i1} \eta_1 \delta_i^1 + m_2 m_{i2} \eta_2 \delta_i^1 + \dots + m_n m_{in} \eta_n \delta_i^1 \right)$
= $\sum_{i=1}^n m_{i1} \left(M_n \circ m_1 \right) \eta_0 \delta_i^1 + m_{i2} \left(M_n \circ m_2 \right) \eta_0 \delta_i^1 + \dots + m_{in} \left(M_n \circ m_n \right) \eta_0 \delta_i^1$
= $\sum_{i=1}^n \sum_{j=1}^n \left(M_n \circ m_j \right) \eta_0 m_{ij} \delta_i^1$
= $\sum_{i=1}^n \left(M_n \circ m_i \right) \eta_0 \delta_i^2$

With induction, one can show that

$$\left(M_n \circ \eta_0\right)^k \beta_0 X_n = \sum_{i=1}^n \left(M_n \circ m_i\right) \eta_0 \delta_i^{k-1}$$

Thus,

$$E(D_n|X) = X_n\beta_0 + \left(M_n \circ X_n\right)(\beta_0\eta_0) + \sum_{i=1}^n \left(M_n \circ m_i\right)\eta_0\delta_i^\infty$$

where $\delta_i^{\infty} = \sum_{j=1}^{\infty} \delta_i^j$.

When M_n is the adjacency matrix, the *j*th column of $(M_n \circ m_i)$ is 0 if $m_{ij} = 0$ or *i* is not connect with *j*; and is equal to m_j if $m_{ij} = 1$

Thus

where

$$\begin{split} E(D_n) &= X_n \beta_0 + \left(m_1 x_1 \eta_1 \beta_0 + m_2 x_2 \eta_2 \beta_0 + \dots + m_n x_n \eta_n \beta_0 \right) \\ &+ \sum_{i=1}^n \left(m_1 m_{i1} \eta_1 \delta_i^\infty + m_2 m_{i2} \eta_2 \delta_i^\infty + \dots + m_n m_{in} \eta_n \delta_i^\infty \right) \\ &= X_n \beta_0 + \left(m_1 x_1 \eta_1 \left(\beta_0 + \sum_{i=1}^n \frac{m_{i1} \delta_i^\infty}{x_1} \right) + m_2 x_2 \eta_2 \left(\beta_0 + \sum_{i=1}^n \frac{m_{i2} \delta_i^\infty}{x_2} \right) + \dots \\ &+ m_n x_n \eta_n \left(\beta_0 + \sum_{i=1}^n \frac{m_{in} \delta_i^\infty}{x_n} \right) \right) \\ &= X_n \beta_0 + \left(M_n \circ X_n \right) \tilde{\eta} \\ \tilde{\eta}_j &= \eta_j \left(\beta_0 + \sum_{i=1}^n \frac{m_{ij} \delta_i^\infty}{x_j} \right). \end{split}$$

As a result, using $(M_n \circ X_n)$ and X_n are sufficient to determine the influential individuals.

B.4 (8)

$$D_{n} = (M_{n} \circ D_{n})\eta_{0} + \gamma M_{n}D_{n} + X_{n}\beta_{0} + \epsilon_{n}$$

$$\Leftrightarrow D_{n} = (M_{n} \circ \eta_{0})D_{n} + \gamma M_{n}D_{n} + X_{n}\beta_{0} + \epsilon_{n}$$

$$\Leftrightarrow (I_{n} - (M_{n} \circ \eta_{0}) - \gamma M_{n})D_{n} = X_{n}\beta_{0} + \epsilon_{n}$$

$$\Leftrightarrow D_{n} = (I_{n} - (M_{n} \circ \eta_{0}) - \gamma M_{n})^{-}(X_{n}\beta_{0} + \epsilon_{n})$$

$$\Leftrightarrow D_{n} = \sum_{i=0}^{\infty} (M_{n} \circ \eta_{0} + \gamma M_{n})^{i}(X_{n}\beta_{0} + \epsilon_{n})$$

C Multiple Networks Assumptions

Assumption* 1. Among n individuals in q_n networks, let S_n^j be the set of influential individuals in network j. Let $s_n^j = |S_n^j|$ be the number of elements in S_n^j .

$$s_n^j = o\left(\frac{\sqrt{n}}{\log n}\right), \quad as \ n \to \infty$$
$$s_g = \sum_{j=1}^{q_n} \mathbf{1}(s_n^j \neq 0) = o\left(\frac{n}{\log q_n}\right), \quad as \ n \to \infty$$

Notice same individual from different networks are counted as different elements in S_n Assumption* 2.

- There exists an $\eta_{\max} < 1$ such that $\sum_{j=1}^{q} \|\eta_0^j\|_{\infty} \leq \eta_{\max}$
- The ϵ_j are i.i.d with 0 mean and variance σ^2
- The regressors x_i in X_n are uniformly bounded constants for all n. $\lim_{n\to\infty} X'_n X_n/n$ exists and is nonsingular

Apply the same algebra:

$$D_n = \sum_{j=1}^q \left(M_n^j \circ D_n \right) \eta_0^j + X_n \beta_0 + \epsilon_n$$

$$\Leftrightarrow D_n = \sum_{j=1}^q \left(M_n^j \circ \eta_0^j \right) D_n + X_n \beta_0 + \epsilon_n$$

$$\Leftrightarrow \left(I - \sum_{j=1}^q \left(M_n^j \circ \eta_0^j \right) \right) D_n = X_n \beta_0 + \epsilon_n$$

$$\Leftrightarrow D_n = \left(I - \sum_{j=1}^q \left(M_n^j \circ \eta_0^j \right) \right)^- (X_n \beta_0 + \epsilon_n)$$

$$\Leftrightarrow D_n = \sum_{i=0}^\infty \left(\sum_{j=1}^q M_n^j \circ \eta_0^j \right)^i (X_n \beta_0 + \epsilon_n)$$

