# ESTIMATING SOCIAL NETWORK MODELS WITH LINK MISCLASSIFICATION

ARTHUR LEWBEL

Department of Economics, Boston College


XI QU

Antai College of Economics and Management, Shanghai Jiao Tong University


XUN TANG

Department of Economics, Rice University

We propose an adjusted 2SLS estimator for social network models when the network links reported in samples are subject to two-sided misclassification errors (due, e.g., to recall errors by survey respondents, or lapses in data input). Misclassified links make all covariates endogenous, and add a new source of correlation between the structural errors and peer outcomes (in addition to simultaneity), thus invalidating conventional estimators used in the literature. We resolve these issues by adjusting endogenous peer outcomes with estimates of the misclassification rates and constructing new instruments that exploit properties of the noisy network measures. Simulation results confirm our adjusted 2SLS estimator corrects the bias from a naive, unadjusted 2SLS estimator which ignores misclassification and uses conventional instruments. We apply our method to study peer effects in household decisions to participate in a microfinance program in Indian villages.

KEYWORDS: Social Network, Link Misclassification.

# 1. INTRODUCTION

In many social and economic environments, an individual's behavior or outcome (such as a consumption choice or a test score) depends not only on his or her own characteristics, but also on the behavior and characteristics of other individuals. Call such dependence between two individuals a *link*. A *social network* consists of a group of individuals, some of whom are linked to others. The econometrics literature on social networks has largely focused on disentangling various channels of social effects based on observed outcomes and characteristics of network members. These include identifying the effects on each individual's outcome of (i) the individual's own characteristics (*individual effects*), (ii) the characteristics of people linked to the individual (*contextual effects*), and (iii) the outcomes of people linked to the individual (*peer effects*). See Blume et al. (2011) and Graham (2020) for extensive surveys about identifying such effects in social network models.

A popular approach for estimating social network models is to use two-stage least squares (2SLS). This requires researchers to construct instruments for the endogenous peer outcomes, using *perfect knowledge* of the network structure, as given by the *adjacency* matrix (i.e., the matrix that lists all links in the network). See, for example, Bramoullé et al. (2009), Kelejian and Prucha (1998), Lee (2007), and Lin (2010). In practice, samples of network links are often collected from survey responses. Such samples may be subject to an issue of misclassification in link status, due, e.g., to recall errors or misunderstandings by survey respondents, or lapses in data input. These misclassification errors can be *two-sided*: an existing link between two individuals may be misclassified as non-existent, or the sample may erroneously record links between those who are not linked.

Misclassification of links in the sample poses major methodological challenges for estimators like 2SLS. To see this, consider a data-generating process (DGP) from which a large number of independent networks (i.e., groups) are drawn. Each group consists of $n$ individual members.[1] Suppose that in each group, a vector of individual outcomes $y \in \mathbb{R}^n$ is determined by a structural model:

$$y = \lambda G y + X \beta + \varepsilon, \text{ where } E(\varepsilon | X, G) = 0.$$

---

[1] We later allow for a single growing network, but our results are easiest to illustrate in the context of many independent, identically sized groups.

In this model, the adjacency matrix, $G$, is an $n$-by-$n$ matrix of dummy variables that describes the group's network: the element in row $j$ and column $k$ of $G$ equals one if individual $j$ is linked to member $k$, and zero otherwise.[2] Here $X$ is an $n$-by-$K$ matrix of exogenous covariates, and $\varepsilon$ is an $n$-vector of structural errors. The random arrays $y$, $G$, $X$, and $\varepsilon$ all vary across the groups in the sample, while the coefficients $\lambda$ and $\beta$ are the same across groups. We drop group subscripts for clarity.

For simplicity we have for now omitted contextual effects, i.e., a term defined as $GX\gamma$, and any group-level fixed effects. Extensions of our results that deal with these features are provided later in Section 5.

The regressors in the model are $Gy$ and $X$. While $X$ is exogenous, the regressors $Gy$ are correlated with $\varepsilon$. The issue of simultaneity arises here, because any one individual's outcome depends on, and is determined simultaneously with, the outcomes of other group peers. A simple estimator of the peer effect $\lambda$ and individual effects $\beta$ that deals with this simultaneity problem is 2SLS, using $GX$ or $G^2X$ as instruments for $Gy$, as in Bramoullé et al. (2009).[3]

But now suppose that, in each group, a researcher does not observe $G$ perfectly, but instead observes a noisy measure $H$, which differs from $G$ by randomly misclassifying some links in the data-generating process. The goal now is to estimate $\lambda$ and $\beta$ from a "*feasible*" structural form like:

$$y = \lambda Hy + X\beta + u, \tag{1}$$

where $u \equiv [\varepsilon + \lambda(G - H)y]$ is a vector of *composite* errors.

The misclassified links in $H$ aggravate endogeneity issues in (1) in three important ways. First, they lead to correlation between $X$ and the error $u$ through $\lambda(G - H)y$, a component of $u$ that is due to the measurement error in the adjacency matrix. This component contains

---

[2]This is a "local-aggregate" network model, where the endogenous effect depends on the *aggregate* outcome of those linked to an individual. It differs from a "local-average" network model, where the endogenous effect is represented by the *average* outcome of those linked peers.

[3]If the model includes contextual effects $GX\gamma$ in its structural form, then $G^2X$ can be used as instruments for $Gy$; otherwise use of $GX$ as instruments suffices.

4

$y$ which by the model is correlated with $X$. This means that the regressors $X$ are no longer exogenous!

Second, these misclassified links cause an additional source of endogeneity in $Hy$. Like $Gy$, the feasible $Hy$ is correlated with the model error $\varepsilon$ due to simultaneity. But in addition, $Hy$ is also correlated with $u$ through the measurement error term $\lambda(G - H)y$.

Third, misclassification means that, unlike using $GX$ or $G^2X$ as instruments when $G$ is perfectly reported in the sample, 2SLS estimates based on the *feasible* instruments $HX$ or $H^2X$ would be inconsistent because $HX$ correlates with $\lambda(G - H)y$, resulting in a failure of instrument exogeneity.

For all these reasons, conventional 2SLS estimators of this model become inconsistent in the presence of misclassification errors in the links.[4]

In this paper, we introduce an *adjusted-2SLS* estimator, which resolves these endogeneity issues and consistently estimates $(\lambda, \beta)$ using alternative valid instruments constructed from $H$ despite the misclassification errors in the links. We first introduce the main idea for a benchmark case, where an observed $H$ differs from the true $G$ due to random, two-sided misclassification errors at unknown rates $p_0, p_1 \in (0, 1)$. Here, $p_1$ is the probability that any existing link is missing in the sample, while $p_0$ is the probability that a non-existent link is erroneously recorded as existing in the sample. Later, in Section 5, we extend our method to allow the misclassification rates $p_0(X), p_1(X)$ to depend on covariates.

Our method is based on a series of new insights that have not been explored in the literature. *First*, we observe that by adjusting the noisy measure of peer outcomes $Hy$ using the misclassification rates $(p_0, p_1)$, we restore the exogeneity of $X$ in an *adjusted* feasible structural form. Formally, this means if we replace (1) with:

$$y = \lambda W_{(H,p_0,p_1)}y + X\beta + v, \tag{2}$$

where $W_{(H,p_0,p_1)}$ is a properly designed adjustment of the network measure $H$, then the adjusted composite errors $v \equiv \varepsilon + \lambda[G - W_{(H,p_0,p_1)}]y$ in (2) satisfy $E(v|X, G) = 0$. This holds regardless of how the actual network $G$ is formed, as long as $E(\varepsilon|X, G) = 0$.

---

[4]While we focus on the 2SLS estimator in this paper, the same arguments apply to show that conventional maximum likelihood, and the generalized least squares estimators based on (1) are also inconsistent when there are misclassification errors in the links.

*Second*, despite the restored exogeneity of $X$ in (2), conventional instruments such as $HX$ or $H^2X$ remain invalid, because the adjusted errors $v$ still depend on $H$. To resolve this issue, we provide alternative functions of $H$ and $X$ that are valid instruments. For example, we show that if $H$ is an unsymmetrized measure of the actual adjacency matrix $G$, then $H'X$ is uncorrelated with $v$ (where $H'$ is the transpose of $H$), *despite* misclassification errors in $H$. This result holds regardless of whether $G$ is symmetric (i.e., with all links being undirected) or asymmetric (i.e., consisting of directed links). Therefore, we can use $H'X$ as valid instruments in an adjusted-2SLS where peer outcomes are adjusted by $W_{(H,p_0,p_1)}$. To the best of our knowledge, no other paper in the literature has proposed $H'X$ as instruments.

Another scenario, which works regardless of whether the observed or actual adjacency matrices are symmetric or not, is when we observe two noisy measures of the same actual $G$. An example is our empirical application, where we observe two different reports of who visits whom. This means we observe two different $H$ matrices with independent misclassification errors. We show 2SLS becomes valid if we use one of these matrices to construct the adjustment term $W_{(H,p_0,p_1)}$ and the other to construct instruments.

Our *third* contribution is to show that under either scenario above (i.e., when the sample reports either a single unsymmetrized noisy measure $H$, or two independent measures that may or may not be symmetrized), we can provide simple methods to identify and estimate the unknown misclassification rates $(p_0, p_1)$.[5]

Building on these insights, we construct adjusted 2SLS estimators for $(\lambda, \beta)$, and provide their limiting distribution as the number of groups in the sample grows to infinity. This estimator essentially applies 2SLS to the adjusted peer outcomes $W_{(H,p_0,p_1)}y$ in (2), using our new instruments and a closed-form, sample analog estimator for the misclassification rates $(p_0, p_1)$. The estimator is easy to implement, does not require any numerical searches, and Monte Carlo simulations demonstrate its good performance in finite samples.

---

[5]The approach we take in this step differs from, and is simpler than, other papers that use multiple measures to deal with misclassification in discrete explanatory variables (e.g. Mahajan (2006), Lewbel (2007), and Hu (2008)). This is because, for implementing our adjusted-2SLS, it is only necessary to estimate the rates $(p_0, p_1)$, rather than the distribution of outcomes conditional on the actual $G$.

We then generalize the model and our estimator in several directions. We show how to include contextual effects (a term defined as $GX\gamma$) as well as group-level fixed effects into the structural form in (2). We also allow the misclassification rates $(p_0, p_1)$ to be heterogeneous and depend on covariates in $X$. Furthermore, we extend our method to the case of a single large network where the sample can partitioned into *approximate* groups.

Finally, we apply our method to estimate peer effects in household decisions to participate in a microfinance program in Indian villages, using data from Banerjee et al. (2013). We match the individual survey to the household survey there, yielding a sample of 4,149 households from 43 villages in South India. The parameter of interest is the peer (endorsement) effect, which reflects how a household's decision is influenced by the microfinance program participation of other households to which it is linked. Survey information about visits between the households provides two symmetrized noisy measures of undirected links (i.e., two symmetrized $H$ measures). We estimate the misclassification rates in each of these two measures using our method, and apply these estimated rates in our adjusted-2SLS procedure to estimate the peer effects.

We find that participation by another linked household increases a household's own participation rate by around 5.1%. This effect is economically significant, compared to the average participation rate of 18.9% in the sample. We also find that ignoring the issue of link misclassification in the noisy measures and applying conventional 2SLS estimation results in an upward bias in the estimates of these peer effects (Monte Carlo simulations show that this bias can be large, though it turns out to be modest in our application).

**Roadmap.** Section 2 reviews the related literature, and explains our contribution in its context. Section 3 specifies the model, and illustrates the main ideas in a benchmark model with independent and identical misclassification rates. Section 4 defines a closed-form estimator for the misclassification rates, and provides an adjusted-2SLS estimator for the social effects. Section 5 extends the method to settings with contextual effects, heterogeneous misclassification rates, or group fixed effects, and to the setting of a single large network. Section 6 presents Monte Carlo simulation results. Section 7 applies our method to analyze peer effects in microfinance participation in India. Proofs are collected in the Online Appendix.

## 2. RELATED LITERATURE

Models with misclassified binary or discrete variables have been studied extensively in the econometrics literature. Aigner et al. (1973), Klepper (1988), Bollinger (1996), and Molinari (2008) point-identify or set-identify such models using various restrictions on the misclassification rates; Mahajan (2006), Lewbel (2007), and Hu (2008) exploit exogenous instruments to deal with misclassified explanatory variables.

Estimation of peer effects in social networks with measurement errors in the links is an increasingly important topic. Butts (2003) proposes a hierarchical Bayesian model to infer social structure in the presence of measurement errors. Shalizi and Rinaldo (2013) note the challenge of dealing with missing network links in Random Graph Models. Advani and Malde (2018) show that even a relatively low misreporting rate can lead to large bias in causal effect estimates. Chandrasekhar and Lewis (2011) show how egocentrically sampled network data can be used to predict the full network in a graphical reconstruction process. Liu (2013a) shows that when the adjacency matrix is not row-normalized, instrumental variable estimators based on an out-degree distribution can be valid.

Hardy et al. (2019) estimate treatment effects on a social network when the reported links are a noisy representation of true spillover pathways. They use a mixture model that accounts for missing links as unobserved network heterogeneity, and estimate it using an Expectation-Maximization algorithm. This approach requires a parametric model of how links are determined and treatment is assigned, and requires enumerating the likelihood conditional on all possible treatment exposures (which in turn depends on the latent unobserved network). Auerbach (2022) studies a network model where links are correctly measured but both peer and contextual effects interact with unobserved individual heterogeneity that affects link formation. In contrast with these papers, we focus on estimating social effects in linear social networks while fully exploiting implications of randomly misclassified links. Our method does not require modeling the formation of actual links; our estimator is an adjusted 2SLS, which has a closed form and is easy to compute.

Liu (2013b) estimates a social network model when the data consists of a *subset* of individuals sampled randomly from a larger group in the population. In his setting, the links and outcomes among this sampled subset of group members are perfectly measured

8

while those of all others are not reported in the data.[6] In comparison, we do not study the inference of sampled networks; instead, we let the group memberships be fixed and known, and allow every individual in the sample to have randomly misclassified links. As noted above, this imperfect measure of links leads to failure of conventional 2SLS in our setting.

Boucher and Houndetoungan (2020) estimate peer effects when the social networks in the sample are subject to measurement issues, such as missing or misclassified links. They also consider the case when the researcher has access to aggregated relational data only. Their method requires researchers know, or have a consistent estimator of, the distribution of the actual network. They construct instruments by drawing from this distribution, and use 2SLS to estimate the peer effects. In comparison, the method we propose does not require such prior knowledge or estimates of network distribution.

Griffith (2022) studies the case where links are censored in the sample, and character-izes the bias in a reduced-form regression (i.e., when the outcomes in $y$ are regressed on exogenous covariates $X$ and $GX$). For a model with $\lambda = 0$, Griffith (2022) shows the bias can be consistently estimated under an order invariance condition, i.e., the covariance of characteristics of those linked to an individual is invariant to the order in which those links are reported or censored.[7] Griffith and Kim (2023) extend this investigation to in-clude both linear-in-sums (where $G$ has binary entries) and linear-in-means (where $G$ is row-normalized). They show how nonzero, structural peer effects $\lambda$ enter the estimand of the reduced-form regression above, as well as how general misclassification, e.g., due to randomly missing links or censored links, affect these estimands. In comparison, we focus on empirical settings where links are misclassified at random. (This is later generalized to the case with heterogeneous misclassification rates.) We show that conventional 2SLS estimands in this case contain bias in peer effects (e.g., an augmentation bias when misclas-sification is one-sided with $p_0 = 0$ and $p_1 > 0$), and no bias in other individual effects. Bias

------

[6]In our notation, this means some rows in $G$, as well as their corresponding rows in $Y$ and $X$, are not included in the data due to random sampling.

[7]This condition mitigates the issue of endogenous selection of uncensored links, and in this sense plays a similar role to our assumption of randomly misclassified links.

correction in our case is immediate once the misclassification rates are estimated using a simple approach that we provide.

Lewbel et al. (2023) show that if the order of measurement errors in links is sufficiently small (e.g., the number of misclassified links in a single, large network does not grow too fast with the sample size), conventional instrumental variables estimators that ignore these measurement errors remain consistent, and standard asymptotic inference methods remain valid. They also provide specific examples in which the link formation or misclassification rates decrease with the sample size to imply such a small order of measurement errors. In contrast, in this paper we deal with new challenges outside the scope of Lewbel et al. (2023). Namely, we allow the misclassification rates to non-diminishing (fixed) in an asymptotic framework with many independent, finite-sized groups. In such settings, the measurement errors are large enough to invalidate conventional 2SLS estimators.

## 3. MODEL AND IDENTIFICATION

Consider a DGP from which a large number of small, independent networks (groups) are drawn, such as villages or classrooms.[8] We will first identify and estimate a linear social network model when links are randomly misclassified in the sample (we later allow misclassification probabilities to depend on covariates). We establish the asymptotic properties of our estimator as the number of groups in the sample approaches infinity.

The structural form for the vector of individual outcomes $y_s \in \mathbb{R}^{n_s}$ in group $s$ is:

$$y_s = \lambda G_s y + X_s \beta + \varepsilon_s, \tag{3}$$

where the peer effect $\lambda$ and the direct effects $\beta$ are constant parameters of interest, $X_s$ is an $n_s$-by-$K$ matrix of individual- or group-level explanatory variables, and $G_s \in \{0,1\}^{n_s \times n_s}$ is the network (adjacency) matrix for group $s$, with its $(i,j)$-th entry $G_{s,ij} = 1$ if an individual member $i$ is linked to another member $j$ in group $s$, and $G_{s,ij} = 0$ otherwise. The matrix $G_s$ may be asymmetric with directed links ($G_{s,ij} \neq G_{s,ji}$ for some $i \neq j$), or symmetric with undirected links ($G_{s,ij} = G_{s,ji}$ for all $i \neq j$ almost surely).

---

[8]Later in Section 5.4 and the Online Appendix we consider the extension to a single growing network, which includes links both between and within groups. Each group $s$ consists of $n_s \geq 3$ individual members.

Let $I_s$ by an $n_s$-by-$n_s$ identity matrix, and assume $(I_s - \lambda G_s)$ is invertible almost surely. A sufficient condition for this is that $||\lambda G_s|| < 1$ for *any* matrix norm $|| \cdot ||$ almost surely. Solving equation (3) for $y_s$ gives the reduced form for outcomes:

$$y_s = M_s(X_s\beta + \varepsilon_s), \text{ where } M_s \equiv (I_s - \lambda G_s)^{-1}. \tag{4}$$

We do not observe the $G_s$ matrices. Instead, for each group $s$, the sample reports a noisy measure $H_s \in \{0,1\}^{n_s \times n_s}$ of the actual adjacency matrix $G_s$. That is, for some unknown pairs of individuals $i \neq j$, $G_{s,ij}$ is randomly misclassified as $H_{s,ij} = 1 - G_{s,ij}$. By convention, let $G_{s,ii} = 0$ and $H_{s,ii} = 0$ for all $i$ and $s$.

To simplify exposition, we let the group sizes $n_s = n$ be fixed across groups $s = 1, 2, ..., S$ for now. This allows us to drop the group subscript $s$ while presenting our identification argument. We will later add back these group subscripts and allow for variation in group sizes when we define our estimator in Section 4.

### 3.1. *Assumptions*

We maintain the following conditions on the noisy measure $H$ throughout Section 3:

(A1) $E(H_{ij}|G, X) = E(H_{ij}|G_{ij}, X)$ for all $i$ and $j$;

(A2) $E(H_{ij}|G_{ij} = 1, X) = 1 - p_1$, $E(H_{ij}|G_{ij} = 0, X) = p_0$, and $p_0 + p_1 < 1$ for all $i \neq j$;

(A3) $E(\varepsilon|G, X, H) = 0$.

Condition (A1) states the incidence of misclassifying a link between individuals $i$ and $j$ is conditionally independent from the actual status of all other links. Under (A2), misclassification probabilities conditional on actual link status are fixed at $p_0$ and $p_1$ respectively, and are independent from $X$ (we will later allow these probabilities to depend on $X$). With $\Pr\{G_{ij} = 1\} < 1$, the inequality constraint "$p_0 + p_1 < 1$" is equivalent to "$H_{ij}$ and $G_{ij}$ are positively correlated." That is, the noisy measure is positively correlated with the actual link status despite the misclassification error. This is a standard condition in the literature on misclassified regressors, e.g., Bollinger (1996), Hausman et al. (1998), and ensures the relevance of the instrumental variable constructed by $H$. Condition (A3) rules out endogeneity in link formation, assuming $(G, X, H)$ are exogenous to structural errors $\varepsilon$.

Conditions (A1) and (A2) hold jointly in two common scenarios. In the first scenario, which we refer to as *unsymmetrized* measures, each $(i,j)$-th entry in $H$ is an *independent* measure of $G_{ij}$. For example, $H_{ij}$ (or $H_{ji}$) reports individual $i$'s (or $j$'s) binary response to a survey question about whether a link exists between $i$ and $j$. A measure $H$ constructed this way is flexible in that it allows the researcher to remain agnostic about whether the actual $G$ is symmetric with undirected links or not. This is also an intuitive way to construct $H$ when the actual $G$ is *known* to be asymmetric with directed links. In this scenario, if misclassification of $G_{ij}$ happens independently at rates $p_0$ or $p_1$ across links (depending on whether $G_{ij} = 1$ or $0$), then (A1) and (A2) are satisfied. To reiterate, (A1) and (A2) hold in this first scenario, regardless of whether the actual $G$ is symmetric or not.

In the second scenario, which we refer to as *symmetrized* measures, the actual $G$ is *known* to be symmetric with undirected links, and hence the researcher chooses to symmetrize $H$ by combining independent measures of entries in $G$. For example, the researcher asks $i$ and $j$ whether they have an undirected link, and records their responses respectively. The researcher then constructs a symmetrized measure by setting $H_{ij}$ and $H_{ji}$ both to 1 if *either $i$ or $j$* responds positively, and both to $0$ otherwise. Suppose the responses from $i$ or $j$ independently misclassify an *existing* link at rate $\varphi_1 > 0$ (say, due to idiosyncratic recall errors). Then $\Pr\{H_{ij} = 0 | G_{ij} = 1\} \equiv p_1 = \varphi_1^2$. Likewise, if $i$ and $j$ independently misclassify a *non-existent* link at rate $\varphi_0$, then $\Pr\{H_{ij} = 1 | G_{ij} = 0\} \equiv p_0 = 1 - (1 - \varphi_0)^2$. Thus, in this second scenario, (A1) and (A2) hold with $\Pr\{H_{ij} = H_{ji}\} = 1$ and with the two entries sharing the same misclassification rates $p_1$ and $p_0$ specified above.

On the other hand, (A1) *does* rule out a third, empirically less plausible scenario, in which the actual $G$ is asymmetric with directed links but researchers mistakenly impose a symmetrized $H$ using independent measures of $G_{ij}$ and $G_{ji}$ as in the second scenario. In this case, the equality in (A1) fails in general because $E(H_{ij} | G_{ij} = 1, G_{ji} = 1) = 1 - \varphi_1^2$ while $E(H_{ij} | G_{ij} = 1, G_{ji} = 0) = \varphi_0 + (1 - \varphi_1) - \varphi_0(1 - \varphi_1)$.

A clear advantage of the method we propose is that it allows researchers to consistently estimate social effects while being agnostic about whether the actual links in $G$ are directed or not. Our method only requires the noisy measure $H$ satisfy (A1)-(A3), which, as explained above, is not confined to the (a)symmetry of $G$ or $H$. We therefore recommend a simple guideline for practitioners collecting link data: if a researcher is unsure about whether the actual links in $G$ are directed or undirected, then a safe approach is to con-

struct an unsymmetrized measure $H$ as in the first scenario, and apply our method in this paper to deal with possible misclassification of the links.

It is important to note that (A1)-(A3) do *not* specify how the actual links in $G$ are formed. These conditions do not impose any known information about the actual adjacency matrix, except for its exogeneity in (A3). Nor do they impose any structure that can be used to derive a conditional likelihood for the actual network, which is $\Pr\{G|H,X\} = \frac{\Pr(H|G,X)\Pr(G|X)}{\sum_{G'}\Pr(H|G',X)\Pr(G'|X)}$. Constructing such a likelihood would require specifying the likelihood of the actual network $\Pr\{G|X\}$, which we refrain from doing in this paper. Our method therefore differs qualitatively from alternative methods which either use graphical reconstructions such as Chandrasekhar and Lewis (2011), or require knowledge of the distribution of actual adjacency matrix such as Boucher and Houndetoungan (2020).

Define an *infeasible, adjusted* measure of the adjacency matrix:

$$W \equiv W_{(H,p_0,p_1)} \equiv \frac{H - p_0(\iota\iota' - I)}{1 - p_0 - p_1},$$

where $\iota$ is a vector of ones and $(\iota\iota' - I)$ is a square matrix with all off-diagonal entries being 1 and all diagonal entries being 0. For the rest of this paper, we suppress subscripts indicating the arguments $(H,p_0,p_1)$ in $W$ to simplify notation. Then, $W_{ij} = (H_{ij} - p_0)/(1 - p_0 - p_1)$ for $i \neq j$, and $W_{ii} = H_{ii} = 0$. Under (A1) and (A2), $E(W_{ij}|G,X) = 1$ whenever $G_{ij} = 1$, and $E(W_{ij}|G,X) = 0$ whenever $G_{ij} = 0$ (including the case with $i = j$). Thus,

$$E(W|G,X) = G. \tag{5}$$

In the next subsection, we exploit this property in (5) to establish a useful intermediate result: despite link misclassification, $(\lambda,\beta)$ could be consistently estimated by an adjusted 2SLS if the misclassification rates $p_0, p_1$ were known.

## 3.2. *Infeasible two-stage least squares*

We write a new *adjusted* structural form using $W$:

$$y = \lambda W y + X\beta + \underbrace{\varepsilon + \lambda(G - W)y}_{\equiv v}. \tag{6}$$

This form is infeasible because $W$ is a function of the unknown misclassification rates $p_0$ and $p_1$. Lemma 1 shows $X$ is uncorrelated with its composite errors $v$, despite link misclassification.

LEMMA 1: *Under (A1), (A2), and (A3), $E(v|G, X) = 0$.*

This lemma is fundamental for our method; it restores exogeneity of $X$ by adjusting the structural form properly to account for link misclassification. Such exogeneity then allows us to construct instruments that depend on $X$.

The importance of Lemma 1 is best illustrated in contrast with the *naive* structural form in (1), i.e., $y = \lambda Hy + X\beta + u$, which ignores misclassification errors and simply uses $Hy$ as peer outcomes on the right-hand side. The composition errors in (1) are:

$$u = \varepsilon + \lambda(G - H)y = v + \lambda(W - H)y$$

$$= v + \left(\frac{p_0 + p_1}{1 - p_0 - p_1}\right)\lambda Hy - \left(\frac{p_0}{1 - p_0 - p_1}\right)\lambda(\iota\iota' - I)y. \tag{7}$$

While $E(v|G, X) = 0$ by Lemma 1, the second and third terms on the right-hand side of (7) do not satisfy such mean independence. Therefore, in a simple, feasible structural form that uses $Hy$ instead of $Wy$, the covariates in $X$ are generally endogenous due to the ignored misclassification errors. Later we show such endogeneity leads to an "*augmentation bias*" in the 2SLS estimation of (1) when misclassification is one-sided ($p_0 = 0$). To reiterate, Lemma 1 shows that the adjustment in $W$ is crucial for restoring exogeneity of $X$ in (6).

Lemma 1 may seem surprising ex ante, because one would expect $(G, X)$ to be correlated with the composite error $v$ which depends on $y$. The intuition for the exogeneity in this lemma is as follows. Once we condition on the actual adjacency $G$ and $X$, randomness in individual outcomes $y$ is solely due to the actual structural errors $\varepsilon$, which are uncorrelated with both $X$ and $(H, G)$ under (A3). As a result, any potential correlation between $v$ and $(G, X)$ could only be due to the measurement error $\lambda(G - W)y$. But the property established in (5) and the exogeneity of $\varepsilon$ in (A3) imply this measurement error is mean-independent from $(G, X)$. A formal proof of Lemma 1 is in the Online Appendix.

Note that we can not use the exogeneity established in Lemma 1 alone to construct GMM estimators for $(\lambda, p_0, p_1)$, because it does not suffice for the joint identification of these parameters. This can be easily seen in the special case where the misclassification

14

is one-sided ($p_0 = 0$). In that case, the conditional moment due to Lemma 1 simplifies to $E(y - \frac{\lambda}{1-p_1}Hy - X\beta|G, X) = 0$, which is not sufficient for recovering $\lambda$ and $p_1$ separately *even if* $G$ were perfectly observed in the data-generating process.

Our goal for the rest of Section 3 is to combine the exogeneity attained in Lemma 1 with further information, such as instruments and multiple measures $H$, to identify all model parameters, including the misclassification rates. First off, note the term $Wy$ in (6) remains endogenous, *even if* the misclassification rates were known and used to construct the adjusted measure $W$. This is because $E[(Wy)' v] \neq 0$ in general.[9]

We next consider 2SLS estimation of equation (6). Let $R \equiv (Wy, X)$. Suppose that we had a set of instruments $Z$ for $R$, i.e., instruments that we could use to estimate equation (6). By Lemma 1, $Z$ can include $X$, so we only need an additional instrument for $Wy$. We will later provide some possible instruments for $Wy$. But for now, just consider what properties any such matrix of instruments $Z$ must satisfy. $Z$ must be an $n$-by-$L$ matrix with $L \geq K + 1$ such that $E(Z'v) = 0$ and the following rank condition holds:

(IV-R) $E(Z'R)$ and $E(Z'Z)$ have full column rank.

Let $\Pi \equiv [E(Z'Z)]^{-1} E(Z'R)$. By (6) and Lemma 1,

$$\Pi' E(Z'y) = \Pi' E(Z'R)(\lambda, \beta')' + \Pi' E(Z'v)$$

$$\Rightarrow (\lambda, \beta')' = [\Pi' E(Z'R)]^{-1} [\Pi' E(Z'y)]. \tag{8}$$

PROPOSITION 1: *Suppose (A1), (A2), and (A3) hold, and that (IV-R) holds for instruments $Z$. The two-stage least-squares estimand using $Z$ for (6) is $(\lambda, \beta')'$.*

Using $Wy$ instead of $Hy$ as the first regressor in $R$ is crucial for consistency in Proposition 1. To see why, suppose one applies 2SLS to (1) using $Hy$, so the regressors are $\check{R} \equiv (Hy, X)$ and the resulting model errors are $u$ as defined in (7). Then the 2SLS estimand would be $(\lambda, \beta')' + [\check{\Pi}' E(Z'\check{R})]^{-1}[\check{\Pi} E(Z'u)]$, where $\check{\Pi}$ is similar to $\Pi$ only with

---

[9]Under (A1) and (A2), $E(W'G|G, X) = G'G$, but $E(W'W|G, X) \neq G'G$ in general. This is because the $i$-th diagonal entry in $W'W$ is $\sum_k W_{ki}^2$ while its $(i, j)$-th off-diagonal entry is $\sum_k W_{ki}W_{kj}$. It then follows from (A3) and the law of iterated expectation that $E(y'W'Wy) \neq E(y'W'Gy)$ in general.

$R$ replaced by $\check{R}$. Endogeneity bias arises as $E(Z'u) \neq 0$ in general. This is because $u$ is correlated with the latter two terms on the right-hand side of (7) through $y$.

In the special case with one-sided misclassification (i.e., $p_0 = 0$ and $p_1 > 0$ so that actual links are missing at random, but the sample never reports links that do not exist), such endogeneity bias takes a simple analytical form. By plugging in the expression of $u$ from (7) and setting $p_0 = 0$, we can show $E(Z'u) = \left(\frac{p_1}{1-p_1}\right) E[Z'\check{R}(\lambda, \mathbf{0})']$, where $\mathbf{0}$ is a row-vector of $K$ zeros. Consequently, the 2SLS estimand in this case is $(\frac{\lambda}{1-p_1}, \beta')'$, indicating an "*augmentation*" bias in the peer effect estimator.

Based on Proposition 1, we have two main requirements for estimating the model. First, we need to construct a valid instrument for $Wy$. One possibility, based on Lemma 1, is that nonlinear functions of regressors $X$ could serve as instruments, provided they satisfy the rank condition in (IV-R), possibly by correlating with the unknown process by which links are formed in $G$. However, nonlinear functions of $X$ may be weak instruments since the structural model is linear in $X$. So instead in Section 3.3 we show how to construct valid instruments using $H$ and $X$.

The second remaining requirement for estimating the model is that we need to identify and estimate the unknown misclassification rates $p_0$ and $p_1$ in order to construct the adjusted measure $W$. We address this question in Section 3.4.

### 3.3. *Constructing instruments from network measures*

We now focus on the question about how to construct instruments using noisy network measures in the sample and subject to misclassification errors. We propose two options for constructing IVs, depending on the number of measures available and whether the measures are symmetrized.

3.3.1. *Instruments using a single unsymmetrized measure*

First, consider a setting in which the sample reports a single, unsymmetrized network measure $H$. Assume:

(A4) Conditional on $(G, X)$, $H_{ij}$ and $H_{kl}$ are independent whenever $i \neq k$ or $j \neq l$.

This condition states that different links are misclassified independently conditional on the actual link status. It is important to note that this condition does *not* restrict whether the

16

actual network $G$ is symmetric or not. For example, $H$ may be an *unsymmetrized* measure of $G$ as defined in the first scenario under (A1)-(A2) in Section 3.1). In this case, (A4) holds when $H_{ij}$ and $H_{ji}$ are independent measures of $G_{ij}$ and $G_{ji}$ respectively, *regardless of* whether $G_{ij} = G_{ji}$ in the actual $G$.

On the other hand, (A4) fails when $H$ is a *symmetrized* measure, because in this case $H_{ij}$ and $H_{ji}$ are identical by construction and hence cannot be independent. To deal with this case of symmetrized measures, we give an alternative method for constructing instruments in Section 3.3.2.

We propose to construct instruments using $H$ and $X$ in the following proposition.

PROPOSITION 2: *Suppose (A1), (A2), (A3), and (A4) hold. Then $E(Z'v) = 0$ for $Z \equiv (H'X, X)$ or $Z \equiv (W'X, X)$.*

Proposition 2 suggests using $H'X$ or $W'X$ as instruments for $Wy$. There is a simple interpretation of these instruments: the $i$-th component (row) of $H'X$ is the sum of characteristics of all individuals who report links with $i$ in the sample.

Recall that $GX$ would be valid instruments for $Gy$ if $G$ were perfectly observed in the sample. Therefore, one may wonder why we use $H'X$ instead of $HX$ as instruments here. The reason is that $H'X$ are valid instruments while $HX$ are not. To give some intuition why, observe that the composite error $v$ in (6) contains $\lambda(G - W)$ and so includes $H$ through $W$. The covariance of this error with $HX$ contains the conditional variance of $H$, which can't be zero. Therefore the error $v$ is correlated with $HX$. In contrast, the corresponding terms in the covariance of $v$ with $H'X$ are conditional covariances of $H_{ij}$ with $H_{ji}$, which by (A4) are zero. And condition $p_0 + p_1 < 1$ in (A2) ensures the relevance of instrument. Hence $H'X$ satisfies instrument exogeneity while $HX$ does not. The same logic holds for using $W'X$ but not $WX$ as instruments.

In addition to validity, the set of instruments $Z$ needs to also satisfy the rank condition (IV-R). The next proposition specifies sufficient conditions for $Z \equiv (W'X, X)$ to satisfy (IV-R). These conditions are primitive, i.e., they are expressed just in terms of moments of functions of $(X, G)$.[10]

---

[10]We can use the same steps as in the proof of Proposition 3 to derive similar conditions for (IV-R) when the instruments are $H'X$. Those conditions are omitted from the text for brevity.