Consider X_n as a one dimensional vector:

$$E(D_n) = \sum_{i=0}^{\infty} \left(\sum_{j=1}^q M_n^j \circ \eta_0^j \right)^i \beta_0 X_n$$

= $\beta_0 X_n + \beta_0 \left(\sum_{j=1}^q M_n^j \circ \eta_0^j \right) X_n + \sum_{i=2}^{\infty} \left(\sum_{j=1}^q M_n^j \circ \eta_0^j \right)^i \beta_0 X_n$
= $X_n \beta_0 + \sum_{j=1}^q \left(M_n^j \circ X_n \right) (\beta_0 \eta_0^j) + \sum_{i=2}^{\infty} \left(\sum_{j=1}^q M_n^j \circ \eta_0^j \right)^i \beta_0 X_n$

 $X_n, (M_n^1 \circ X_n), (M_n^2 \circ X_n), \cdots, (M_n^q \circ X_n)$ are valid instruments.

Apply the same algebra

$$E(D_n) = X_n \beta_0 + \sum_{j=1}^q \left(M_n^j \circ X_n \right) \tilde{\eta}^j$$

Assumption* 3. $\left[X_n, \left(M_n^1 \circ X_n\right)_S, \left(M_n^2 \circ X_n\right)_S, \cdots, \left(M_n^q \circ X_n\right)_S\right]$ is full rank with probability equals to 1.

We can use Group Lasso to identify those influential individuals and the networks that deliver the influence. For Group Lasso to achieve consistent selection, we need the following assumption:

Assumption* 4.

(Group Irrepresentable Condition) There exists $N \in \mathbb{N}$: $\forall n \geq N$, there is a $\vartheta \in (0,1)$ such that

$$P\left(\max_{u:\|u\|_{2}\leq 1}\max_{1\leq j\leq q}\left\|diag(\hat{D}_{n})\left(\tilde{\Sigma}_{2,1,n}\Sigma_{1,1,n}^{-1}diag\left(\left[(\hat{D}_{n})_{S_{1}},\cdots,(\hat{D}_{n})_{S_{q}}\right]\right)^{-1}u\right)^{j}\right\|_{2}\leq\vartheta\right)=1$$

(Beta Min Condition) There exists $N \in \mathbb{N}$: $\forall n \geq N$, there is a m > 0 such that

$$\min|(\eta_0)_S| \ge m/\sqrt{n}$$

Assumption^{*} 5.

(Maximum Neighbors Condition)

$$\|M_n^{j'}\mathbf{1}_n\|_{\infty} \le O(\log n) \quad for \ all \ j$$

(Variance Condition)

$$\frac{1}{n}M_n^{0'}W_n(I - \sum_{j=1}^q M_n^j \circ \eta_0^j)^{-1}(I - \sum_{j=1}^q M_n^j \circ \eta_0^j)^{-1'}W_nM_n^0 \to \Omega_2$$

where $M_n^0 = [M_n^1, M_n^2, \cdots, M_n^q]$, and $W_n = (I - X_n (X'_n X_n)^{-1} X'_n)$

Define
$$\Sigma_{1,1,n} = \frac{1}{n} \left[\left(M_n^1 \right)_S, \cdots, \left(M_n^q \right)_S \right]' \left[\left(M_n^1 \right)_S, \cdots, \left(M_n^q \right)_S \right]$$

Define $\Sigma_{2,1,n} = \frac{1}{n} \left[\left(M_n^1 \right)_{S^c}, \cdots, \left(M_n^q \right)_{S^c} \right]' \left[\left(M_n^1 \right)_S, \cdots, \left(M_n^q \right)_S \right].$
Define $\tilde{\Sigma}_{2,1,n} = \frac{1}{n} \left[\left(\tilde{M}_{S^c}^1 \right), \cdots, \left(\tilde{M}_{S^c}^q \right) \right]' \left[\left(\tilde{M}_S^1 \right), \cdots, \left(\tilde{M}_S^q \right) \right].$

where $\tilde{M}_{S^c}^j$ is defined as M_n^j with all non-influential individuals columns being replaced with 0s

D Adjacency Matrix for Influential Individuals

We use the following adjacency matrix for influential individuals when there are five of them:

 $\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$

We use the following adjacency matrix for influential individuals when there are ten of them:

_									-
0	0	0	1	0	0	0	0	0	1
0	0	0	1	1	1	0	0	0	0
0	0	0	1	0	0	0	0	0	0
1	1	1	0	0	0	0	0	0	0
0	1	0	0	0	0	1	0	0	0
0	1	0	0	0	0	1	0	1	0
0	0	0	0	1	1	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0	0	0
1	0	0	0	0	0	0	0	0	0

E Centraility

Denote a graph as G = (V, E), where V represents the set for vertex and E represents the set for edges. Define a measure $d : (x, y) \to \mathbb{R}$ as the length of the shortest path between the node x and y. And define M as the adjacency matrix for graph G. I consider the following centrality measures:

• Degree centrality

The degree centrality $C_D(x)$ of a vertex V is defined as the number of edges connected to node v.