PROPOSITION 3: *Suppose (A1), (A2), (A3), and (A4) hold, and $E(X'X)$ is non-singular. Let $M \equiv (I - \lambda G)^{-1}$. Then (IV-R) holds for $Z \equiv (W'X, X)$ if*

$$\begin{pmatrix} E(X'X) & E(X'M^{-1}X) \\ E(X'MX) & E(X'X) \end{pmatrix} \text{ and } \begin{pmatrix} E(X'G^2X) & E(X'GX) \\ E(X'GX) & E(X'X) \end{pmatrix} \text{ are non-singular.} \quad (9)$$

These primitive conditions are weak restrictions on the distribution of $(G, X)$; they only serve to rule out "knife-edge" cases where the link formation process is aligned with the regressor distribution in such a pathological way that the rank of moments above is reduced. Our simulation shows (9) holds even for restrictive cases where dyadic links are formed as i.i.d. Bernoulli, and independent from $X$. On the other hand, (9) fails in some other special cases. One example is the linear-in-means social interactions model, where $G$ is proportional to a linear combination of $I$ and a square matrix of ones. Note this linear-in-means model would not be identified even if $G$ were correctly observed, due to the "reflection" problem as defined in Manski (1993). See, e.g., Bramoullé et al. (2009), who require that $I$, $G$, and $G^2$ be perfectly observed and linearly independent.

### 3.3.2. *Instruments using multiple measures*

The method for constructing instruments in Section 3.3.1 assumes the sample reports a single *unsymmetrized* network measure $H$. In this section, we provide an alternative, complementary method for constructing instruments when the sample provides two (or more) measures of $G$, regardless of whether the measures are symmetrized or not.

For example, Banerjee et al. (2013) provide multiple measures of symmetrized links between households in rural villages across the State of Karnataka, India. Two such measures involve visiting between households. For each pair of households, the survey asks which households you visited, and which ones visited you. Banerjee et al. (2013) symmetrize each of these two measures, yielding symmetric matrices we call $H^{(1)}$ and $H^{(2)}$. These two matrices are both measures of the same underlying symmetric network $G$ (where $G_{ij}$ is one if either $i$ visited $j$ or $j$ visited $i$, and zero otherwise). However, as we show later, these two matrices empirically differ substantially, indicating that they are different noisy measures of $G$.

Suppose we observe two measures of the adjacency matrix, $H^{(1)}$ and $H^{(2)}$, which satisfy (A1), (A2), (A3), and

(A4') Conditional on $(G, X)$, $H_{ij}^{(1)}$ and $H_{kl}^{(2)}$ are independent for all $i \neq k$ or $j \neq l$.

These two measures $H^{(1)}$ and $H^{(2)}$ have their own misclassification rates, denoted $(p_0^{(t)}, p_1^{(t)})$ for $t = 1, 2$ respectively. Condition (A4') is plausible when these distinct measures are constructed independently using responses from separate survey questions.

Define

$$W^{(t)} \equiv W_{(H,p_0,p_1)}^{(t)} \equiv \frac{H^{(t)} - p_0^{(t)}(\iota\iota' - I)}{1 - p_0^{(t)} - p_1^{(t)}}.$$

Using either $W^{(1)}$ or $W^{(2)}$, we can construct a structural form. That is, for $t = 1, 2$,

$$y = \lambda W^{(t)} y + X\beta + v^{(t)}, \text{ where } v^{(t)} = \varepsilon + \lambda \left[ G - W^{(t)} \right] y. \tag{10}$$

Under (A1)-(A3) and (A4') and by an argument similar to Proposition 2, we can show that $W^{(2)}X$ and $H^{(2)}X$ satisfy instrument exogeneity with regard to $v^{(1)}$:

$$E\left[ (W^{(2)}X)'v^{(1)} \right] = \frac{1}{1 - p_0^{(2)} - p_1^{(2)}} E\left[ (H^{(2)}X)'v^{(1)} \right] = 0,$$

and likewise with $W^{(2)}$ replaced by $H^{(2)}$. A symmetric result holds by swapping the indexes $t = 1, 2$ in the display above. (See the Online Appendix for details.) We can therefore use either $H^{(1)}X$ or $W^{(1)}X$ as instruments for $W^{(2)}y$ or use either $W^{(2)}X$ or $H^{(2)}X$ as instruments for $W^{(1)}y$. In Section 4, we discuss how to construct 2SLS estimators that combine these multiple network measures.

Note that unlike the instruments in Section 3.3.1 that required an asymmetric $H$, the use of multiple $H^{(t)}$ matrices described here works regardless of whether each $H^{(t)}$ is symmetric or not.

### 3.4. *Recovering misclassification rates*

To construct $W$ and apply 2SLS, we still need to identify and estimate the unknown misclassification rates $p_0$ and $p_1$. We will show how to recover these rates from the observation

of noisy network measures. The main idea is to leverage variation in $X$ that affects true link formation.

### 3.4.1. *Using two conditionally independent measures*

We start with the case where the sample reports two independent measures $H^{(1)}$ and $H^{(2)}$ with misclassification rates $\left(p_0^{(t)}, p_1^{(t)}\right)$ for $t = 1, 2$ respectively, and satisfy (A1), (A2), (A3), and (A4') as before. [11] Our goal is to estimate these misclassification rates.

Assume that we can construct some function of $X$ that is correlated with network formation. Specifically, assume we can define a function $\phi_{ij}(X)$ that is related in some way to the probability that $G_{ij}$ equals zero vs one. In the simplest case $\phi_{ij}(X)$ would be binary valued, with $G_{ij}$ having a different unknown probability of equalling one when $\phi_{ij}(X) = 0$ than when $\phi_{ij}(X) = 1$.

Note this construction imposes no restriction on the true link formation process other than being correlated in some way with $X$. For example, we can accommodate polar extreme cases, such as endogenous network formation based on pairwise stability, where $G_{ij}$ depends on the demographics of all group members $X$, vs dyadic link formation models where $G_{ij}$ depends only on pair-specific demographics $(X_i, X_j)$.

To illustrate, in our empirical application in Section 7 we define $\phi_{ij}(X) \equiv 1\{X_{i,1} = X_{j,1}\}$, where $1\{\cdot\}$ is the indicator function and $X_{i,1}$ is $i$'s caste. So $\phi_{ij}(X) = 1$ if $i$ and $j$ are from the same caste, otherwise $\phi_{ij}(X) = 0$. In this example the required assumption is that two people of the same caste have a different probability of forming a link than two people from different castes.

The intuition for our identification is as follows. Let $\pi_1$ denote the unknown average probability that a cell $G_{ij}$ equals one, conditional on $\phi_{ij}(X) = 1$. If we then consider the average probability (which we can estimate) that a cell $H_{ij}^{(t)}$ equals one, conditional on $\phi_{ij}(X) = 1$, this probability will be a known function of $\pi_1$, $p_0^{(t)}$, and $p_1^{(t)}$ for $t = 1, 2$. This provides two equations (one for each value of $t$) in the unknown misclassification probabilities and in $\pi_1$. The same construction conditioning on $\phi_{ij}(X) = 0$ gives two more equations

---

[11] It is worth emphasizing that this case is flexible enough to accommodate *both scenarios* in Section 3.1. That is, the two independent measures $H^{(1)}, H^{(2)}$ may either be *unsymmetrized* or *symmetrized*, as introduced in Section 3.1. Recall that in the first scenario researchers do not know whether the actual adjacency $G$ is symmetric or not, while in the second scenario researchers do know the actual $G$ is symmetric with undirected links.

in the unknown misclassification probabilities and in $\pi_0$. Finally, looking at the conditional average probability that the product $H_{ij}^{(1)} H_{ij}^{(2)}$ equals one gives two more equations for identification.

Making this logic precise, define $\pi_1 \equiv \frac{1}{n(n-1)} \sum_{i \neq j} \Pr\{G_{ij} = 1 | \phi_{ij}(X) = 1\}$. Consider the following set of three conditional moments of $H_{ij}^{(1)}$ and $H_{ij}^{(2)}$:

$$\frac{1}{n(n-1)} \sum_{i \neq j} E\left[ H_{ij}^{(1)} H_{ij}^{(2)} \Big| \phi_{ij}(X) = 1 \right] = \left(1 - p_1^{(1)}\right)\left(1 - p_1^{(2)}\right) \pi_1 + p_0^{(1)} p_0^{(2)} (1 - \pi_1);$$

$$\frac{1}{n(n-1)} \sum_{i \neq j} E\left[ H_{ij}^{(t)} \Big| \phi_{ij}(X) = 1 \right] = \left(1 - p_1^{(t)}\right) \pi_1 + p_0^{(t)} (1 - \pi_1) \text{ for } t = 1, 2. \quad (11)$$

Note these are three distinct equations because the second applies for both $t = 1$ and $t = 2$. We obtain three more equations (six in total) by replacing $\phi_{ij}(X) = 1$ with $\phi_{ij}(X) = 0$ and replacing $\pi_1$ with $\pi_0$. The left-hand side of each of these six equations can be estimated from our observations of $H^{(1)}$, $H^{(2)}$, and $X$, while the right-hand sides are functions of six unknown parameters: $\pi_1, \pi_0$ and $p_1^{(t)}, p_0^{(t)}$ for $t = 1, 2$. Assume that $\pi_1 \neq \pi_0$, meaning that $\phi_{ij}(X)$ does affect the probability of true link formation. Then despite the nonlinearity of these equations, we show that they can be uniquely solved for these six parameters, and in particular we provide closed-form expressions for the misclassification rates $p_1^{(t)}, p_0^{(t)}$ for $t = 1, 2$. See the proof in the Online Appendix for details.

This identification requires choosing a function $\phi_{ij}(\cdot)$ such that the probability of link formation is different for the event $\{\phi_{ij}(X) = 1\}$ than when $\{\phi_{ij}(X) = 0\}$ so $\pi_1 \neq \pi_0$. In other words, these conditioning events provide exogenous variation in population moments that assist in identifying the misclassification rates.

It should also be noted that our focus here is just on recovering the misclassification rates. We treat $\pi_1, \pi_0$ as "nuisance" parameters that are identified as an intermediate step in our constructive identification of $p_1^{(t)}, p_0^{(t)}$ for $t = 1, 2$. We do not exploit knowledge of $\pi_1, \pi_0$ for estimation, or to infer anything about the link formation process.

We can generalize the identification argument above to broader settings with other choices of $\phi_{ij}(\cdot)$. For instance, $\phi_{ij}(X)$ may be a continuous measure of the difference between demographic features of $i$ and $j$. In this case, one can partition the support of $\phi_{ij}(X)$ into mutually exclusive subsets, denoted by $\phi^0$ and $\phi^1$. Then define $\pi_0 \equiv$

$\frac{1}{n(n-1)} \sum_{i \neq j} \Pr\{G_{ij} = 1 | \phi_{ij}(X) \in \boldsymbol{\phi^0}\}$; define $\pi_1$ analogously conditioning on $\{\phi_{ij}(X) \in \boldsymbol{\phi^1}\}$. The constructive identification strategy above applies with events $\{X_{i,1} = X_{j,1}\}$, $\{X_{i,1} \neq X_{j,1}\}$ replaced by $\{\phi_{ij}(X) \in \boldsymbol{\phi^1}\}$, $\{\phi_{ij}(X) \in \boldsymbol{\phi^0}\}$ respectively.

### 3.4.2. *Using a single, unsymmetrized measure*

The identification method of the previous section can be readily modified to recover the misclassification probabilities in the case with a *single, unsymmetrized* measure $H$ when the actual $G$ is *known* to be symmetric with undirected links. Suppose $H$ satisfies (A1), (A2), (A3), and (A4) with misclassification rates $p_1, p_0$. For any *unordered* pair $i \neq j$, construct two noisy measures for $G_{ij}$ as $H^{(1)}_{\{i,j\}} \equiv H_{ij}$ and $H^{(2)}_{\{i,j\}} \equiv H_{ji}$. (The adoption of new subscripts for $H^{(t)}$, i.e., $\{i,j\}$, only serves as a reminder that these two measures are symmetrized by construction.) We then obtain a system of equations similar to (11), only with $\frac{1}{n(n-1)}$, $\sum_{i \neq j}$, $H^{(t)}_{ij}$, $\phi_{ij}$ replaced by $\frac{2}{n(n-1)}$, $\sum_{i>j}$, $H^{(t)}_{\{i,j\}}$, $\phi_{\{i,j\}}$ respectively, and with identical rates across the measures, i.e. $p_1^{(t)} = p_1$ and $p_0^{(t)} = p_0$ for $t = 1, 2$. The same argument then identifies $\pi_1, \pi_0, p_1, p_0$ using variation in $\phi_{\{i,j\}}(X)$.

### 3.5. *Concluding remarks about identification*

The methods proposed in Section 3 are flexible enough to accommodate various scenarios regarding whether $G$ is symmetric or not, and whether the observed network measure(s) is(are) symmetrized or not. The table below summarizes the solutions of adjusted 2SLS that we propose for each one of those scenarios.

<div align="center">Reported Network Measures</div>

|  | Single, unsym'zed | | Multiple, sym'zed | | Multiple, unsym'zed | |
|---|---|---|---|---|---|---|
|  | (IV) | (MR) | (IV) | (MR) | (IV) | (MR) |
| Sym. $G$ | Sec 3.3.1 | Sec 3.4.2 | Sec 3.3.2 | Sec 3.4.1 | Sec 3.3 | Sec 3.4 |
| Asym. $G$ | Sec 3.3.1 | see text | violates (A1) | | Sec 3.3.2 | Sec 3.4.1 |

Each one of the six cells in last two rows of the table represents a particular scenario, defined by the (a)symmetry of the actual adjacency $G$ *as well as* the number and property of network measures $H$ available. Solutions for estimating $\lambda$ and $\beta$ in each scenario consist of two parts: construction of instruments (IV), and recovery of misclassification rates (MR).

22

For instance, if the actual $G$ is symmetric and the sample reports a single, unsymmetrized measure, one can recover MRs using Section 3.4.2 and construct IVs using Section 3.3.1. Likewise, if the actual $G$ is asymmetric and the sample reports multiple, unsymmetrized measures, one can recover MRs using Section 3.4.1 and construct IVs using Section 3.3.2. If the actual $G$ is symmetric and the sample reports multiple, unsymmetrized measures, then one can recover MRs using *either* approach in Section 3.4, and construct IVs using *either* approach in Section 3.3.

For the scenario with an asymmetric $G$ and a single, unsymmetrized measure, our paper presents a valid way to construct instruments, but does not propose a way to identify the MRs. To perform the latter task, one might be able to adopt a method from Hausman et al. (1998) to a dyadic link formation model. We do not elaborate on that method in this paper, because it would require researchers to specify a link formation model, which we have intentionally refrained from doing throughout this paper.

Some additional remarks about our use of multiple, noisy network measures in Section 3.3.2 and 3.4.1 are in order. There is a broad and growing econometrics literature that uses repeated noisy measures to estimate nonlinear models with errors in variables, e.g., Li (2002), Chen et al. (2005) and Hu and Sasaki (2017) or unobserved heterogeneity, e.g., Hu (2008) and Bonhomme et al. (2016). Hu and Lin (2018) use repeated measurement to estimate a binary choice model with misclassification and social interactions. These papers typically apply mathematical tools such as deconvolution, and eigenvalue or LU decomposition to the distribution of repeated measures.

In contrast, we use the repeated measures in a different way that does not require any deconvolution or matrix decomposition. Focusing on linear social networks, we exploit the identifying power from repeated measures by a standard 2SLS in Section 3.3.2, and apply a closed-form algebraic argument to recover the misclassification rates in Section 3.4.1.

Finally, note that our 2SLS estimators are unlikely to suffer from weak instrument issues, because Assumption (A2) ensures correlation between mismeasures $H$ and $G$, and our instruments are constructed from $H$.

## 4. TWO-STEP ESTIMATION

We now propose adjusted 2SLS estimators for the coefficients of structural effects $(\lambda, \beta')'$, which require an initial step for estimating the misclassification rates.

Consider a sample of $S$ independent groups. (In the Online Appendix we consider extensions to a single growing network instead of many independent groups.) For each group $s = 1, ..., S$, the sample reports an $n_s$-by-1 vector of outcomes $y_s$, an $n_s$-by-$K$ matrix of regressors $X_s$, and either an $n_s$-by-$n_s$ unsymmetrized measure $H_s$, or two $n_s$-by-$n_s$ conditionally independent symmetrized measures $H_s^{(1)}$ and $H_s^{(2)}$.

### 4.1. *Closed-form estimation of misclassification rates*

To estimate misclassification rates, we apply the analog principle to the constructive proof of identification. We include closed-form estimates in the text for completeness; the logic for these estimators is self-evident as presented in the Online Appendix.

First, consider the case in Section 3.4.1, where the sample reports two conditionally independent measures $H_s^{(1)}, H_s^{(2)}$. To exploit identifying power from their joint distribution, let $H_{s,ij}^{(3)} \equiv \max\left\{ H_{s,ij}^{(1)}, H_{s,ij}^{(2)} \right\}$ for each $(i,j)$-th entry in $H_s^{(t)}$. For $t = 1, 2, 3$, define $\hat{\psi}_1^{(t)}$:

$$\hat{\psi}_1^{(t)} \equiv \frac{\sum_s \left[ \frac{1}{n_s(n_s-1)} \left( \sum_{i \neq j} H_{s,ij}^{(t)} 1\{\phi_{s,ij} = 1\} \right) \right]}{\sum_s \left[ \frac{1}{n_s(n_s-1)} \left( \sum_{i \neq j} 1\{\phi_{s,ij} = 1\} \right) \right]}, \tag{12}$$

where $\phi_{s,ij}$ is short for $\phi_{ij}(X_s)$. And define $\hat{\psi}_0^{(t)}$ by replacing $\{\phi_{s,ij} = 1\}$ with $\{\phi_{s,ij} = 0\}$.

For instance, in our application, we define $\phi_{ij}(X_s)$ as a simple function $1\{X_{s,i,k} = X_{s,j,k}\}$, where $X_{s,i,k}$ is the $k$-th component in $X_{s,i}$ that reports the individual $i$'s caste. In this case, $\hat{\psi}_1^{(t)}$ and $\hat{\psi}_0^{(t)}$) are, respectively, the fraction of same caste and different caste pairs that are linked according to the measures $H_s^{(t)}$ for $t = 1, 2, 3$. It is straightforward to generalize this identification argument to broader settings with other choices of $\phi_{ij}(\cdot)$.

Using the sample moments, we define a vector of coefficients:

$$\widehat{\mathcal{C}}_2 \equiv \frac{\hat{\psi}_0^{(1)} - \hat{\psi}_1^{(1)}}{\hat{\psi}_0^{(2)} - \hat{\psi}_1^{(2)}}, \ \widehat{\mathcal{C}}_1 \equiv \hat{\psi}_1^{(1)} - 1 + \frac{\hat{\psi}_0^{(3)} - \hat{\psi}_1^{(3)}}{\hat{\psi}_0^{(2)} - \hat{\psi}_1^{(2)}} - (1 - \hat{\psi}_1^{(2)})\widehat{\mathcal{C}}_2,$$

$$\widehat{\mathcal{C}}_0 \equiv \hat{\psi}_1^{(1)} + \hat{\psi}_1^{(2)} - \hat{\psi}_1^{(1)}\hat{\psi}_1^{(2)} - \hat{\psi}_1^{(3)}.$$

24

Our closed-form estimators for misclassification rates are then:

$$\hat{p}_0^{(1)} \equiv \hat{\psi}_1^{(1)} - \widehat{\mathcal{C}}_2\hat{\xi}, \ \hat{p}_0^{(2)} \equiv \hat{\psi}_1^{(2)} - \hat{\xi},$$

where

$$\hat{\xi} \equiv (2\widehat{\mathcal{C}}_2)^{-1}\left(\widehat{\mathcal{C}}_1 + \sqrt{\left(\widehat{\mathcal{C}}_1\right)^2 + 4\widehat{\mathcal{C}}_2\widehat{\mathcal{C}}_0}\right);$$

and

$$\hat{p}_1^{(t)} \equiv 1 - \hat{p}_0^{(t)} - \frac{\hat{\psi}_1^{(t)} - \hat{p}_0^{(t)}}{\hat{\pi}_1} \text{ for } t = 1, 2,$$

where

$$\hat{\pi}_1 = \frac{\left(\hat{\psi}_1^{(1)} - \hat{p}_0^{(1)}\right)\left(\hat{\psi}_1^{(2)} - \hat{p}_0^{(2)}\right)}{\left(1 - \hat{p}_0^{(1)}\right)\left(\hat{\psi}_1^{(2)} - \hat{p}_0^{(2)}\right) + \left(1 - \hat{p}_0^{(2)}\right)\left(\hat{\psi}_1^{(1)} - \hat{p}_0^{(1)}\right) - \left(\hat{\psi}_1^{(3)} - \hat{p}_0^{(3)}\right)},$$

with $\hat{p}_0^{(3)} \equiv \hat{p}_0^{(1)} + \hat{p}_0^{(2)} - \hat{p}_0^{(1)}\hat{p}_0^{(2)}$ by construction.

Next, consider the case in Section 3.4.2, where the sample reports a single, unsymmetrized measure $H$ with misclassification rates $(p_0, p_1)$ while the actual $G$ is known to be symmetric. Estimation of $(p_0, p_1)$ in this case follows from almost identical steps. For any *unordered* pair $\{i, j\}$, define $H_{s,\{i,j\}}^{(1)} \equiv H_{s,ij}$ and $H_{s,\{i,j\}}^{(2)} \equiv H_{s,ji}$. By construction, $p_1^{(t)} = p_1$ and $p_0^{(t)} = p_0$ do not vary between $t = 1, 2$. Construct $H_{\{i,j\}}^{(3)} = \max\{H_{\{i,j\}}^{(1)}, H_{\{i,j\}}^{(2)}\}$; define $\hat{\psi}_1^{(t)}$ and $\hat{\psi}_0^{(t)}$ in this case by replacing $\frac{1}{n_s(n_s-1)}$, $\sum_{i \neq j}$ and $H_{s,ij}^{(t)}$ in (12) with $\frac{2}{n_s(n_s-1)}$, $\sum_{i>j}$, and $H_{s,\{i,j\}}^{(t)}$ respectively. Replace $\widehat{\mathcal{C}}_2$ with 1, and replace $\hat{\psi}_1^{(1)}, \hat{\psi}_1^{(2)}$ with their average in $\widehat{\mathcal{C}}_1, \widehat{\mathcal{C}}_0$ and all subsequent expressions. These lead to a single pair of estimates $(\hat{p}_0, \hat{p}_1)$.

We derive the limiting distribution of these estimators using a standard delta method. Consider the case with a single unsymmetrized measure in Section 3.4.2. For each group $s$, define $v_{1s,1} \equiv \frac{2}{n_s(n_s-1)}\sum_{i>j}H_{s,\{i,j\}}1\{\phi_{s,\{i,j\}} = 1\}$ and $v_{2s,1} \equiv \frac{2}{n_s(n_s-1)}\sum_{i>j}1\{\phi_{s,\{i,j\}} = 1\}$; define $v_{1s,0}, v_{2s,0}$ analogously be replacing $\phi_{s,\{i,j\}} = 1$ with $\phi_{s,\{i,j\}} = 0$. Let $v_s \equiv (v_{1s,1}, v_{2s,1}, v_{1s,0}, v_{2s,0})'$. The estimator $\hat{p} = (\hat{p}_0, \hat{p}_1)$ is a closed-form function of $v_s$; it has an asymptotic linear presentation:

$$\sqrt{S}(\hat{p} - p) = \frac{1}{\sqrt{S}}\sum_s \underbrace{\mathcal{J}_0 \times [v_s - E(v_s)]}_{\equiv \tau_s} + o_p(1),$$

where $\mathcal{J}_0$ denotes the Jacobian matrix of $\hat{p}$ w.r.t. the sample averages of $\upsilon_s$, evaluated at population mean of $\upsilon_s$. Thus $\sqrt{S}(\widehat{p} - p)$ converges in distribution to a multivariate normal distribution with zero means and a covariance matrix $E(\tau_s \tau_s')$. Limiting distribution for the case with two measures in Section 3.4.1 follows from the same type of arguments.

## 4.2. *Adjusted 2SLS using a single unsymmetrized measure*

With estimates of misclassification rates, we can now construct adjusted 2SLS estimators for the peer and individual effects $\lambda$ and $\beta$. As noted in Section 3.3, construction of instruments depends on the number and nature of network measures available.

First consider the setting in Section 3.3.1, where the sample reports a single *unsymmetrized* measure $H_s$ for each group. Let $p \equiv (p_0, p_1)'$. For each group $s$, define:

$$R_s(p) \equiv (W_s(p)y_s, X_s) \text{ and } Z_s \equiv \left(H_s' X_s, X_s\right),$$

where $W_s(p) \equiv [H_s - p_0(\iota_s \iota_s' - I_s)]/(1 - p_0 - p_1)$. Let $N \equiv \sum_{s=1}^S n_s$, and $Y$ be an $N$-by-1 vector that stacks $y_s$ for $s = 1, ..., S$. Let $\mathbf{R}(p)$ be an $N$-by-$(K+1)$ matrix that stacks $R_s(p)$ for all group $s$, and $\mathbf{Z}$ an $N$-by-$2K$ matrix that stacks $Z_s$ for all $s$. Our adjusted 2SLS estimator for $\theta \equiv (\lambda, \beta')'$ is:

$$\widehat{\theta} \equiv \left(\mathbf{A}'\mathbf{B}^{-1}\mathbf{A}\right)^{-1} \mathbf{A}'\mathbf{B}^{-1}\left(\mathbf{Z}'Y\right), \tag{13}$$

where $\mathbf{A} \equiv \mathbf{Z}'\mathbf{R}(\widehat{p})$ and $\mathbf{B} \equiv \mathbf{Z}'\mathbf{Z}$, with $\hat{p} \equiv (\hat{p}_0, \hat{p}_1)'$.

We now present the limiting distribution of $\widehat{\theta}$ as $S \to \infty$. Define

$$\Sigma_0 \equiv \left(A_0' B_0^{-1} A_0\right)^{-1} A_0' B_0^{-1},$$

where $A_0 \equiv \lim_{S \to \infty} \frac{1}{S} \sum_{s=1}^S E\left[Z_s' R_s(p)\right]$ and $B_0 \equiv \lim_{S \to \infty} \frac{1}{S} \sum_{s=1}^S E(Z_s' Z_s)$. For each group $s$ and individual $i \leq n_s$, let $R_{s,i}(p)$ denote the corresponding row in $\mathbf{R}(p)$, and $\nabla_p R_{s,i}(p)$ be the $(K+1)$-by-2 Jacobian of $R_{s,i}(p)$ with respect to $p$.[12]

Let $\nabla_p \left[R_s(p)\theta\right]$ denote an $n_s$-by-2 matrix with each row $i \leq n_s$ being $\theta' \nabla_p R_{s,i}(p)$; let $\nabla_p \left[\mathbf{R}(p)\theta\right]$ be an $N$-by-2 matrix formed by stacking these $n_s$-by-2 matrices over

---

[12]The last $K$ rows in $\nabla_p R_{s,i}(p)$ are zeros; its first row is the $i$-th row in $\left(\frac{H_s - (1-p_1)(\iota_s \iota_s' - I_s)}{(1-p_0-p_1)^2} y_s, \frac{H_s - p_0(\iota_s \iota_s' - I_s)}{(1-p_0-p_1)^2} y_s\right)$.

26

$s = 1, 2, ..., S$. Define

$$\kappa_s \equiv Z_s' v_s - F_0 \tau_s,$$

where $v_s$ is the $n_s$-by-1 vector of composite errors in the feasible structural form (6), and

$$F_0 \equiv \lim_{S \to \infty} S^{-1} \sum_{s=1}^{S} E \left\{ Z_s' \nabla \left[ R_s(p)\theta \right] \right\}.$$

Intuitively, $F_0$ illustrates how the moment condition in this adjusted 2SLS depends on mis-classification rates $p$, and the added term "$-F_0\tau_s$" in the influence function accounts for the first-stage estimation error in $\widehat{p}$.

PROPOSITION 4: *Suppose (A1), (A2), (A3), and (A4) hold, and (IV-R) is satisfied with* $Z \equiv (H'X, X)$. *Then*

$$\sqrt{S} \left( \widehat{\theta} - \theta \right) \xrightarrow{d} \mathcal{N}(0, \Sigma_0 E(\kappa_s \kappa_s') \Sigma_0'),$$

*under the regularity conditions (REG) in the Online Appendix.*

Note that this limiting distribution includes group level clustering. The conditions in (*REG*), presented in the Online Appendix, are needed for applying the law of large numbers, the central limit theorem, and the delta method to observations from independent groups with heterogeneous sizes.

Standard errors for $\widehat{\theta}$ (which are clustered at the group level) are calculated by replacing $A_0$, $B_0$, $F_0$, and $E(\kappa_s \kappa_s')$ with their sample analogs:

$$\widehat{A} = \tfrac{1}{S} \sum_s Z_s' R_s(\widehat{p}), \ \ \widehat{B} = \tfrac{1}{S} \sum_s Z_s' Z_s, \ \ \widehat{\kappa}_s = Z_s' \left( y_s - R_s(\widehat{p})\widehat{\theta} \right) - \widehat{F}\widehat{\tau}_s.$$

Instead of the above limiting distribution, one could combine the two steps in Section 4.1 and 4.2 into a single GMM step by stacking the moments used in these two sections. This would allow one to estimate $\theta$ jointly with $(p_0, p_1)$ instead of sequentially, and standard GMM limiting distribution theory could be applied. However, this GMM would require numerically solving a nonlinear optimization problem. In contrast, the two-step method we propose here yields a closed-form estimator that is straightforward to compute with no numerical searching, thus providing a computational advantage over the GMM alternative with numerical stability.

## 4.3. *Adjusted 2SLS using multiple measures*

We now apply the same idea for estimation under the other setting in Section 3.3.2, where the sample reports two conditionally independent measures $H_s^{(t)}$ for $t = 1, 2$, with misclassification rates $p_0^{(t)}, p_1^{(t)}$ for $t = 1, 2$ respectively. These measures may either be *symmetrized* or *unsymmetrized*. To reiterate, when $H_s^{(t)}$ are *unsymmetrized*, our estimation method applies *regardless of* whether the actual adjacency $G$ is symmetric or not; on the other hand, when $H_s^{(t)}$ are symmetrized, (A1) holds only if $G$ is symmetric.

As noted in Section 3.3.2, these measures lead to two feasible structural forms:

$$y_s = R_s^{(t)}\theta + v_s^{(t)} \text{ for } t = 1, 2, \tag{14}$$

where $\theta \equiv (\lambda, \beta')'$, $R_s^{(t)} \equiv \left(W_s^{(t)}y_s, X_s\right)$ and $v_s^{(t)} \equiv \varepsilon_s + \lambda\left(G_s - W_s^{(t)}\right)y_s$, with $W_s^{(t)} \equiv \frac{H_s^{(t)} - p_0^{(t)}(\iota_s\iota_s' - I_s)}{1 - p_0^{(t)} - p_1^{(t)}}$. This leads to two sets of moment conditions:

$$E\left[(H_s^{(3-t)}X, X)'(y_s - \lambda W_s^{(t)}y_s - X_s\beta)\right] = E\left[(H_s^{(3-t)}X, X)'v_s^{(t)}\right] = 0 \text{ for } t = 1, 2,$$

with instruments $Z_s^{(t)} \equiv \left(H_s^{(3-t)}X_s, X_s\right)$ for $t = 1, 2$. Stack the moments by defining:

$$\tilde{Z}_s \equiv \begin{pmatrix} Z_s^{(1)} & 0 \\ 0 & Z_s^{(2)} \end{pmatrix}; \tilde{y}_s \equiv \begin{pmatrix} y_s \\ y_s \end{pmatrix}; \tilde{R}_s \equiv \begin{pmatrix} R_s^{(1)} \\ R_s^{(2)} \end{pmatrix}.$$

Instrument exogeneity then implies:

$$E\left[\tilde{Z}_s'(\tilde{y}_s - \tilde{R}_s\theta)\right] = 0.$$

This moment condition identifies $\theta$, provided $E(\tilde{Z}_s'\tilde{R}_s)$ has full rank. Using arguments similar to Proposition 3 in Section 3.3.1, we can derive analogous sufficient conditions for this rank condition. We omit the details here so as to avoid repetition.

We define a system, or stacked adjusted two-stage least squares (S2SLS) estimator as follows. Let $\tilde{\mathbf{Z}}$ denote a $2N$-by-$4K$ matrix that is constructed by vertically stacking $S$ matrices $(\tilde{Z}_s)_{s \leq S}$. Likewise, construct a $2N$-by-$(K+1)$ matrix $\tilde{\mathbf{R}}$ by stacking $(\tilde{R}_s)_{s \leq S}$, where $p_0^{(t)}$ and $p_1^{(t)}$ are replaced by estimates $\hat{p}_0^{(t)}$ and $\hat{p}_1^{(t)}$, and construct a $2N$-by-1 vector $\tilde{\mathbf{y}}$ by stacking $(\tilde{y}_s)_{s \leq S}$. The S2SLS estimator is

$$\tilde{\theta} \equiv [\tilde{\mathbf{R}}'\tilde{\mathbf{Z}}(\tilde{\mathbf{Z}}'\tilde{\mathbf{Z}})^{-1}\tilde{\mathbf{Z}}'\tilde{\mathbf{R}}]^{-1}\tilde{\mathbf{R}}'\tilde{\mathbf{Z}}(\tilde{\mathbf{Z}}'\tilde{\mathbf{Z}})^{-1}\tilde{\mathbf{Z}}'\tilde{\mathbf{y}}. \tag{15}$$

This provides us with a single estimator that exploits both sets of instruments in the two structural forms in (14). Similar to $\hat{\theta}$ in (13), we can construct the standard error for $\tilde{\theta}$ that accounts for estimation error in $\hat{p}_0^{(t)}, \hat{p}_1^{(t)}$ for $t = 1, 2$. We omit details here for brevity.

## 5. EXTENSIONS

We now extend the method in Section 3 and 4 to more general settings with contextual effects, heterogeneous misclassification rates, or group fixed effects. In each case, we focus on extending the ideas for constructive identification. Estimation in each case follows from an analog principle and similar steps as in Section 4.

As in Section 3, to simplify exposition, we let group sizes $n_s = n$ be fixed throughout the remainder of this section. This allows us to suppress group subscripts $s$ in notation.

### 5.1. *Contextual effects*

Suppose the structural form, based on perfect observation of the actual adjacency $G$, is:

$$y = \lambda G y + X\beta + GX\gamma + \varepsilon,$$

where $\gamma$ are contextual effects showing how individual outcomes are directly influenced by the characteristics of others linked to the individual. The feasible structural form, based on $H$ and subject to misclassification errors, is:

$$y = \lambda W y + X\beta + WX\gamma + \eta,$$

where $W \equiv [H - p_0(\iota\iota' - I)]/(1 - p_0 - p_1)$ as before, and the composite error $\eta$ is:

$$\eta \equiv \varepsilon - \lambda \left( W - G \right) y - \left( W - G \right) X\gamma.$$

Under the same conditions and by the same arguments as in the case with no contextual effects in Section 3.2, we can show that the new composite error $\eta$ is mean-independent from $(X, G)$. Similarly, we can construct instruments using network measures $H$ as before. Our next proposition establishes these results. For generality, let $\zeta(X) \in \mathbb{R}^{n \times L}$ be any generic function of $X$ with $L \geq K$.

PROPOSITION 5: *Suppose (A1), (A2), and (A3) hold. Then $E(\eta|X,G) = 0$. If in addition (A4) holds, then $E[\zeta\left(X\right)' W\eta] = E[\zeta\left(X\right)' H\eta] = 0$.*

This proposition implies that $H'\zeta(X)$ and $W'\zeta(X)$ satisfy instrument exogeneity for generic functions of $X$. In fact, a stronger result holds under (A1)-(A4): $E(W\eta|G, X) = E(H\eta|G, X) = 0$. The intuition is the same as in Proposition 2. Thus we can apply 2SLS as before to consistently estimate $(\lambda, \beta', \gamma')'$ using $(H'X, X, H'\zeta(X))$ as instruments for $(Wy, X, WX)$, provided the appropriate rank conditions hold.