• Closeness centrality

The Closeness centrality is the average length of the shortest path between the node and all other nodes in the graph.

$$C_C(v) = \frac{1}{\sum_{y \in V} d(v, y)}$$

• Betweenness centrality

Betweenness centrality measures the number of times a node acts as a bridge along the shortest path between two other nodes.

$$C_B(v) = \sum_{x \neq v \neq y \in V} \frac{\sigma_{x,y}(v)}{\sigma_{x,y}}$$

where $\sigma_{x,y}$ is total number of shortest paths from node x to node y. $\sigma_{x,y}(v)$ is the number of those paths that pass through v.

• Eigenvector centrality

Eigenvector centrality is defined as the left-hand eigenvector of the adjacency matrix M associated with the largest eigenvalue λ :

$$\lambda x = xM$$

And the vth entry of x is the eigenvector centrality of v.

F Tables

Table	8:	Simulation

		0.1				0.2	
Network Size	50	200	500		50	200	500
Avgcov S_0	0.9780	0.9560	0.9380	0.	9770	0.9480	0.9580
Avglength S_0	2.9420	3.6734	2.6136	1.	0179	3.3098	2.0386
Avgcov S_0^c	0.9222	0.9861	0.9846	0.	9920	0.9861	0.9884
Avglength S_0^c	18.9664	8.1006	2.5444	21	.4923	3.1052	1.9782
A 2	0.0700	0.0700	0.0050	0	0500	0.0050	0.0000
Avgcov β	0.8700	0.9700	0.9650	0.	9500	0.9650	0.9800
Avglength β	4.0056	0.4890	0.2959	0.	9773	0.7905	0.5209
Domon 1	0.9140	0 1 2 0 0	0.4650	0	5970	0.9590	0 1770
Power ¹	0.2140	0.1800	0.4000	0.	5870	0.2520	0.1770
FDR ²	0.0147	0.0001	0.0000	0.	0017	0.0000	0.0030
Awgoon 1	0.0000	0.0000	0.0000	0	0500	0.0000	0.0850
AvgCOV 1	0.9900	0.9000 E 6446	0.9900	0.	9000 7909	4 8020	1.0251
Avglength 1	3.0138	0.0440	0.8000	0.	1202	4.8020	1.9551
Avecov 2	0 9700	0 9700	0 8600	0	9850	0 9650	0 9550
Avglength 2	5.0537	2 6632	1 5617	1	5222	2 7718	2 51/2
Avgiengen 2	0.0001	2.0052	1.5017	1.	0222	2.1110	2.0142
Avgcov 3	0.9800	0.9400	0.9150	0.	9900	0.9250	0.9900
Avglength 3	2.4503	4.8645	4.3329	0.	7604	4.2660	1.5686
0 0							
Avgcov 4	0.9600	0.9800	0.9650	0.	9950	0.9500	0.9950
Avglength 4	1.8805	3.6298	4.1772	0.	8212	3.6354	1.6983
Avgcov 5	0.9900	0.9900	0.9600	0.	9650	1.0000	0.8650
Avglength 5	2.3115	1.5647	2.0417	1.	2652	1.0741	2.4768

This table summarizes the results simulated on Erdos-Renyi type random graphs. When a node is added into the graph, it has probability p = 0.1 or p = 0.2 to form a link with all existing nodes.

The reported coverage is from 200 simulations from 50, 200 and 500 nodes graphs. The active set S_0 contains 5 nodes and coverage for each is reported as Avgcov 1-5. Nonactive set S_0^c contains all remaining nodes.

1. Power represents the averaged percentage in the active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

2. FDR reports the averaged percentage in the non-active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

Table	9:	Simulation

		0.1			0.2	
Network Size	50	200	500 3	50	200	500 3
Avgcov S_0	0.9730	0.9870	0.8805	0.6905	0.9870	0.8330
Avglength S_0	11.8263	1.5870	4.5802	0.8400	3.6207	2.1104
Avgcov S_0^c	0.9942	0.9905	0.9638	0.9827	0.9972	0.9733
Avglength S_0^c	23.2425	2.5128	4.0562	9.3871	2.9423	5.5000
Avgcov β	0.9800	0.9700	0.9300	0.9500	0.9950	0.9950
Avglength p	2.0520	0.5205	0.9008	1.2024	0.5201	0.7915
Power ¹	0.0475	0.1680	0.5725	0.8175	0.1620	0.4710
FDR 2	0.0000	0.0001	0.0000	0.0102	0.0000	0.0004
Avgcov 1	0.9250	1.0000	0.9200	0.8150	0.9950	0.8800
Avglength 1	6.8760	0.9064	1.6038	0.5357	0.8415	2.0727
Avgcov 2	0.9600	1.0000	0.9350	0.8350	1.0000	0.8550
Avglength 2	5.4525	1.0753	17.1486	0.3364	0.8152	1.2254
Avgcov 3	0.9950	0.9900	0.9300	0.9750	0.9650	0.6600
Avglength 3	12.2042	2.8843	2.5877	2.0093	1.8027	1.4287
0.0						
Avgcov 4	0.9600	0.9600	0.9450	0.4150	0.9750	0.8300
Avglength 4	17.9822	1.8353	1.5234	0.3987	2.1970	1.7792
Avgcov 5	0.9900	1.0000	0.8600	1.0000	1.0000	0.8150
Avglength 5	3.0831	1.0112	1.1228	0.4478	1.0124	0.7712
Avecov 6	0.9750	0.0000	0.8850	0 5250	1.0000	0.9100
Avglength 6	39 1903	2 1300	14 2821	0.3260	20 5070	3 9918
				0.0001		
Avgcov 7	0.9750	0.9600	0.8500	0.2800	0.9600	0.8250
Avglength 7	19.0146	2.5613	3.2858	0.8466	5.3259	5.0870
Avgcov 8	0.9650	1.0000	0.9100	0.7800	0.9850	0.7950
Avglength 8	3.2928	1.2513	1.2913	0.3709	1.4046	1.3949
	0.0000	0.0700	0.0050	0.0050	0.0050	0.0500
Avgcov 9 Avglopath 0	0.9900	0.9700	0.9250	0.9250	0.9950	0.9500
Avgiengtin 9	ə.əə29	1.1019	1.4152	2.9014	1.4000	1.2109
Avgcov 10	0.9950	1.0000	0.6450	0.3550	0.9950	0.8100
Avglength 10	7.3340	1.1136	1.5418	0.2474	0.8831	2.1366