### 5.2. *Heterogeneous misclassification rates*

We now extend our methods to allow the misclassification rates $p_0, p_1$ to vary with individual characteristics $X$. To fix ideas, we return to the case with no contextual effects as in (6). Generalization to include contextual effects, using the results from the preceding sub-section, is immediate. Suppose we relax (A2) as follows:

(A2') $E(H_{ij}|G_{ij} = 1, X) = 1 - p_{ij,1}(X)$ and $E(H_{ij}|G_{ij} = 0, X) = p_{ij,0}(X) \; \forall i \neq j$.

Define

$$W_{ij}(X) \equiv \frac{H_{ij} - p_{ij,0}(X)}{1 - p_{ij,0}(X) - p_{ij,1}(X)} \text{ if } i \neq j, \text{ and } W_{ii}(X) = 0.$$

Under (A2'), $E[W_{ij}(X)|G, X] = 1$ for $G_{ij} = 1$, and $E[W_{ij}(X)|G, X] = 0$ for $G_{ij} = 0$. Hence $E(W(X)|G, X) = G$.

To recover misclassification rates $p_{ij,1}(\cdot)$ and $p_{ij,0}(\cdot)$, we can apply methods in Section 3.4 to pairwise links $G_{ij}$ and conditioning on $X$. In practice, we can mitigate the curse of dimensionality by specifying the rates $p_{ij,1}(X)$ and $p_{ij,0}(X)$ as functions of $X_i$ and $X_j$.

With knowledge of these heterogeneous misclassification rates, we can use adjusted 2SLS to consistently estimate $(\lambda, \beta')'$ from a feasible structural form:

$$y = \lambda W(X)y + X\beta + \underbrace{\varepsilon + \lambda[G - W(X)]y}_{v^*}.$$

Under (A2') and (A3),

$$E(v^*|G, X) = \lambda\{GE(y|G, X) - E[W(X)y|G, X]\}$$

$$= \lambda\{GMX\beta - E[W(X)|G, X]MX\beta\} = \lambda(G - G)MX\beta = 0. \qquad (16)$$

Let $R^* \equiv (W(X)y, X)$ and $Z^* \equiv (\zeta(X), X)$ where $\zeta(X) \in \mathbb{R}^{n \times L}$ is a nonlinear function of $X$ with $L \geq K$ (e.g., $\zeta(X) \equiv X \circ X$, where $\circ$ denotes the Hadamard product of matrices).

30

Then (16) implies $E(Z^{*\prime}v^*) = 0$. If $E(R^{*\prime}Z^*)$ and $E(Z^{*\prime}Z^*)$ have full rank, we can use this adjusted 2SLS to consistently estimate $(\lambda, \beta')'$.

## 5.3. *Group fixed effects*

Suppose there is group-level unobserved heterogeneity $\alpha$ in the data-generating process:

$$y = \lambda G y + X\beta + \alpha + \varepsilon.$$

We can implement the "with-in" transformation on the adjusted network measure $W$, as in fixed-effect estimation of linear panel data models, to get:

$$\dot{\mathcal{W}} \equiv \left[ I - \iota\iota'/n \right] W.$$

Essentially, this transformation just corresponds to demeaning $W$ within groups. Similarly, define with-in transformations on $y, \varepsilon, X, G$ to obtain $\dot{y}, \dot{\varepsilon}, \dot{X}, \dot{G}$ respectively. The resulting demeaned version of the structural form is

$$\dot{y} = \lambda \dot{\mathcal{W}} y + \dot{X}\beta + \underbrace{\dot{\varepsilon} + \lambda(\dot{G} - \dot{\mathcal{W}})y}_{\equiv \dot{v}}.$$

As $\dot{G}$ and $\dot{\mathcal{W}}$ are linear in $G$ and $H$ respectively, the same argument as Lemma 1 in Section 3 applies to show that $E(\dot{v}|X, G) = E(\dot{v}|\dot{X}, \dot{G}) = 0$. Note that the presence of group fixed effects does not affect our method for recovering the misclassification rates in Section 3.4.1. With multiple network measures $H^{(t)}$ for $t = 1, 2$, we can apply adjusted 2SLS as in Section 3.3.2 to estimate $(\lambda, \beta')'$, using $\dot{H}^{(2)}X$ as instruments for $\dot{\mathcal{W}}^{(1)}y$.

## 5.4. *A single large network*

We examine a setting where the sample is partitioned into *approximate* groups, a.k.a. blocks. Sparse links (with diminishing formation rates) exist *between* these blocks, but are not recorded in the sample; links within the blocks can be dense and are randomly misclassified with constant rates. In this case, we show that our adjusted 2SLS estimator, when pooled over all individuals in the sample, still converges to the intended estimand without any endogeneity bias. Detailed setting and proofs are in the Online Appendix.

## 6. SIMULATION

We now use monte carlo simulation to examine the finite-sample performance of the adjusted 2SLS estimator proposed in Section 4.

For the data-generating process, we use a structural form with group-level fixed effects:

$$y_s = \lambda G_s y_s + X_s \beta + \alpha_s + \varepsilon_s, \text{ with } s = 1, ..., S.$$

Each member $i$ in group $s$ has two individual characteristics $X_{s,i} \equiv (X_{s,i,1}, X_{s,i,2}) \in \mathbb{R}^2$, drawn independently across $i$ and $s$, from a Bernoulli with success probability 0.5 and a standard normal $N(0,1)$ respectively. The error term $\varepsilon_{s,i}$ is also drawn from a standard normal $N(0,1)$ independently across $i$ and $s$. The coefficients for social effects are $\lambda = 0.05$ and $\beta = (\beta_1, \beta_2) = (1, 2)$. The group-level fixed effect is $\alpha_s = 5\overline{X}_s \beta - 1.5 + e_s$, where $\overline{X}_s$ is the group average of $X_s$ and $e_s$ is drawn from standard normal $N(0,1)$ independently across $i$ and $s$. This construction allows the fixed effects $\alpha_s$ to be correlated with group demographics $\overline{X}_s \beta$. The dyadic link formation rates are

$$\pi_1 = \Pr\{G_{s,ij} = 1 | X_{s,i,1} = X_{s,j,1}\} = 0.2,$$

$$\pi_0 = \Pr\{G_{s,ij} = 1 | X_{s,i,1} \neq X_{s,j,1}\} = 0.1.$$

For $t = 1, 2$, we generate the following measure $H_s^{(t)}$ with link misclassification:

$$H_{s,ij}^{(t)} = m_{ij,1}^{(t)} \cdot 1\{G_{s,ij} = 1\} + (1 - m_{ij,0}^{(t)}) \cdot 1\{G_{s,ij} = 0\},$$

where $m_{ij,0}^{(t)}$ and $m_{ij,1}^{(t)}$ are drawn independently across ordered pairs $(i,j)$ from Bernoulli distributions with success probabilities $1 - p_0^{(t)}$ and $1 - p_1^{(t)}$ respectively.

To see how various estimators behave in the presence of misclassified links, we use two sets of misclassification rates. In the first set, the misclassification rates are *small*:

$$p_0^{(1)} = \Pr\{H_{s,ij}^{(1)} = 1 | G_{s,ij} = 0, X\} = 0.10, \ p_1^{(1)} = \Pr\{H_{s,ij}^{(1)} = 0 | G_{s,ij} = 1, X\} = 0.20;$$

$$p_0^{(2)} = \Pr\{H_{s,ij}^{(2)} = 1 | G_{s,ij} = 0, X\} = 0.08, \ p_1^{(2)} = \Pr\{H_{s,ij}^{(2)} = 0 | G_{s,ij} = 1, X\} = 0.16.$$

In the second set, we specify *large* misclassification rates that are twice as high:

$$p_0^{(1)} = \Pr\{H_{s,ij}^{(1)} = 1 | G_{s,ij} = 0, X\} = 0.20, \ p_1^{(1)} = \Pr\{H_{s,ij}^{(1)} = 0 | G_{s,ij} = 1, X\} = 0.40;$$

$$p_0^{(2)} = \Pr\{H_{s,ij}^{(2)} = 1 | G_{s,ij} = 0, X\} = 0.16, \ p_1^{(2)} = \Pr\{H_{s,ij}^{(2)} = 0 | G_{s,ij} = 1, X\} = 0.32.$$

Each group has the same size $n_s = n$. We experiment with group sizes $n \in \{25, 50, 100\}$ and the number of groups in the sample $S \in \{50, 100\}$. The total sample size is $nS$. For each combination of $\{n, S\}$, we generate $Q = 100$ samples. For each combination of $\{n, S\}$, Table 1(a) reports the mean and the standard deviation (s.t.d.) of the estimates for $\pi_0, \pi_1, p_0^{(1)}, p_1^{(1)}, p_0^{(2)}, p_1^{(2)}$ based on their empirical distribution across these 100 samples.

From Table 1(a), we can see the misclassification rates $(p_0^{(t)}, p_1^{(t)})$, as well as the network parameters $(\pi_0, \pi_1)$, are accurately estimated in all settings. For a fixed group number $S$, the s.t.d. decreases at the rate $n$. For a fixed groups size $n$, the s.t.d. decreases at the rate $\sqrt{S}$. This is because the size of the sample used for estimation is $S \times n^2$. The standard deviations of these estimates are also larger when the misclassification rates are higher.

**Table 1(a): Estimates of Misclassification Rates and Network Parameters**

| Small | $\pi_1 = 0.2$ | $\pi_0 = 0.1$ | $p_0^{(1)} = 0.1$ | $p_1^{(1)} = 0.2$ | $p_0^{(2)} = 0.08$ | $p_1^{(2)} = 0.16$ |
|---|---|---|---|---|---|---|
| $S = 50$ | $\widehat{\pi}_1$ | $\widehat{\pi}_0$ | $\widehat{p}_0^{(1)}$ | $\widehat{p}_1^{(1)}$ | $\widehat{p}_0^{(2)}$ | $\widehat{p}_1^{(2)}$ |
| $n = 25$ | 0.2009 | 0.1015 | 0.0990 | 0.2020 | 0.0792 | 0.1638 |
| | (0.0123) | (0.0081) | (0.0061) | (0.0301) | (0.0059) | (0.0349) |
| $n = 50$ | 0.1996 | 0.0998 | 0.1002 | 0.2000 | 0.0800 | 0.1573 |
| | (0.0063) | (0.0042) | (0.0031) | (0.0150) | (0.0031) | (0.0186) |
| $n = 100$ | 0.2000 | 0.1002 | 0.1000 | 0.2007 | 0.0798 | 0.1573 |
| | (0.0030) | (0.0021) | (0.0014) | (0.0075) | (0.0015) | (0.0086) |
| $S = 100$ | | | | | | |
| $n = 25$ | 0.1994 | 0.0997 | 0.0996 | 0.1968 | 0.0804 | 0.1588 |
| | (0.0099) | (0.0060) | (0.0042) | (0.0241) | (0.0047) | (0.0245) |
| $n = 50$ | 0.2006 | 0.1006 | 0.0997 | 0.2011 | 0.0798 | 0.1608 |
| | (0.0043) | (0.0029) | (0.0020) | (0.0099) | (0.0019) | (0.0112) |
| $n = 100$ | 0.2002 | 0.1002 | 0.0999 | 0.2001 | 0.0800 | 0.1609 |
| | (0.0025) | (0.0017) | (0.0011) | (0.0054) | (0.0011) | (0.0067) |
| Large | $\pi_1 = 0.2$ | $\pi_0 = 0.1$ | $p_0^{(1)} = 0.2$ | $p_1^{(1)} = 0.4$ | $p_0^{(2)} = 0.16$ | $p_1^{(2)} = 0.32$ |
| $S = 50$ | $\widehat{\pi}_1$ | $\widehat{\pi}_0$ | $\widehat{p}_0^{(1)}$ | $\widehat{p}_1^{(1)}$ | $\widehat{p}_0^{(2)}$ | $\widehat{p}_1^{(2)}$ |
| $n = 25$ | 0.2032 | 0.1039 | 0.1994 | 0.4012 | 0.1586 | 0.3191 |
| | (0.0370) | (0.0260) | (0.0092) | (0.0442) | (0.0112) | (0.0654) |
| $n = 50$ | 0.1987 | 0.0994 | 0.2005 | 0.3990 | 0.1602 | 0.3137 |
| | (0.0174) | (0.0122) | (0.0045) | (0.0224) | (0.0052) | (0.0330) |
| $n = 100$ | 0.2004 | 0.1006 | 0.1998 | 0.4004 | 0.1598 | 0.3206 |
| | (0.0084) | (0.0059) | (0.0023) | (0.0100) | (0.0025) | (0.0155) |
| $S = 100$ | | | | | | |
| $n = 25$ | 0.1987 | 0.0993 | 0.1995 | 0.3943 | 0.1604 | 0.3142 |
| | (0.0257) | (0.0173) | (0.0062) | (0.0322) | (0.0075) | (0.0452) |
| $n = 50$ | 0.2011 | 0.1012 | 0.1998 | 0.4013 | 0.1594 | 0.3189 |
| | (0.0123) | (0.0090) | (0.0032) | (0.0159) | (0.0039) | (0.0216) |
| $n = 100$ | 0.2004 | 0.1003 | 0.1999 | 0.4003 | 0.1599 | 0.3201 |
| | (0.0059) | (0.0042) | (0.0017) | (0.0073) | (0.0017) | (0.0112) |

Note: standard deviations based on 100 simulated samples are reported in parentheses.

Then, we compare five estimators based on three versions of 2SLS estimation: naive, adjusted, and oracle (infeasible). The naive 2SLS uses the noisy measure $H$ in place of the true network $G$, which means it uses $H_s y_s$ as an endogenous regressor and $H_s X_s$ as its instrument. The adjusted 2SLS estimator is what we propose in Section 4. It requires two steps. First, estimate the misclassification rates based on $(H^{(1)}, H^{(2)}, X)$. Then, construct $W_s^{(t)} = \frac{H_s^{(t)} - \widehat{p}_0^{(t)}(\iota_n \iota_n' - I_n)}{1 - \widehat{p}_0^{(t)} - \widehat{p}_1^{(t)}}$ for $t = 1, 2$, based on the first-step estimates $\widehat{p}_0^{(t)}$ and $\widehat{p}_1^{(t)}$, and apply 2SLS using $W_s^{(t)} y$ as an endogenous regressor and $W_s^{(t')} X$ as its instrument where $t \neq t'$. The oracle (infeasible) 2SLS uses the peer outcomes based on the actual network, i.e., $G_s y_s$, as an endogenous regressor, and uses $G_s X_s$ as its instrument.

### Table 1(b): Peer Effects Estimation: Small Misclassification

| | $S = 50$ | | | | | $S = 100$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Naive | | Adjusted | | Oracle | Naive | | Adjusted | | Oracle |
| Reg. | $H^{(1)}y$ | $H^{(2)}y$ | $W^{(1)}y$ | $W^{(2)}y$ | $Gy$ | $H^{(1)}y$ | $H^{(2)}y$ | $W^{(1)}y$ | $W^{(2)}y$ | $Gy$ |
| IV | $H^{(1)}X$ | $H^{(2)}X$ | $H^{(2)}X$ | $H^{(1)}X$ | $GX$ | $H^{(1)}X$ | $H^{(2)}X$ | $H^{(2)}X$ | $H^{(1)}X$ | $GX$ |
| $n = 25$ | Expected # of peers 3.75 | | | | | | | | | |
| $\lambda = 0.05$ | 0.0259 | 0.0307 | 0.0490 | 0.0467 | 0.0508 | 0.0283 | 0.0324 | 0.0517 | 0.0511 | 0.0489 |
| s.t.d | (0.007) | (0.006) | (0.012) | (0.014) | (0.005) | (0.005) | (0.005) | (0.008) | (0.009) | (0.007) |
| $\beta_1 = 1$ | 1.0613 | 1.0523 | 1.0113 | 1.0131 | 1.0108 | 1.0614 | 1.0540 | 1.0102 | 1.0117 | 1.0112 |
| s.t.d | (0.078) | (0.081) | (0.079) | (0.086) | (0.062) | (0.064) | (0.066) | (0.062) | (0.064) | (0.078) |
| $\beta_2 = 2$ | 1.9978 | 1.9983 | 1.9950 | 1.9951 | 2.0018 | 2.0064 | 2.0058 | 2.0041 | 2.0027 | 1.9946 |
| s.t.d | (0.046) | (0.046) | (0.047) | (0.047) | (0.031) | (0.032) | (0.032) | (0.034) | (0.032) | (0.046) |
| $n = 50$ | Expected # of peers 7.5 | | | | | | | | | |
| $\lambda = 0.05$ | 0.0274 | 0.0312 | 0.0492 | 0.0497 | 0.0499 | 0.0274 | 0.0310 | 0.0495 | 0.0493 | 0.0499 |
| s.t.d | (0.003) | (0.004) | (0.006) | (0.006) | (0.003) | (0.002) | (0.003) | (0.005) | (0.004) | (0.003) |
| $\beta_1 = 1$ | 1.1001 | 1.0836 | 1.0029 | 0.9971 | 1.0019 | 1.1021 | 1.0897 | 1.0010 | 1.0059 | 0.9988 |
| s.t.d | (0.068) | (0.064) | (0.067) | (0.060) | (0.043) | (0.047) | (0.047) | (0.047) | (0.046) | (0.060) |
| $\beta_2 = 2$ | 2.0036 | 2.0032 | 2.0021 | 2.0008 | 1.9991 | 2.0017 | 2.0013 | 1.9990 | 1.9983 | 2.0010 |
| s.t.d | (0.032) | (0.031) | (0.035) | (0.032) | (0.020) | (0.021) | (0.020) | (0.022) | (0.021) | (0.030) |
| $n = 100$ | Expected # of peers 15 | | | | | | | | | |
| $\lambda = 0.05$ | 0.0277 | 0.0313 | 0.0504 | 0.0504 | 0.0500 | 0.0278 | 0.0313 | 0.0503 | 0.0500 | 0.0501 |
| s.t.d | (0.001) | (0.001) | (0.003) | (0.003) | (0.001) | (0.001) | (0.001) | (0.002) | (0.002) | (0.001) |
| $\beta_1 = 1$ | 1.2544 | 1.2210 | 0.9984 | 1.0039 | 1.0060 | 1.2589 | 1.2197 | 1.0051 | 0.9999 | 1.0008 |
| s.t.d | (0.072) | (0.065) | (0.070) | (0.064) | (0.026) | (0.048) | (0.041) | (0.047) | (0.045) | (0.041) |
| $\beta_2 = 2$ | 2.0002 | 2.0004 | 1.9983 | 1.9988 | 1.9979 | 2.0017 | 2.0010 | 1.9983 | 1.9973 | 1.9993 |
| s.t.d | (0.026) | (0.022) | (0.035) | (0.028) | (0.013) | (0.019) | (0.017) | (0.023) | (0.019) | (0.020) |

34

Across the simulated samples indexed by $q = 1, 2, ..., Q$, we record the empirical distribution of these estimates of $(\lambda, \beta_1, \beta_2)$. Tables 1(b) and (c) report the average estimates, and sample s.t.d. based on this empirical distribution under different misclassification rates.

Results in Tables 1(b) and (c) demonstrate the following patterns. First, the naive method that ignores the misclassification in $H$ has serious bias in estimating the peer effects $\lambda = 0.05$. With lower misclassification rates, it estimates $\lambda$ at around 0.028 using $H^{(1)}$ and around 0.031 using $H^{(2)}$; with higher misclassification rates, it estimates $\lambda$ at around 0.013 using $H^{(1)}$ and around 0.018 using $H^{(2)}$. When estimating $\beta$, the naive estimation also shows bias, but not smaller than the bias in $\lambda$.

**Table 1(c): Peer Effects Estimation: Large Misclassification**

| Reg. IV | $S = 50$ | | | | | $S = 100$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Naive | | Adjusted | | Oracle | Naive | | Adjusted | | Oracle |
| | $H^{(1)}y$ $H^{(1)}X$ | $H^{(2)}y$ $H^{(2)}X$ | $W^{(1)}y$ $H^{(2)}X$ | $W^{(2)}y$ $H^{(1)}X$ | $Gy$ $GX$ | $H^{(1)}y$ $H^{(1)}X$ | $H^{(2)}y$ $H^{(2)}X$ | $W^{(1)}y$ $H^{(2)}X$ | $W^{(2)}y$ $H^{(1)}X$ | $Gy$ $GX$ |
| $n = 25$ | Expected # of peers 3.75 | | | | | | | | | |
| $\lambda = 0.05$ | 0.0118 | 0.0180 | 0.0460 | 0.0437 | 0.0489 | 0.0136 | 0.0195 | 0.0532 | 0.0500 | 0.0508 |
| s.t.d | (0.007) | (0.007) | (0.020) | (0.027) | (0.007) | (0.005) | (0.004) | (0.019) | (0.020) | (0.005) |
| $\beta_1 = 1$ | 1.0813 | 1.0733 | 1.0117 | 1.0173 | 1.0112 | 1.0822 | 1.0722 | 1.0005 | 1.0189 | 1.0108 |
| s.t.d | (0.081) | (0.081) | (0.101) | (0.095) | (0.078) | (0.068) | (0.068) | (0.085) | (0.078) | (0.062) |
| $\beta_2 = 2$ | 1.9967 | 1.9980 | 1.9951 | 1.9937 | 1.9946 | 2.0045 | 2.0059 | 2.0023 | 2.0027 | 2.0018 |
| s.t.d | (0.047) | (0.046) | (0.054) | (0.054) | (0.046) | (0.033) | (0.032) | (0.042) | (0.035) | (0.031) |
| $n = 50$ | Expected # of peers 7.5 | | | | | | | | | |
| $\lambda = 0.05$ | 0.0132 | 0.0188 | 0.0510 | 0.0510 | 0.0499 | 0.0133 | 0.0184 | 0.0491 | 0.0486 | 0.0499 |
| s.t.d | (0.003) | (0.003) | (0.014) | (0.020) | (0.003) | (0.002) | (0.002) | (0.009) | (0.011) | (0.003) |
| $\beta_1 = 1$ | 1.1431 | 1.1273 | 0.9942 | 0.9865 | 0.9988 | 1.1458 | 1.1348 | 0.9956 | 1.0111 | 1.0019 |
| s.t.d | (0.072) | (0.068) | (0.097) | (0.088) | (0.060) | (0.050) | (0.051) | (0.067) | (0.071) | (0.043) |
| $\beta_2 = 2$ | 2.0011 | 2.0027 | 1.9987 | 1.9995 | 2.0010 | 2.0000 | 2.0010 | 1.9967 | 1.9976 | 1.9991 |
| s.t.d | (0.030) | (0.031) | (0.046) | (0.036) | (0.030) | (0.022) | (0.021) | (0.030) | (0.022) | (0.017) |
| $n = 100$ | Expected # of peers 15 | | | | | | | | | |
| $\lambda = 0.05$ | 0.0133 | 0.0185 | 0.0504 | 0.0500 | 0.0501 | 0.0135 | 0.0185 | 0.0500 | 0.0506 | 0.0500 |
| s.t.d | (0.001) | (0.001) | (0.008) | (0.008) | (0.001) | (0.001) | (0.001) | (0.005) | (0.006) | (0.001) |
| $\beta_1 = 1$ | 1.3679 | 1.3357 | 0.9936 | 1.0079 | 1.0008 | 1.3726 | 1.3358 | 1.0079 | 0.9860 | 1.0060 |
| s.t.d | (0.092) | (0.086) | (0.136) | (0.115) | (0.041) | (0.060) | (0.055) | (0.096) | (0.087) | (0.026) |
| $\beta_2 = 2$ | 1.9983 | 1.9996 | 1.9982 | 1.9986 | 1.9993 | 2.0007 | 2.0015 | 1.9995 | 1.9988 | 1.9979 |
| s.t.d | (0.027) | (0.026) | (0.061) | (0.045) | (0.020) | (0.210) | (0.019) | (0.046) | (0.035) | (0.014) |

Second, our proposed adjusted 2SLS can estimate $(\lambda, \beta_1, \beta_2)$ with high accuracy. The average estimates are very close to the oracle estimates, albeit with larger standard devia-

tions. This is of course due to the noise from link misclassification as well as estimation errors in the initial estimates of the misclassification rates.

Third, when we fix the group size $n$, and increase the group number from $S = 50$ to $100$, the s.t.d. decreases by around $1/\sqrt{2}$, consistent with our theory of $\sqrt{S}$ asymptotics.

## 7. APPLICATION: MICROFINANCE PARTICIPATION IN INDIA

We apply our method to study how peer effects influence household decisions to participate in a microfinance program in India. The sample was collected by Banerjee et al. (2013) using survey questionnaires from the State of Karnataka, India between 2006-2007. Banerjee et al. (2013) impute a social network structure in the sample by aggregating several network measures that were inferred from the survey responses. They study how the dissemination of information about a microfinance program, Bharatha Swamukti Samsthe, or *BSS*, depended on the network position of the households that were the first to be informed about the program. Banerjee et al. (2013) use a binary response model with social interactions to disentangle the effect of information diffusion from the peer effects, a.k.a. *endorsement* effects. In contrast, we use two of the multiple measures in Banerjee et al. (2013) as noisy measures for an actual network, and apply our method to estimate peer effects in a linear social network model.

### 7.1. *Institutional background and data*

The sample was collected by Banerjee et al. (2013) through survey questionnaires from $S = 43$ villages in the State of Karnataka, India.[13] These villages are largely linguistically homogeneous but heterogeneous in terms of caste. The sample contains information about the socioeconomic status and some demographic characteristics of 9,598 households. On average, there were about 223 households in each village, with a minimum of 114, a maximum of 356, and a standard deviation of 56.2.

We merge the information from a full-scale household census and an individual-level survey in Banerjee et al. (2013). The household census gathered demographic information and data on a variety of amenities, such as roofing material, type of latrine, and quality of access to electric power. The individual survey was administered to a randomly selected

---

[13]The data are publicly available at: http://economics.mit.edu/faculty/eduflo/social.

sub-sample of villagers, which covered 46% of all households in the census. Individual questionnaires collected demographic information, such as age, caste and sub-caste, education, language, and having a ration card or not, but does not include explicit financial information. We merged the information about the head of household from the individual survey with the household information from the census. This yields a sample of 4,149 households.

**Table 2(a): Summary of Dependent and Explanatory Variables**

(Number of obs.: 4,149)

| Variable | definition | mean | s.d. | min | max |
|----------|------------|------|------|-----|-----|
| $y$ | dummy for participation | 0.1894 | 0.3919 | 0 | 1 |
| $room$ | number of rooms | 2.4389 | 1.3686 | 0 | 19 |
| $bed$ | number of beds | 0.9229 | 1.3840 | 0 | 24 |
| $age$ | age of household head | 46.057 | 11.734 | 20 | 95 |
| $edu$ | education of household head | 4.8383 | 4.5255 | 0 | 15 |
| $lang$ | whether to speak other language | 0.6799 | 0.4666 | 0 | 1 |
| $male$ | whether the hh head is male | 0.9161 | 0.2772 | 0 | 1 |
| $leader$ | whether it has a leader | 0.1393 | 0.3463 | 0 | 1 |
| $shg$ | whether in any saving group | 0.0513 | 0.2207 | 0 | 1 |
| $sav$ | whether to have a bank account | 0.3840 | 0.4864 | 0 | 1 |
| $election$ | whether to have an election card | 0.9525 | 0.2127 | 0 | 1 |
| $ration$ | whether to have a ration card | 0.9012 | 0.2985 | 0 | 1 |

Table 2(a) reports summary statistics for the dependent variable ($y = 1$ if participates in the microfinance program) as well as a few continuous and binary explanatory variables. Summary statistics for additional categorical variables, such as religion, caste, property ownership, access to electricity, etc, are reported in Table 2(b). The individual-level survey in Banerjee et al. (2013) also collected information of social interactions between households, including (i) individuals whose homes the respondent visited, and (ii) individuals who visited the respondent's home. Banerjee et al. (2013) construct graphs with undirected links by symmetrizing the data.[14] In other words, the sample in Banerjee et al. (2013) contains two symmetrized measures for the same latent network, based on the responses to (i) and (ii) respectively. These two measures, reported as "visitGo" and "visitCome" matrices

---

[14]Two households $i$ and $j$ are considered connected by an undirected link if an individual from either household mentioned the name of someone from the other household in response to question (i). Likewise, a second symmetric network measure is constructed based on responses to (ii).

in the sample and denoted as $H^{(1)}$ and $H^{(2)}$ in our notation, lend themselves to application of our method in Section 3.3.2.[15]

**Table 2(b): Summary of Category Variables**

| Variable | value | obs. | per. | Variable | value | obs. | per. |
|---|---|---|---|---|---|---|---|
| *religion* | | | | *latrine* | | | |
| - | Hinduism | 3943 | 95.04 | - | Owned | 1195 | 28.80 |
| - | Islam | 198 | 4.77 | - | Common | 20 | 0.48 |
| - | Christianity | 7 | 0.19 | - | None | 2934 | 70.72 |
| *roof* | | | | *property* | | | |
| - | Thatch | 82 | 1.98 | - | Owned | 3727 | 89.83 |
| - | Tile | 1388 | 33.45 | - | Owned & shared | 32 | 0.77 |
| - | Stone | 1172 | 28.25 | - | Rented | 390 | 9.40 |
| - | Sheet | 868 | 20.92 | | | | |
| - | RCC | 475 | 11.45 | | | | |
| - | Other | 164 | 3.95 | | | | |
| *electricity* | | | | *caste* | | | |
| - | No power | 243 | 5.86 | - | Scheduled caste | 1139 | 27.54 |
| - | Private | 2662 | 64.18 | - | Scheduled tribe | 221 | 5.34 |
| - | Government | 1243 | 29.97 | - | OBC | 2253 | 54.47 |
| | | | | - | General | 523 | 12.65 |

**Table 3: Degree Distribution in Two Network Measures**

| Degree | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $H^{(1)}$ | 2 | 21 | 110 | 227 | 357 | 505 | 526 | 546 | 506 | 379 | 269 |
| $H^{(2)}$ | 4 | 24 | 112 | 245 | 384 | 522 | 534 | 577 | 491 | 386 | 255 |

| Degree | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | $\geq 21$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $H^{(1)}$ | 224 | 145 | 90 | 74 | 54 | 33 | 27 | 15 | 9 | 6 | 24 |
| $H^{(2)}$ | 179 | 137 | 102 | 59 | 46 | 28 | 22 | 13 | 9 | 3 | 17 |

Table 3 reports the empirical distribution of the degrees of $H^{(1)}$ and $H^{(2)}$. As these measures are symmetric, there is no distinction between the degrees of in-bound or out-bound

---

[15]Banerjee et al. (2013) aggregate responses from 12 questions, including (i) and (ii), to construct a single symmetric network, which is considered as the actual adjacency matrix $G$, in the absence of link misclassification. In contrast, we take a different approach by interpreting responses to questions (1) and (2) as two noisy measures of a single, actual adjacency matrix.

links. We pool all households across 43 villages into a single, large network. There are no links between households from different villages in the sample, so the observed network structure is block-diagonal. Our estimator allows for the possibility that the unobserved true structure may include links between blocks, using our results from section 5.4.

Each column of Table 3 reports the number of households in $H^{(1)}$ and in $H^{(2)}$ that report the number of links given by the degree column heading. Table 3 shows large differences between the two matrices in the number of reported connections between households. If there were no misclassification of actual undirected links in these measures, we would expect the two matrices $H^{(1)}$ and $H^{(2)}$ to be identical, and therefore have the same degree distribution. The fact that they differ substantially is indicative of substantial link misclassification in the measures, possibly due to the respondents' recall errors, or differences in how they interpret the questions regarding visits.

### 7.2. *Empirical strategy for estimating peer effects*

We use the following specification for the adjusted feasible structural form:

$$y = \lambda W^{(t)} y + X\beta + villageFE + v^{(t)} \text{ for } t = 1, 2, \tag{17}$$

where $y$ is a binary variable indicating whether the household participated in the microfinance program (BSS), $X$ is a matrix of household characteristics, and $villageFE$ are village fixed effects. Definitions and summary statistics of regressors in $X$ are listed in Table 2. Note that (17) provides *two* different feasible structural forms (of the same actual structural model), corresponding to $t = 1, 2$ respectively.

To implement our adjusted 2SLS estimator, we define $\phi_{ij} \equiv \phi_{ij}(X) = 1$ if $i$ and $j$ have the same caste, and $0$ otherwise. Then, based on our two network matrices $H^{(1)}$ (visit-go) and $H^{(2)}$ (visit-come), we get the following estimates:

$$\widehat{\pi}_1 = E(G_{ij}|\phi_{ij} = 1) = 0.0357, \ \widehat{\pi}_0 = E(G_{ij}|\phi_{ij} = 0) = 0.0144,$$

$$\widehat{p}_0^{(1)} = \Pr\{H_{ij}^{(1)} = 1|G_{ij} = 0\} = 0.0020, \ \widehat{p}_1^{(1)} = \Pr\{H_{ij}^{(1)} = 0|G_{ij} = 1\} = 0.1425,$$

$$\widehat{p}_0^{(2)} = \Pr\{H_{ij}^{(2)} = 1|G_{ij} = 0\} = 0.0001, \ \widehat{p}_1^{(2)} = \Pr\{H_{ij}^{(2)} = 0|G_{ij} = 1\} = 0.1079.$$

Let $n_s$ be the group size of village $s$. We then construct the adjusted measures

$$W_s^{(t)} = \frac{H_s^{(t)} - \widehat{p}_0^{(t)}(\iota_{n_s}\iota'_{n_s} - I_{n_s})}{1 - \widehat{p}_0^{(t)} - \widehat{p}_1^{(t)}}, \text{ for } s = 1, 2, ..., S, \text{ and } t = 1, 2$$

and apply our adjusted 2SLS estimator. The estimation results are reported in Table 4, whose columns are defined as follows:

- OLS: regression of a simple, linear model that ignores network effects by setting $\lambda = 0$.
- (a): naive 2SLS that uses $H^{(1)}$ in place of the actual $G$, by including $H^{(1)}y$ as an endogenous regressor and using $H^{(1)}X$ as its instruments.
- (b): adjusted 2SLS for the structural form with $t = 1$ in (17), using $H^{(2)}X$ as instruments for the *adjusted* endogenous regressor $W^{(1)}y$.
- (c): naive 2SLS analogous to (a), only with $H^{(1)}$ replaced by $H^{(2)}$.
- (d): adjusted 2SLS analogous to (b), only swapping the roles of $H^{(1)}$ and $H^{(2)}$. That is, using $H^{(1)}X$ as instruments for $W^{(2)}y$ in (17) with $t = 2$.
- (e): S2SLS as defined in (15). This is a "combined" estimator that stacks the moments and associated IVs from both structural forms in (b) and (d).

In summary, columns (a) and (c) report estimators that a researcher would use if he or she ignored the issue of link misclassification, and treated either $H^{(1)}$ or $H^{(2)}$, respectively, as if it were the true adjacency matrix $G$, applying a standard 2SLS estimator in the literature. In contrast, columns (b), (d) and (e) report the adjusted 2SLS estimators we propose to remove the estimation bias due to link misclassification.[16] Column (e) combines the information used for the estimators in (b) and (d), and so is our preferred estimator.

## 7.3. *Empirical results*

Table 4 reports that our adjusted 2SLS estimates for the peer effect $\widehat{\lambda}$ are 0.0499 when using $W^{(1)}y$ in the structural form (column (b)), 0.0542 using $W^{(2)}y$ (column (d)), and 0.0515 using both measures and S2SLS (column (e)). Standard errors are clustered at the village level. These estimates are all significant at the 1% level, and the differences between

---

[16]We need two network measures in this particular context because the measures in the sample are symmetric. As we noted in Section 3.3.1, we can also apply the adjusted 2SLS when the sample reports a single *asymmetric* network measure.

them are small relative to their standard errors. These estimates imply the likelihood of a household to participate in the microfinance program is increased by about 5.15% when the household is linked to one more participating household on the network (note for this calculation that our model does not row-normalize the network measures). With the average participation rate being 18.9% in the sample, these estimates suggest that peer effects, a.k.a. "endorsement effects" in Banerjee et al. (2013), are substantial.