This table summarizes the results simulated on Erdos-Renyi type random graphs. When a node is added into the graph, it has probability p = 0.1 or p = 0.2 to form a link with all existing nodes.

The reported coverage is from 200 simulations from 50, 200 and 500 nodes graphs. The active set S_0 contains 10 nodes and coverage for each is reported as Avgcov 1-10. Nonactive set S_0^c contains all remaining nodes.

1. Power represents the averaged percentage in the active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

2. FDR reports the averaged percentage in the non-active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

 $3.\,$ For 500 cases, lasso tuning parameter is chosen using rule of thumb instead of cross-validation

	0.04			0.08		
Network Size	50	200	500	50	200	500
Avgcov S_0	0.9180	0.8490	0.9920	0.9410	0.8310	0.9860
Avglength S_0	5.7298	1.6333	1.6646	5.2069	3.8309	0.9132
Avgcov S_0^c	0.9543	0.9646	0.9809	0.9577	0.9581	0.9949
Avglength S_0^c	7.8860	5.2748	4.3686	3.4985	2.9435	3.4044
Avgcov β	0.9900	0.9350	0.9933	0.9350	0.9650	0.9950
Avglength β	0.8524	0.4044	0.9067	0.7532	0.5382	1.4130
Power ¹	0.0340	0.4350	0.1013	0.1640	0.1020	0.5470
FDR 2	0.0026	0.0000	0.0056	0.0039	0.0000	0.0000
Avgcov 1	0.8450	0.7650	1.0000	0.9700	0.7950	1.0000
Avglength 1	22.9085	1.1822	2.7360	12.8441	1.1567	0.5102
Avgcov 2	0.9550	0.7200	0.9933	0.9400	0.8400	1.0000
Avglength 2	1.6373	5.2736	2.2481	10.5694	8.0125	0.5399
Avgcov 3	0.9150	0.8900	0.9933	0.8850	0.8100	1.0000
Avglength 3	1.3111	0.6948	2.0804	0.7724	0.8774	0.5530
Avgcov 4	0.9500	0.8900	0.9933	0.9550	0.8600	1.0000
Avglength 4	1.1413	0.3973	0.6932	1.0133	1.9001	0.5196
Avgcov 5	0.9250	0.9800	0.9800	0.9550	0.8500	0.9300
Avglength 5	1.6509	0.6189	0.5653	0.8354	7.2078	2.4434

Table 10: Simulation: small world

This table summarizes the results simulated on small-world type random graphs. Given the number of node N = 50,200,500, the mean degree for each node is 0.04N and 0.08N. The rewriting probability is fixed at 0.4.

The reported coverage is for 200 simulations from 50, 200 and 500 nodes graphs. The active set S_0 contains 5 nodes and coverage for each is reported as Avgcov 1-5. Nonactive set S_0^c contains all remaining nodes.

1. Power represents the averaged percentage in the active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

2. FDR reports the averaged percentage in the non-active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

		0.1			0.2	
Network Size	50	200	500^{5}	50	200	500^{5}
Avgcov S_0	0.9860	0.9940	0.9990	0.8950	0.9910	1.0000
Avglength S_0	10.2325	0.6168	1.4046	6.1020	0.7531	1.1056
Avgcov S_0^c	0.9923	0.9884	0.9868	0.9893	0.9909	0.9945
Avglength S_0^c	9.2108	2.2820	4.0344	5.9378	1.7050	1.9574
Avgcov β	0.9900	0.9650	1.0000	0.9750	0.9650	0.9900
Avglength β	8.2338	0.5556	1.1923	6.4205	0.6026	1.0265
0 0 /						
Power ¹	0.3480	0.6810	0.2300	0.5480	0.5620	0.1340
FDR 2	0.0037	0.0003	0.0001	0.0076	0.0026	0.0019
Network 1:						
probability	0.8050	0.8950	0.3700	0.7400	0.8500	0.2300
# identified 4	2.3540	3.8547	3.2973	4.0743	3.3314	2.9778
Network 2						
probability ³	0.0450	0.0550	0.0300	0.1300	0.0350	0.0100
# identifed 4	4.3333	1.0000	2.1667	3.4615	1.0000	1.0000
Avgcov 1	0.9750	0.9950	1.0000	0.8850	0.9950	1.0000
Avglength 1	17.4022	0.7528	0.9309	17.1760	0.9289	1.1231
Avgcov 2	0.9950	1.0000	1.0000	0.9450	0.9950	1.0000
Avglength 2	6.7053	0.4484	1.6877	1.7011	0.5777	1.3113
Avgcov 3	0.9900	0.9800	1.0000	0.9750	0.9850	1.0000
Avglength 3	15.5009	0.7919	1.3359	7.9879	1.0755	0.9475
0 0 0						
Avgcov 4	0.9800	0.9950	1.0000	0.8250	0.9850	1.0000
Avglength 4	9.6077	0.5766	1.3279	2.7826	0.7049	0.9092
Avgcov 5	0.9900	1.0000	0.9950	0.8450	0.9950	1.0000
Avglength 5	1.9465	0.4742	1.7408	0.8622	0.4783	1.2387