The signs of estimated marginal effects by individual or household characteristics are plausible. Column (e) suggests the head of household being a "leader" (e.g. a teacher, a leader of a self-help group, or a shopkeeper) increases the participation rate by around 3.8%. These households with "leaders" were the first ones to be informed about the program, and were asked to forward information about the microfinance program to other potentially interested villagers. These leaders had received first-hand, detailed information about the program from its administrator, which could be conducive to higher participation rates. Households with younger heads are more likely to participate, but the magnitude of this age effect is less substantial. Being 10 years younger increases the participation rate by 1.7%. Having a ration card increases the participation rate by around 4.2%. Compared to households using private electricity, households using government-supplied electricity have a 3.3% higher participation rate. These two factors indicate that, holding other factors equal, households in poorer economic conditions are more inclined to participate in the microfinance program.

Table 4 also shows that, if we had ignored the issue of misclassified links in network measures, and had done 2SLS using $H^{(t)}X$ as instruments for the (un-adjusted) endogenous peer outcomes $H^{(t)}y$, then the estimator would have been biased. In (a), where we use $H^{(1)}X$ as instruments for $H^{(1)}y$, the estimate for $\lambda$ is 0.0523. In comparison, in (b) where we correct for misclassified link bias by using $H^{(2)}X$ as instruments for $W^{(1)}y$, then the estimated $\lambda$ is 0.0499. The upward bias resulted from ignoring the misclassified links is about 4.8% (as 0.0523/0.0499=1.048). Likewise, in (c) where we erroneously use $H^{(2)}X$ as instruments for $H^{(2)}y$, we get an upward bias about 1.5% in the peer effect estimate compared with the correct estimate in (d) (as 0.0550/0.0542=1.015).

**Table 4: Adjusted Two-stage Least Square Estimates**

| | OLS | (a) | (b) | (c) | (d) | (e) |
|---|---|---|---|---|---|---|
| R.h.s. Endogeneity Instruments | | $H^{(1)}y$ $H^{(1)}X$ | $W^{(1)}y$ $H^{(2)}X$ | $H^{(2)}y$ $H^{(2)}X$ | $W^{(2)}y$ $H^{(1)}X$ | $W^{(t)}y$ Combined |
| $\widehat{\lambda}$ | | 0.0523*** | 0.0499*** | 0.0550*** | 0.0542*** | 0.0515*** |
| | | (0.0079) | (0.0086) | (0.0097) | (0.0082) | (0.0083) |
| $leader$ | 0.0515*** | 0.0371** | 0.0355** | 0.0414** | 0.0403** | 0.0379** |
| | (0.0175) | (0.0187) | (0.0188) | (0.0184) | (0.0184) | (0.0185) |
| $age$ | -0.0012*** | -0.0017*** | -0.0017*** | -0.0016*** | -0.0017*** | -0.0017*** |
| | (0.0005) | (0.0005) | (0.0005) | (0.0005) | (0.0005) | (0.0005) |
| $ration$ | 0.0502** | 0.0438** | 0.0430** | 0.0420** | 0.0412** | 0.0422** |
| | (0.0212) | (0.0201) | (0.0202) | (0.0195) | (0.0194) | (0.0198) |
| $electricity - gov$ | 0.0441** | 0.0338** | 0.0326** | 0.0349** | 0.0339** | 0.0333** |
| | (0.0152) | (0.0157) | (0.0158) | (0.0156) | (0.0155) | (0.0156) |
| $electricity - no$ | 0.0162 | 0.0226 | 0.0233 | 0.0240 | 0.0248 | 0.0240 |
| | (0.0275) | (0.0296) | (0.0296) | (0.0300) | (0.0298) | (0.0297) |
| $caste - tribe$ | -0.0411 | -0.0278 | -0.0263 | -0.0270 | -0.0255 | -0.0260 |
| | (0.0294) | (0.0309) | (0.0305) | (0.0301) | (0.0298) | (0.0301) |
| $caste - obc$ | -0.0822*** | -0.0505** | -0.0468** | -0.0472** | -0.0435*** | -0.0456*** |
| | (0.0163) | (0.0217) | (0.0214) | (0.0218) | (0.0210) | (0.0212) |
| $caste - gen$ | -0.1142*** | -0.0718*** | -0.0669*** | -0.0669*** | -0.0620** | -0.0650*** |
| | (0.0239) | (0.0238) | (0.0244) | (0.0244) | (0.0235) | (0.0241) |
| $religion - Islam$ | 0.1225*** | 0.0967*** | 0.0938*** | 0.0880*** | 0.0843*** | 0.0895*** |
| | (0.0332) | (0.0325) | (0.0325) | (0.0346) | (0.0349) | (0.0335) |
| $religion - Chri$ | 0.1569 | 0.1427 | 0.1410 | 0.1462 | 0.1450 | 0.1431 |
| | (0.1440) | (0.1295) | (0.1279) | (0.1310) | (0.1299) | (0.1287) |
| $Controls$ | √ | √ | √ | √ | √ | √ |
| $Village FE$ | √ | √ | √ | √ | √ | √ |
| $R^2$ | 0.0862 | 0.1339 | 0.1353 | 0.1356 | 0.1366 | 0.1358 |
| Obs | 4134 | 4134 | 4134 | 4134 | 4134 | 4134 |

Note: s.e. clustered at village level are in parentheses. ***, **, and * indicate 1%, 5% and 10% significant. Controls include $male, roof, room, bed, latrine, edu, lang, shg, sav, election, own$.

As explained in Section 3.2, the bias in (a) and (c) is due to the correlation between $H^{(t)}X$ and the composite errors $\varepsilon + \lambda[G - H^{(t)}]y$. The magnitude of this bias is determined in part by the misclassification rates $(p_0^{(t)}, p_1^{(t)})$, which affect the correlation between the composite errors and the traditional instruments $H^{(t)}X$ for endogenous peer outcomes

42

$H^{(t)}y$ in a naive 2SLS. This is evident from (7): if both $p_0$ and $p_1$ were close to zero, then the r.h.s. side of (7) would be almost reduced to $v$, which is mean independent from $X$ under Lemma 1. In that case, $H^{(3-t)}X$ and non-linear functions of $X$ would function as valid IVs for $H^{(t)}y$ even without making adjustments in $W^{(t)}$.

The fact that estimates in (a) and (c) are fairly close to those in (b), (d) and (e) indicate the impact of link misclassification on peer effects is low in this application. However, our Monte Carlo simulations sometimes showed much larger impacts from misclassification, which suggests that in other empirical environments, we may expect larger bias when misclassification of links is not accounted for in estimation. The method we propose in this paper provides an easy remedy for this issue.

We conclude this section with model validation results in Table 5, which shows how the predicted values of $E(y|X)$ fit with the sample data. The Probit and Logit models use the same set of regressors as in Table 4. We report the summary statistics of the fitted values $\widehat{E(y|X)}$ under different models. Columns (a) through (d) of Table 5 are the fitted values of the feasible structural models used in each of the corresponding columns in Table 4.

**Table 5: Model Validation: Predicted Microfinance Participation**

| $\widehat{E(y\|X)}$ | Probit | Logit | OLS | (a) | (b) | (c) | (d) | (e) |
|---|---|---|---|---|---|---|---|---|
| $mean$ | 0.1894 | 0.1894 | 0.1894 | 0.1894 | 0.1894 | 0.1894 | 0.1894 | 0.1894 |
| $s.t.d$ | 0.1176 | 0.1181 | 0.1151 | 0.1357 | 0.1403 | 0.1372 | 0.1416 | 0.1405 |
| min | 0.0103 | 0.0166 | -0.0953 | -0.1062 | -0.1107 | -0.1282 | -0.1316 | -0.1314 |
| max | 0.7490 | 0.7673 | 0.6895 | 0.7911 | 0.8159 | 0.7370 | 0.7615 | 0.8286 |
| $< 0$ | 0% | 0% | 2.95% | 4.96% | 5.32% | 5.06% | 5.56% | 5.41% |
| $I\{\widehat{E(y\|X)}> 0.5\}$ | | | | | | | | |
| underpred. (1 to 0) | 17.76% | 17.66% | 18.34% | 17.27% | 17.05% | 17.30% | 17.08% | 17.10% |
| overpred. (0 to 1) | 0.92% | 1.11% | 0.27% | 0.94% | 1.14% | 0.87% | 1.92% | 1.04% |
| correct | 81.33% | 81.23% | 81.40% | 81.79% | 81.81% | 81.83% | 81.91% | 81.86% |

In all but one of the models in Table 5, the sample mean of the predicted participation probability $\widehat{E(y|X)}$ is 0.1894, which is equal to the sample mean of $y$ in the 4,134 observations used in the regression. The standard deviation of the predicted participation probability varies across different models. Predictions of linear probability models (LPM), reported under the column of "OLS" and (a)-(e), are mostly within the unit interval $[0, 1]$.

LPM predictions are strictly less than 1 for all observations in the sample; only 2.95% to 5.56% of the households in the sample end up with negative LPM predictions. That is, about 95% all LPM predictions in the sample are indeed within the unit interval.

Based on $\widehat{E(y|X)}$, we use the indicator $1\{\widehat{E(y|X)} > 0.5\}$ to predict whether an individual participates in the microfinance program, and calculate prediction rates. Predictions in our linear social network models in columns (a)-(e) generally outperform the OLS, Probit and Logit models in terms of the percentage of correct predictions.

## 8. CONCLUSION

This paper proposes adjusted-2SLS estimators that consistently estimate structural parameters, including peer effects, in social networks when the links reported in a sample are subject to random misclassification errors. By adjusting the endogenous peer outcomes and applying new instruments constructed from noisy network measures, our estimators resolve the additional endogeneity issues caused by link misclassification. As an initial step of our method, we propose simple, closed-form estimators for the misclassification rates in the network measures.

We apply our method to analyze the peer (endorsement) effects in households' decisions to participate in a microfinance program in Indian villages, using the data collected by Banerjee et al. (2013). Consistent with our theoretical results, our empirical estimates show that ignoring the issue of misclassified links in 2SLS estimation of social network models leads to an upward bias of up to 5% in the estimates of peer effects. A Monte Carlo analysis shows that in other applications, the bias from failing to account for link misclassification can be much larger.

## REFERENCES

ADVANI, ARUN AND BANSI MALDE (2018): "Credibly identifying social effects: Accounting for network formation and measurement error," *Journal of Economic Surveys*, 32 (4), 1016–1044. [7]

AIGNER, DENNIS J ET AL. (1973): "Regression with a binary independent variable subject to errors of observation," *Journal of Econometrics*, 1 (1), 49–59. [7]

AUERBACH, ERIC (2022): "Identification and estimation of a partially linear regression model using network data," *Econometrica*, 90 (1), 347–365. [7]

BANERJEE, ABHIJIT, ARUN G CHANDRASEKHAR, ESTHER DUFLO, AND MATTHEW O JACKSON (2013): "The diffusion of microfinance," *Science*, 341 (6144), 1236498. [6, 17, 35, 36, 37, 40, 43]

44

BLUME, LAWRENCE E, WILLIAM A BROCK, STEVEN N DURLAUF, AND YANNIS M IOANNIDES (2011): "Identification of social interactions," in *Handbook of social economics*, Elsevier, vol. 1, 853–964. [2]

BOLLINGER, CHRISTOPHER R (1996): "Bounding mean regressions when a binary regressor is mismeasured," *Journal of Econometrics*, 73 (2), 387–399. [7, 10]

BONHOMME, STÉPHANE, KOEN JOCHMANS, AND JEAN-MARC ROBIN (2016): "Non-parametric estimation of finite mixtures from repeated measurements," *Journal of the Royal Statistical Society: Series B: Statistical Methodology*, 211–229. [22]

BOUCHER, VINCENT AND ARISTIDE HOUNDETOUNGAN (2020): *Estimating peer effects using partial network data*, Centre de recherche sur les risques les enjeux économiques et les politiques. [8, 12]

BRAMOULLÉ, YANN, HABIBA DJEBBARI, AND BERNARD FORTIN (2009): "Identification of peer effects through social networks," *Journal of econometrics*, 150 (1), 41–55. [2, 3, 17]

BUTTS, CARTER T (2003): "Network inference, error, and informant (in) accuracy: a Bayesian approach," *social networks*, 25 (2), 103–140. [7]

CHANDRASEKHAR, ARUN AND RANDALL LEWIS (2011): "Econometrics of sampled networks," *Unpublished manuscript, MIT.[422].* [7, 12]

CHEN, XIAOHONG, HAN HONG, AND ELIE TAMER (2005): "Measurement error models with auxiliary data," *The Review of Economic Studies*, 72 (2), 343–366. [22]

GRAHAM, BRYAN S (2020): "Network data," in *Handbook of Econometrics*, Elsevier, vol. 7, 111–218. [2]

GRIFFITH, ALAN (2022): "Name your friends, but only five? the importance of censoring in peer effects estimates using social network data," *Journal of Labor Economics*, 40 (4), 779–805. [8]

GRIFFITH, ALAN AND JUNGYOUN KIM (2023): "The Impact of Missing Links on Linear Reduced-form Network-Based Peer Effects Estimates," . [8]

HARDY, MORGAN, RACHEL M HEATH, WESLEY LEE, AND TYLER H MCCORMICK (2019): "Estimating spillovers using imprecisely measured networks," *arXiv preprint arXiv:1904.00136.* [7]

HAUSMAN, JERRY A, JASON ABREVAYA, AND FIONA M SCOTT-MORTON (1998): "Misclassification of the dependent variable in a discrete-response setting," *Journal of econometrics*, 87 (2), 239–269. [10, 22]

HU, YINGYAO (2008): "Identification and estimation of nonlinear models with misclassification error using instrumental variables: A general solution," *Journal of Econometrics*, 144 (1), 27–61. [5, 7, 22]

HU, YINGYAO AND ZHONGJIAN LIN (2018): "Misclassification and the Hidden Silent Rivalry," . [22]

HU, YINGYAO AND YUYA SASAKI (2017): "Identification of paired nonseparable measurement error models," *Econometric Theory*, 33 (4), 955–979. [22]

KELEJIAN, HARRY H AND INGMAR R PRUCHA (1998): "A generalized spatial two-stage least squares procedure for estimating a spatial autoregressive model with autoregressive disturbances," *The Journal of Real Estate Finance and Economics*, 17 (1), 99–121. [2]

KLEPPER, STEVEN (1988): "Bounding the effects of measurement error in regressions involving dichotomous variables," *Journal of Econometrics*, 37 (3), 343–359. [7]

LEE, LUNG-FEI (2007): "Identification and estimation of econometric models with group interactions, contextual factors and fixed effects," *Journal of Econometrics*, 140 (2), 333–374. [2]

LEWBEL, ARTHUR (2007): "Estimation of average treatment effects with misclassification," *Econometrica*, 75 (2), 537–551. [5, 7]

LEWBEL, ARTHUR, XI QU, AND XUN TANG (2023): "Ignoring Measurement Errors in Social Networks," *The Econometrics Journal*, utad028. [9]

LI, TONG (2002): "Robust and consistent estimation of nonlinear errors-in-variables models," *Journal of Econometrics*, 110 (1), 1–26. [22]

LIN, XU (2010): "Identifying peer effects in student academic achievement by spatial autoregressive models with group unobservables," *Journal of Labor Economics*, 28 (4), 825–860. [2]

LIU, XIAODONG (2013a): "Estimation of a local-aggregate network model with sampled networks," *Economics Letters*, 118 (1), 243–246. [7]

——— (2013b): "Estimation of a local-aggregate network model with sampled networks," *Economics Letters*, 118 (1), 243–246. [7]

MAHAJAN, APRAJIT (2006): "Identification and estimation of regression models with misclassification," *Econometrica*, 74 (3), 631–665. [5, 7]

MANSKI, CHARLES F (1993): "Identification of endogenous social effects: The reflection problem," *The review of economic studies*, 60 (3), 531–542. [17]

MOLINARI, FRANCESCA (2008): "Partial identification of probability distributions with misclassified data," *Journal of Econometrics*, 144 (1), 81–117. [7]

SHALIZI, COSMA ROHILLA AND ALESSANDRO RINALDO (2013): "Consistency under sampling of exponential random graph models," *Annals of statistics*, 41 (2), 508. [7]

# Online Appendix: Estimating Social Network Models with Link Misclassification

Arthur Lewbel, Xi Qu, and Xun Tang

June 11, 2024

## A. Proofs in Sections 3-4

### A1. Proofs in Sections 3.1-3.5

*Proof of Lemma 1.* Under (A3), we have $E(Gy|G,X) = E[GM(X\beta + \varepsilon)|G,X] = GMX\beta$, and $E(Wy|G,X) = E[WME(X\beta + \varepsilon|H,G,X)|G,X] = E(W|G,X)MX\beta$. Under (A1) and (A2), $E(W|G,X) = G$. It follows from the definition of $v$ in (6) that $E(v|G,X) = 0$. $\square$

*Proof of Proposition 2.* By (A1), (A2), (A4), conditional mean of $(i,j)$-th entry in $W^2$ is

$$
\begin{aligned}
E\left[(W^2)_{ij}|G,X\right] &= E\left(\sum_{k\neq i,j} W_{ik}W_{kj}\Big|G,X\right) = \sum_{k\neq i,j} E\left(W_{ik}W_{kj}|G,X\right) \\
&= \sum_{k\neq i,j} E\left(W_{ik}|G_{ik},X\right)E\left(W_{kj}|G_{kj},X\right) \\
&= \sum_{k\neq i,j} G_{ik}G_{kj} = \left(G^2\right)_{ij}. \tag{1}
\end{aligned}
$$

1

It then follows that

$$
\begin{aligned}
E[(W'X)'v|G,X] &= E(X'W\varepsilon|G,X) + \lambda E\left[X'W\left(G-W\right)y\middle|\,G,X\right] \\
&= \lambda E\left[X'W\left(G-W\right)MX\beta\middle|\,G,X\right] \\
&= \lambda X'\left[E(W|G,X)G - E(W^2|G,X)\right]MX\beta \\
&= \lambda X'\left(G^2 - G^2\right)MX\beta = 0,
\end{aligned}
\tag{2}
$$

where the first two equalities are due to (A3) and the reduced form of $y$, and the last due to (1) and the fact that $E\left[W|G,X\right] = G$ under (A1) and (A2).