This table summarizes the results simulated on two Erdos-Renyi type random graphs. One of the network (Network 1) passes the endogenous effects while the other one (Network 2) is irrelevant to the decision.

The reported coverage is for 200 simulations from 50, 200 and 500 nodes graphs. The active set S_0 contains 5 nodes and coverage for each is reported as Avgcov 1-5. Nonactive set S_0^c contains all remaining nodes.

1. Power represents the averaged percentage in the active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

2. FDR reports the averaged percentage in the non-active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

3. Probability reports the empirical probability that at least one regressor in the group is significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

4. # identified reports the averaged number of significant regressors in the group conditioning on at least one regressor in the group is significant. False discover rate is controlled at 5% using Benjamini-Hochberg method.

5. For 500 cases, lasso tuning parameter is chosen using rule of thumb instead of cross-validation $% \left({{{\rm{T}}_{\rm{s}}}} \right)$

		0.1			0.2	
Network Size	50	200	500^{5}	50	200	500^{5}
Avgcov S_0	0.9670	0.9580	0.9850	0.9610	0.9954	0.9980
Avglength S_0	20.3014	1.3383	1.9988	8.6764	2.0044	4.5572
Avgcov S_0^c	0.9665	0.9883	0.9975	0.9680	0.9926	0.9980
Avglength S_0^c	14.0695	3.4002	4.7511	40.5927	1.6113	4.7505
Avgcov β	0.9800	0.9950	0.9900	0.9750	0.9943	0.9950
Avglength β	2.9138	0.8404	0.5866	1.5054	0.6253	0.6881
Avgcov γ	0.9600	0.9950	0.9950	0.9950	1.0000	1.0000
Avglength γ	0.5683	0.1568	0.0257	0.4235	0.0544	0.0294
test- $\gamma \neq 0$	0.4300	0.3750	1.0000	0.4950	1.0000	1.0000
Power ¹	0.0110	0.1890	0.4680	0.0140	0.7726	0.2550
FDR 2	0.0000	0.0035	0.0000	0.0000	0.0009	0.0000
Avgcov 1	0.9650	0.8100	0.9950	0.9700	1.0000	1.0000
Avglength 1	53.2561	3.1993	1.5312	9.2207	0.2910	5.4525
Avgcov 2	0.9350	0.9950	0.9900	0.9150	0.9886	0.9950
Avglength 2	34.8486	0.6602	0.8785	19.9796	8.4948	0.4050
Avgcov 3	1.0000	0.9950	1.0000	0.9500	1.0000	0.9950
Avglength 3	5.2305	0.9718	4.2919	3.9235	0.3581	3.8538
Avgcov 4	0.9800	0.9950	0.9950	0.9800	0.9943	1.0000
Avglength 4	4.0808	1.0082	2.8069	2.9452	0.5204	2.2894
Avgcov 5	0.9550	0.9950	0.9450	0.9900	0.9943	1.0000
Avglength 5	4.0909	0.8519	0.4855	7.3132	0.3577	10.7855

Table 12: Simulation

This table summarizes the results for Heterogeneous Endogenous Effects Model with Cliques.

The reported coverage is for 200 simulations from 50, 200 and 500 nodes graphs. The active set S_0 contains 5 nodes and coverage for each is reported as Avgcov 1-5. Nonactive set S_0^c contains all remaining nodes.

1. Power represents the averaged percentage in the active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

2. FDR reports the averaged percentage in the non-active set being significant after controlling False discover rate at 5% using Benjamini-Hochberg method.

5. For 500 cases, lasso tuning parameter is chosen using rule of thumb instead of cross-validation $% \left({{{\rm{T}}_{\rm{s}}}} \right)$

village	number of	number of	average	average	household	household	average rooms
	households	villagers	age	family size	having electric	having latrine	per person
Village1	182	843	32.7	4.6	90.7%	21.4%	0.6
Village2	195	877	31.4	4.5	94.4%	41.5%	0.5
Village3	294	1384	30.8	4.7	96.9%	47.3%	0.6
Village4	239	1026	31.3	4.3	98.3%	39.7%	0.5
Village12	175	802	30.7	4.6	90.9%	37.7%	0.6
Village19	204	1134	30.9	5.6	87.3%	14.7%	0.3
Village20	156	716	32.8	4.6	80.8%	25.0%	0.4
Village21	202	1046	28.6	5.2	83.7%	16.8%	0.4
Village23	254	1252	31.7	4.9	87.8%	28.3%	0.4
Village24	163	835	31.9	5.1	93.9%	13.5%	0.5
Village25	252	1313	30.9	5.2	96.4%	30.6%	0.6
Village28	315	1612	31.6	5.1	97.5%	34.0%	0.6
Village29	290	1337	32.2	4.6	84.1%	28.6%	0.6
Village31	153	851	26.1	5.6	97.4%	34.6%	0.5
Village32	241	1181	30.8	4.9	96.7%	20.7%	0.5
Village33	204	843	33.4	4.1	95.1%	6.4%	0.7
Village36	289	1214	33.4	4.2	84.8%	4.8%	0.7
Village39	289	1343	31.8	4.6	93.8%	42.2%	0.7
Village42	192	853	37.7	4.4	89.1%	28.6%	0.7
Village43	198	875	34.1	4.4	97.0%	26.8%	0.7
Village45	222	1076	29.8	4.8	94.6%	34.2%	0.5
Village47	139	687	33.7	4.9	94.2%	38.1%	0.6
Village50	244	999	34.8	4.1	92.2%	26.2%	0.7
Village51	251	1062	33.9	4.2	89.6%	13.1%	0.7
Village52	327	1525	33.8	4.7	91.7%	21.1%	0.7
Village55	257	1180	35.6	4.6	94.9%	4.7%	0.6
Village57	212	956	28.8	4.5	93.9%	3.8%	0.5
Village59	329	1599	31.4	4.9	96.0%	17.9%	0.6
Village62	190	994	32.1	5.2	92.1%	32.6%	0.5
Village65	299	1335	32.8	4.5	93.3%	29.4%	0.7
Village67	193	893	31.8	4.6	96.4%	25.4%	0.6
Village68	153	663	33.0	4.3	88.9%	22.2%	0.7
Village70	205	899	33.1	4.4	95.1%	24.9%	0.7
Village71	298	1388	28.8	4.7	95.0%	42.6%	0.6
Village72	223	999	32.0	4.5	96.9%	30.5%	0.7
Village73	174	870	30.1	5.0	96.6%	20.1%	0.6
Village75	172	831	32.7	4.8	91.3%	27.9%	0.7

Table 13: Descriptive Statistics

	(1)	(2)	(3)
Agriculture labour	-0.0141	0.0476^{*}	0.0672***
	(0.0136)	(0.0286)	(0.0134)
Anganavadi Teacher	0.0386	0.0664	0.1248^{**}
	(0.0602)	(0.1269)	(0.0593)
Bone Specialist	-0.2170	-0.3465	-0.0213
	(0.3314)	(0.6989)	(0.3265)
Blacksmith	-0.0752	-0.2279	0.1606^{*}
	(0.0927)	(0.1954)	(0.0913)
Construction/mud work	0.0050	0.2199^{***}	0.0562^{**}
	(0.0258)	(0.0544)	(0.0254)
Government Official	-0.0608	-0.0217	0.0350
	(0.0506)	(0.1067)	(0.0498)
Cook	0.0346	0.2015^{*}	-0.0168
	(0.0507)	(0.1068)	(0.0499)
Cow/livestock breeding	0.0059	-0.0235	0.0438
	(0.0282)	(0.0595)	(0.0278)
Truck/Tractor Driver	-0.0401	0.0746	0.0415
	(0.0305)	(0.0642)	(0.0300)
Factory worker (bricks/stones/mill)	0.0175	0.1756^{***}	0.0174
	(0.0246)	(0.0518)	(0.0242)
Milk dairy	0.0595	-0.1104	0.0622
	(0.0931)	(0.1965)	(0.0918)
Poultry farm	-0.1927	0.3577	0.0151
	(0.1258)	(0.2654)	(0.1240)
Small business	0.2006^{***}	0.1287^{***}	0.0606***
	(0.0227)	(0.0479)	(0.0224)
Silk/Cotton work	0.0031	0.0542	0.0266
	(0.0296)	(0.0624)	(0.0292)
Tailor Garment worker	0.0903***	0.1169^{*}	0.0309
	(0.0304)	(0.0642)	(0.0300)
Teacher	0.0268	-0.0452	0.0690
	(0.0426)	(0.0898)	(0.0420)

Table 14: Second Stage: who are they

	(1)	(2)	(3)
Daily labourer	-0.0172	0.1239**	0.0390
	(0.0283)	(0.0597)	(0.0279)
Auto driver	0.0113	0.2724**	0.0223
	(0.0548)	(0.1155)	(0.0540)
Police officer	-0.1459	-0.0374	0.3282^{*}
	(0.1917)	(0.4044)	(0.1890)
Waterman	-0.0722	0.0115	0.0715
	(0.0677)	(0.1428)	(0.0667)
Social Worker	-0.1541	-0.2959	-0.0475
	(0.1662)	(0.3505)	(0.1638)
Carpenter	-0.0863	-0.0816	0.0468
	(0.0748)	(0.1578)	(0.0737)
Electronics	0.0711	-0.1140	-0.0337
	(0.0727)	(0.1532)	(0.0716)
Goldsmith	-0.1351	0.2782	-0.0027
	(0.1664)	(0.3510)	(0.1640)
Hotel worker	0.3299***	0.4257^{***}	0.0759
	(0.0750)	(0.1581)	(0.0739)
Poojari	0.3697^{***}	-0.1542	0.1501
	(0.1369)	(0.2887)	(0.1349)
Post man	-0.1708	-0.3427	0.1632
	(0.1253)	(0.2643)	(0.1235)
Veterinary clinic	0.8649^{***}	1.9114^{***}	0.0377
	(0.3314)	(0.6990)	(0.3266)
Mechanic	0.0106	-0.1237	0.1274^{**}
	(0.0634)	(0.1337)	(0.0625)
Painter	-0.0832	0.1570	0.0034
	(0.0746)	(0.1574)	(0.0735)
Real Estate business	0.0158	0.6553^{***}	0.1088
	(0.1108)	(0.2337)	(0.1092)
Skilled labour/work for company	0.0469	0.0252	0.0809^{*}
	(0.0491)	(0.1036)	(0.0484)
Barber/saloon	0.4883^{***}	-0.0036	0.0443

Table 14 Continued: Second Stage: who are they

	(1)	(2)	(3)
	(0.1005)	(0.2119)	(0.0990)
Lawyer	-0.1235	0.0104	-0.1291
	(0.1915)	(0.4039)	(0.1887)
Security guard	-0.0993	0.0081	0.0016
	(0.1352)	(0.2852)	(0.1332)
Librarian	-0.0451	1.7625^{**}	-0.0848
	(0.3301)	(0.6962)	(0.3253)
Student	-0.2236	0.6929	0.1796
	(0.2341)	(0.4938)	(0.2307)
Doctor/Health assistant	0.2691^{**}	0.2703	0.0874
	(0.1053)	(0.2222)	(0.1038)
Fireman	0.0000	0.0000	0.0000
	(0.0000)	(0.0000)	(0.0000)
Photographer	-0.0995	-0.2046	0.2804
	(0.2336)	(0.4926)	(0.2302)
Folk artist	0.3541	-0.4611	0.0144
	(0.2379)	(0.5017)	(0.2344)
Begger	0.0000	0.0000	0.0000
	(0.0000)	(0.0000)	(0.0000)
Wood cutter	-0.0223	0.1942	-0.0000
	(0.0600)	(0.1265)	(0.0591)
Musician/Artist	0.3268	0.0640	-0.0436
	(0.2338)	(0.4931)	(0.2304)
Animal skin business	-0.1053	-0.0327	-0.0294
	(0.2350)	(0.4956)	(0.2316)
Average Age	0.0003	-0.0052***	-0.0003
	(0.0006)	(0.0013)	(0.0006)
Electric	0.0234	-0.0204	0.0309
	(0.0229)	(0.0482)	(0.0225)
Latrine	0.0533***	-0.0882***	0.0148
	(0.0134)	(0.0283)	(0.0132)
# Rooms	0.0315^{***}	-0.0087	0.0132***
	(0.0044)	(0.0094)	(0.0044)

Table 14 Continued: Second Stage: who are they

	(1)	(2)	(3)
Control village fix effect	Y	Y	Y

Table 14 Continued: Second Stage: who are they

Standard errors in parentheses * p < 0.1, ** p < 0.05, *** p < 0.01

design (1) uses whether one is predefined leaders as response variable

design (2) uses whether one joins the micro-finance program as response variable

design (3) uses whether one is selected by lasso as response variable

References

- Acemoglu, D., García-Jimeno, C., and Robinson, J. A. (2012). Finding eldorado: Slavery and long-run development in colombia. NBER WORKING PAPER SERIES.
- Ammermuller, A. and Pischke, J.-S. (2009). Peer effects in european primary schools: Evidence from pirls. Journal of Labor Economics, 27(3):315–348.
- Anselin, L. (1988). Spatial Econometrics: Methods and Models. Boston: Kluwer.
- Ballester, C., Calvó-Armengol, A., and Zenou, Y. (2006). Who's who in networks. wanted: The key player. *Econo-metrica*, (74):1403–1417.
- Bandiera, O., Barankay, I., and Rasul, I. (2009). Social connections and incentives in the workplace: Evidence from personnel data. *Econometrica*, 77(4):1047–1094.
- Banerjee, A., Chandrasekhar, A., Duflo, E., and Jackson, M. (2013). The diffusion of microfinance. Science, 341(6144).
- Belloni, A., Chernozhukov, V., and Hansen, C. (2014a). Inference on treatment effects after selection amongst high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650.
- Belloni, A., Chernozhukov, V., and Kato, K. (2014b). Uniform post selection inference for lad regression and other z-estimation problems.
- Belloni, A., Chernozhukov, V., and Kato, K. (2015). Uniform post selection inference for lad regression and other z-estimation problems.
- Belloni, A., Chernozhukov, V., and Wang, L. (2011). Square-root lasso: Pivotal recovery of sparse signals via conic programming. *Biometrika*, pages 1–18.
- Blume, L. E., Brock, W. A., Durlauf, S. N., and Jayaraman, R. (2015). Linear social interactions models. Journal of Political Economy, 123(2):444–496.
- Bonaldi, P., Hortacsu, A., and Kastl, J. (2015). An empirical analysis of funding costs spillovers in the euro-zone with application to systemic risk.
- Bramoullé, Y., Djebbari, H., and Fortin, B. (2009). Identification of peer effects through social networks. Journal of Econometrics, 150(1):41–55.
- Bühlmann, P. (2013). Statistical significance in high-dimensional linear models. Bernoulli, 41(2):802–837.
- Bühlmann, P. and van de Geer, S. (2011). Statistics for High-Dimensional Data. Springer.
- Calvó-Armengol, A., Patacchini, E., and Zenou, Y. (2009). Peer effects and social networks in education. Review of Economic Studies, (76):1239–1267.
- Chernozhukov, V., Hansen, C., and Spindler, M. (2015). Valid post-selection and post-regularization inference: An elementary, general approach. *Annual Review of Economics*, 7:649–688.
- Clark, A. E. and Loheac, Y. (2007). "it wasn't me, it was them!" social influence in risky behavior by adolescents. Journal of Health Economics, 26:763–784.
- Cliff, A. and Ord, J. K. (1973). Spatial autocorrelation. London: Pion.

- Coelli, T., Rahman, S., and Thirtle, C. (2002). Technical, allocative, cost and scale efficiencies in bangladesh rice cultivation: A nonparametric approach. *Journal of Agricultural Economics*, 53(3):607–626.
- Conley, T. G. and Udry, C. R. (2010). Learning about a new technology: Pineapple in ghana. AMERICAN ECONOMIC REVIEW, 100(1):35–69.
- Cressie, N. A. C. (1993). Statistics for Spatial Data. John Wiley & Sons, Inc.
- de Paula, A., Rasul, I., and Souza, P. C. (2015). Estimating and identifying social interactions.
- Denbee, E., Julliard, C., Li, Y., and Yuan, K. (2015). Network risk and key players: A structural analysis of interbank liquidity.
- Fan, J. and Liao, Y. (2014). Endogeneity in high dimensions. The Annals of Statistics, 42(3):872–917.
- Gautier, E. and TsyBakov, A. B. (2014). High-dimensional instrumental variables regression and confidence sets.
- Guryan, J., Kroft, K., and Notowidigdo, M. J. (2009). Peer effects in the workplace: Evidence from random groupings in professional golf tournaments. *American Economic Journal: Applied Economics*, 1(4):34–68.
- Horracea, W. C., Liu, X., and Patacchini, E. (2016). Endogenous network production functions with selectivity. Journal of Econometrics, 190(2):222–232.
- Jin, F. and Lee, L.-F. (2016). Lasso maximum likelihood estimation of parametric models with singular information matrices.
- Kasy, M. (2015). Uniformity and the delta method. arXiv preprint arXiv:1507.05731.
- Kelejian, H. H. and Prucha, I. R. (1995). A generalized moments estimator for the autoregressive parameter in a spatial model. *INTERNATIONAL ECONOMIC REVIEW*, 40.
- Kelejian, H. H. and Prucha, I. R. (1998). A generalized spatial two-stage least squares procedure for estimating a spatial autoregressive model with autoregressive disturbances. *Journal of Real Estate Finance and Economics*, 17(1):99–121.
- Krauth, B. V. (2005). Peer effects and selection effects on smoking among canadian youth. Canadian Journal of Economics, 38(3):735–757.
- Lee, L. (2002). Consistency and efficiency of least squares estimation for mixed regressive, spatial. *Econometric Theory*, 18(2):252–277.
- Lee, L. (2003). Best spatial two-stage least squares estimators for a spatial autoregressive model with autoregressive. Econometric Reviews, 22(4):305–335.
- Lee, L. (2004). Asymptotic distributions of quasi-maximum likelihood estimators for spatial econometric models. Econometrica, 72:1899–1926.
- Lee, L. and Liu, X. (2010). Efficient gmm estimation of high order spatial autoregressive models with autoregressive disturbances. *Econometric Theory*, 26:187–230.
- Lee, L.-f. and Yu, J. (2010). A spatial dynamic panel data model with both time and individual effects. *Econometric Theory*, 26:564–597.

- Leeb, H. and Potscher, B. M. (2005). Model selection and inference: facts and fiction. *Econometric Theory*, 21(1):21–59.
- Leeb, H. and Potscher, B. M. (2008). Can one estimate the unconditional distribution of post-modelselection estimators? *Econometric Theory*, 24(2):38–376.

Leeb, H. and Potscher, B. M. (2009). Model selection. Handbook of Financial Time Series, pages 889–925.

- Luo, Y. and Chernozhukov, V. (2016). Selecting informative moments via lasso.
- Manresa, E. (2013). Estimating the structure of social interactions using panel data.
- Manski, C. (1993). Identification of endogenous social effects: The reflection problem. The Review of Economic Studies, 60(3):531–542.
- Mas, A. and Moretti, E. (2009). Peers at work. American Economic Review, 99(1):112-145.
- Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *The* Annals of Statistics, 34(1436-1462).
- Nakajima, R. (2007). Measuring peer effects on youth smoking behaviour. *The Review of Economic Studies*, 74(3):897–935.
- Neidell, M. and Waldfogel, J. (2010). Cognitive and noncognitive peer effects in early education. Review of Economics and Statistics, 92(3):562–576.
- Sacerdote, B. (2001). Peer effects with random assignment: Results for dartmouth roommates. *The Quarterly Journal of Economics*, 116(2):681–704.
- Simon, N., Friedman, J., Hastie, T., and Tibshirani, R. (2013). The sparse group lasso. Journal of Computational and Graphical Statistics, 22(2):231–245.
- Upton, G. and Fingleton, B. (1985). Spatial data analysis by example. Volume 1: Point pattern and quantitative data. John Wiley and Sons Ltd.
- van de Geer, S., Buhlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202.
- Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables. Journal of the Royal Statistical Society, B(68):49–67.
- Zhang, C.-H. and Zhang, S. S. (2011). Confidence intervals for low-dimensional parameters in high-dimensional linear models. *Journal of the Royal Statistical Society*, 76(1):217–242.
- Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine Learning Research*, 7:2541–2563.
- Zhu, Y. (2016). Sparse linear models and l1regularized 2sls with high-dimensional endogenous regressors and instruments.