Next, note that $H = (1 - p_0 - p_1)W + p_0(\iota\iota' - I)$. Hence

$$
E[(H'X)'v|G,X] = 0 + E\left\{X'p_0(\iota\iota' - I)v\middle|\,G,X\right\} = 0
$$

where the first equality is due to (2) and the second due to Lemma 1.  $\square$

As noted in Section 3.3.2, we can construct instruments from multiple *symmetrized* measures for $G$, denoted by $H^{(1)}$ and $H^{(2)}$. Suppose $H^{(1)}$ and $H^{(2)}$ both satisfy (A1), (A2), (A3), and are independent in the sense of (A4'). Let $W^{(t)}$ be defined for $t = 1, 2$ as in the text.

We can construct feasible structural forms as in (10) in the main text, and use $W^{(2)}X$ (or $H^{(2)}X$) as instruments for $v^{(1)}$. To see why, note that for all $i$ and $j$ (including the case with $i = j$):

$$
\begin{aligned}
E\left[(W^{(2)}W^{(1)})_{ij}|G,X\right] &= E\left(\sum_{k\neq i,j} W_{ik}^{(2)}W_{kj}^{(1)}\middle|\,G,X\right) \\
&= \sum_{k\neq i,j} E\left(W_{ik}^{(2)}W_{kj}^{(1)}\middle|\,G,X\right) = \sum_{k\neq i,j} E\left(W_{ik}^{(2)}\middle|\,G_{ik},X\right)E\left(W_{kj}^{(1)}\middle|\,G_{kj},X\right) \\
&= \sum_{k\neq i,j} G_{ik}G_{kj} = \left(G^2\right)_{ij}.
\end{aligned}
\tag{3}
$$

Besides, under (A1) and (A2),

$$
E\left(W^{(2)}G|G,X\right) = E(W^{(2)}|G,X)G = G^2.
\tag{4}
$$

2

It then follows that

$$
\begin{aligned}
E[(W^{(2)}X)'v^{(1)}|G,X] &= E(X'W^{(2)}\varepsilon|G,X) + \lambda E\left\{ X'W^{(2)}\left[G - W^{(1)}\right]y\big|\,G,X\right\} \\
&= \lambda E\left[X'W^{(2)}\left(G - W^{(1)}\right)MX\beta\big|\,G,X\right] \\
&= \lambda X'\left[E(W^{(2)}G|G,X) - E(W^{(2)}W^{(1)}|G,X)\right]MX\beta = 0,
\end{aligned}
$$

where the first two equalities are due to (A3), and the last holds because of (3) and (4) under (A1), (A2), and (A4'). Next, by an argument similar to the proof of Proposition 2, $E[(W^{(2)}X)'v^{(1)}|G,X] = 0$ implies $E[(H^{(2)}X)'v^{(1)}|G,X] = 0$.

*Proof of Proposition 3.* Define some $K$-by-$K$ moments involving $(G,X)$:

$$
\begin{aligned}
B_1 &\equiv E(X'G^2MX),\ B_2 \equiv E(X'GMX),\ B_3 \equiv E(X'G^2X), \\
B_4 &\equiv E(X'GX),\ B_5 \equiv E(X'X).
\end{aligned}
$$

Recall $Z \equiv (W'X, X)$ and $R \equiv (Wy, X)$. Under (A1), (A2), (A3), and (A4),

$$
\begin{aligned}
E(Z'R) &= \begin{pmatrix} E(X'W^2y) & E(X'WX) \\ E(X'Wy) & E(X'X) \end{pmatrix} = \begin{pmatrix} E[X'W^2M(X\beta+\varepsilon)] & E(X'\mathcal{H}X) \\ E[X'WM(X\beta+\varepsilon)] & E(X'X) \end{pmatrix} \\
&= \begin{pmatrix} E(X'G^2MX\beta) & E(X'GX) \\ E(X'GMX\beta) & E(X'X) \end{pmatrix} \equiv \begin{pmatrix} B_1\beta & B_4 \\ B_2\beta & B_5 \end{pmatrix}.
\end{aligned}
$$

Suppose the $2K$-by-$(1+K)$ matrix $E(Z'R)$ does not have full rank. By definition this implies the $2K$-by-$2K$ square matrix

$$
\begin{pmatrix} B_1 & B_4 \\ B_2 & B_5 \end{pmatrix} \tag{5}
$$

must be singular. It then follows that non-singularity of the square matrix in (5) implies $E(Z'R)$ has full rank.

As $M - \lambda GM = I$, we have $GM = \lambda^{-1}(M - I)$ and $G^2M = \lambda^{-1}(GM - G) = \lambda^{-2}(M -$

3

$I - \lambda G$). We can write

$$
\begin{pmatrix} B_1 & B_4 \\ B_2 & B_5 \end{pmatrix} = \begin{pmatrix} \lambda^{-1} E[X'(GM - G)X] & E(X'GX) \\ E(X'GMX) & E(X'X) \end{pmatrix}.
$$

Adding the product of the 2nd row and $\left(-\frac{1}{\lambda}\right)$ to the 1st row, we get:

$$
\begin{pmatrix} -\frac{1}{\lambda} E(X'GX) & E(X'GX) - \frac{1}{\lambda} E(X'X) \\ E(X'GMX) & E(X'X) \end{pmatrix}.
$$

Adding the product of the 2nd column and $\left(\frac{1}{\lambda}\right)$ to the 1st column, we get

$$
\begin{pmatrix} -\frac{1}{\lambda^2} E(X'X) & E(X'GX) - \frac{1}{\lambda} E(X'X) \\ E(X'(GM + \frac{1}{\lambda}I)X) & E(X'X) \end{pmatrix} = \begin{pmatrix} -\frac{1}{\lambda^2} E(X'X) & -\frac{1}{\lambda} E(X'M^{-1}X) \\ \frac{1}{\lambda} E(X'MX) & E(X'X) \end{pmatrix}.
$$

The determinant of the matrix on the right-hand side is the product of $\lambda^{-2K}$ and the determinant of $[E(X'X), E(X'M^{-1}X); E(X'MX), E(X'X)]$. Hence, the matrix in (5) is non-singular iff $[E(X'X), E(X'M^{-1}X); E(X'MX), E(X'X)]$ is non-singular.

By the same token, (A1), (A2), and (A4) imply that

$$
E(Z'Z) = \begin{pmatrix} E(X'W^2X) & E(X'WX) \\ E(X'WX) & E(X'X) \end{pmatrix} = \begin{pmatrix} E(X'G^2X) & E(X'GX) \\ E(X'GX) & E(X'X) \end{pmatrix} = \begin{pmatrix} B_3 & B_4 \\ B_4 & B_5 \end{pmatrix}.
$$

Therefore, $E(Z'Z)$ has full rank if and only if $[B_3, B_4; B_4, B_5]$ is non-singular. $\qquad \square$


## A2. Identifying misclassification rates in Section 3.4

Consider the case in Section 3.4.1 where the sample reports two measures $H^{(1)}$ and $H^{(2)}$ with misclassification rates $p_0^{(t)}, p_1^{(t)}$ for $t = 1, 2$ respectively. Assume these two measures satisfy (A1), (A2), (A3), and (A4'). It is convenient to introduce a third measure whose

distribution is determined by the joint distribution of $H_{ij}^{(1)}$ and $H_{ij}^{(2)}$:

$$H_{ij}^{(3)} \equiv \max \left\{ H_{ij}^{(1)}, H_{ij}^{(2)} \right\}.$$

By construction, for $t = 1, 2, 3$, the distribution of $H_{ij}^{(t)}$ is related to $p_0^{(t)}, p_1^{(t)}$ and link formation probability $\pi_1 \equiv E(G_{ij}|\phi_{ij}(X) = 1)$ as follows:

$$\psi_1^{(t)} \equiv E \left[ H_{ij}^{(t)} \middle| \phi_{ij}(X) = 1 \right] = \left( 1 - p_1^{(t)} \right) \pi_1 + p_0^{(t)}(1 - \pi_1) = p_0^{(t)} + \left( 1 - p_1^{(t)} - p_0^{(t)} \right) \pi_1, \quad (6)$$

where (A4') implies:

$$p_0^{(3)} = p_0^{(1)} + p_0^{(2)} - p_0^{(1)} p_0^{(2)}, \tag{7}$$
$$p_1^{(3)} = p_1^{(1)} p_1^{(2)}. \tag{8}$$

Equations similar to (6) hold with "$\phi_{ij}(X) = 1$" and $\pi_1$ replaced by "$\phi_{ij}(X) = 0$" and $\pi_0$ respectively, thus defining $\psi_0^{(t)}$ accordingly.

[**Identifying $p_0^{(1)}$ and $p_0^{(2)}$.**] For convenience, let $\xi_1 \equiv \left( 1 - p_0^{(2)} - p_1^{(2)} \right) \pi_1$ so that

$$\psi_1^{(1)} = p_0^{(1)} + r_{(12)}\xi_1; \quad \psi_1^{(2)} = p_0^{(2)} + \xi_1;$$
$$\psi_1^{(3)} = p_0^{(1)} + p_0^{(2)} - p_0^{(1)} p_0^{(2)} + r_{(32)}\xi_1, \tag{9}$$

where $r_{(t't)} \equiv (\psi_0^{(t')} - \psi_1^{(t')})/(\psi_0^{(t)} - \psi_1^{(t)})$ for $t', t \in \{1, 2, 3\}$. This implies

$$p_0^{(1)} = \psi_1^{(1)} - r_{(12)}\xi_1 \text{ and } p_0^{(2)} = \psi_1^{(2)} - \xi_1. \tag{10}$$

Substituting these into the expression for $\psi_1^{(3)}$ in (9) implies:

$$\psi_1^{(3)} = \left( \psi_1^{(1)} - r_{(12)}\xi_1 - 1 \right) \left( 1 - \psi_1^{(2)} + \xi_1 \right) + 1 + r_{(32)}\xi_1.$$

5

Rearranging terms, we write this quadratic equation in $\xi_1$ as

$$\mathcal{C}_2\xi_1^2 - \mathcal{C}_1\xi_1 - \mathcal{C}_0 = 0, \tag{11}$$

where

$$\mathcal{C}_2 \equiv r_{(12)}, \tag{12}$$
$$\mathcal{C}_1 \equiv \psi_1^{(1)} - 1 + r_{(32)} - r_{(12)}(1 - \psi_1^{(2)}),$$
$$\mathcal{C}_0 \equiv \psi_1^{(1)} + \psi_1^{(2)} - \psi_1^{(1)}\psi_1^{(2)} - \psi_1^{(3)}.$$

By definition, $\mathcal{C}_2 = [1 - p_0^{(1)} - p_1^{(1)}]/[1 - p_0^{(2)} - p_1^{(2)}] > 0$, and

$$\mathcal{C}_0 = \left[1 - p_0^{(1)} - p_1^{(1)}\right]\left[1 - p_0^{(2)} - p_1^{(2)}\right]\pi_1(1 - \pi_1) > 0.$$

Hence $\Delta \equiv (\mathcal{C}_1)^2 + 4\mathcal{C}_2\mathcal{C}_0 > 0$ and $\sqrt{\Delta} > \mathcal{C}_1$. Then (11) admits two solutions in $\xi_1$:

$$\xi_1 = \frac{1}{2\mathcal{C}_2}(\mathcal{C}_1 \pm \sqrt{\Delta}).$$

However, $\xi_a \in (0, 1)$ by definition. Since $\frac{1}{2\mathcal{C}_2}\left(\mathcal{C}_1 - \sqrt{\Delta}\right) < 0$, the only solution in (11) must be $\xi_1 = \frac{1}{2\mathcal{C}_2}\left(\mathcal{C}_1 + \sqrt{\Delta}\right)$. Plugging in this solution of $\xi_1$ into (10) identifies $p_0^{(1)}$ and $p_0^{(2)}$.

[**Identifying $\pi_1$.**] Note that (6) implies

$$p_1^{(t)} = 1 - p_0^{(t)} - \frac{\psi_1^{(t)} - p_0^{(t)}}{\pi_1} \text{ for } t = 1, 2, 3. \tag{13}$$

Plugging in (13) into (8) and using (7), we get

$$\pi_1 = \frac{\left(\psi_1^{(1)} - p_0^{(1)}\right)\left(\psi_1^{(2)} - p_0^{(2)}\right)}{\left(1 - p_0^{(1)}\right)\left(\psi_1^{(2)} - p_0^{(2)}\right) + \left(1 - p_0^{(2)}\right)\left(\psi_1^{(1)} - p_0^{(1)}\right) - \left(\psi_1^{(3)} - p_0^{(3)}\right)}. \tag{14}$$

Recall that $\psi_1^{(t)}$ for $t = 1, 2, 3$ are directly identified from the data. With $p_0^{(t)}$ identified for $t = 1, 2$, we can recover $p_0^{(3)}$ from (7). This implies $\pi_1$ is identified from (14).

**[Identifying $p_1^{(1)}$, $p_1^{(2)}$ and $\pi_0$.]** With $p_0^{(1)}, p_0^{(2)}$, and $\pi_1$ identified above, we can use (13) to recover $p_1^{(t)}$ from $\psi_1^{(t)}$ for $t = 1, 2$. It is worth mentioning that these parameters $\pi_1$ and $p_0^{(t)}, p_1^{(t)}$ are over-identified because the argument above can also be applied to $\psi_0^{(t)}$ instead of $\psi_1^{(t)}$. The final step is to use to definition in (6) to (over-)identify $\pi_0$ as:

$$\pi_0 = \frac{\psi_0^{(t)} - p_0^{(t)}}{\psi_1^{(t)} - p_0^{(t)}} \pi_1 \text{ for } t = 1, 2, 3.$$

**[Single, unsymmetrized measure.]** The same identification argument applies for the case in Section 3.4.2, in which the sample reports a single, unsymmetrized measure $H$ with misclassification rates $p_1, p_0$ when the actual $G$ is known to be symmetric. For each *unordered* pair $\{i, j\}$, define $H_{\{i,j\}}^{(1)} \equiv H_{ij}$, $H_{\{i,j\}}^{(2)} \equiv H_{ji}$, and $H_{\{i,j\}}^{(3)} \equiv \max\{H_{ij}, H_{ji}\}$. There exists a system analogous to (6), with $H_{ij}^{(t)}$ replaced by $H_{\{i,j\}}^{(t)}$. However, in this case, the first two equations for $t = 1, 2$ coincide with each other, as $p_d^{(1)} = p_d^{(2)} = p_d$ for $d \in \{0, 1\}$ by construction. The remaining steps for identification are identical to the case above with two measures $H^{(1)}$ and $H^{(2)}$, except that the closed-form expressions are further simplified due to $r_{(12)} = 1$, $\psi_1^{(1)} = \psi_1^{(2)}$, and $p_d^{(1)} = p_d^{(2)}$ for $d \in \{0, 1\}$.

## A3. Asymptotic property of the adjusted 2SLS estimator

We derive the limiting distribution of our adjusted 2SLS estimator for the structural effects $\hat{\lambda}$ and $\hat{\beta}$ in Proposition 4 of Section 4.2.

Recall from Section 4.1 that we have defined for each group $s$,

$$v_{1s,1} \equiv \frac{2}{n_s(n_s - 1)} \sum_{i > j} H_{s,\{i,j\}} 1\{\phi_{s,\{i,j\}} = 1\},$$

$$v_{2s,1} \equiv \frac{2}{n_s(n_s - 1)} \sum_{i > j} 1\{\phi_{s,\{i,j\}} = 1\},$$

and defined $v_{1s,0}$, $v_{2s,0}$ analogously by replacing $\phi_{s,\{i,j\}} = 1$ with $\phi_{s,\{i,j\}} = 0$. Let $v_s \equiv (v_{1s,1}, v_{2s,1}, v_{1s,0}, v_{2s.0})'$. We maintain the following regularity conditions:

(*REG*) (i) $\exists \delta > 0$ s.t. $\lim_{S \to \infty} \sum_{s=1}^{S} E\left\{\|Z_s' R_s(p)\|^{1+\delta}\right\} / (1 + \delta) < \infty$; similar conditions

hold for $Z_s'Z_s$ and $Z_s'\nabla [R_s(p)\theta]$. (ii) Let $\tau_s$, $\zeta_s$ be defined as in (15) and (17) below. $\exists \delta' > 0$ s.t. $E(||\tau_s||^{2+\delta'}) < \infty$, and $S \times Var\left[S^{-1}\left(\sum_{s=1}^{S}\tau_s\right)\right] > 0$ is bounded away from zero by some positive constants for $S$ large enough; similar conditions hold for $\zeta_s$.

Under these conditions, we can apply appropriate versions of the law of large numbers, the central limit theorem, and the delta method to our sample which consists of observations $y_s, X_s, H_s$ that are independent and potentially heterogeneously distributed (due to the variation in group sizes $n_s$).

First, note our estimator for misclassification rates $\hat{p}$ is a closed-form function of the sample averages of $v_s$. Thus the asymptotic linear presentation of $\hat{p}$ is

$$\sqrt{S}(\hat{p} - p) = \frac{1}{\sqrt{S}}\sum_s \underbrace{\mathcal{J}_0 \times [v_s - E(v_s)]}_{\equiv \tau_s} + o_p(1), \tag{15}$$

where $\mathcal{J}_0$ depends on the Jacobian matrix of $\hat{p}$ w.r.t. the sample averages of $v_s$, evaluated at population counterparts.

Next, note that by construction,

$$\begin{aligned}\sqrt{S}\left(\hat{\theta} - \theta\right) &= \sqrt{S}\left(\mathbf{A}'\mathbf{B}^{-1}\mathbf{A}\right)^{-1}\mathbf{A}'\mathbf{B}^{-1}\mathbf{Z}'\left[Y - \mathbf{R}(\hat{p})\theta\right] \\ &= \left(A_0'B_0^{-1}A_0\right)^{-1}A_0'B_0^{-1}\frac{1}{\sqrt{S}}\mathbf{Z}'\left[Y - \mathbf{R}(\hat{p})\theta\right] + o_p(1), \tag{16}\end{aligned}$$

where the second equality holds as $\mathbf{A}/S \xrightarrow{p} A_0$, $\mathbf{B}/S \xrightarrow{p} B_0$, $S^{-1/2}\mathbf{Z}'\left[Y - \mathbf{R}(\hat{p})\theta\right] = O_p(1)$.

Recall the following definitions from the text:

$$F_0 \equiv \lim_{S \to \infty} S^{-1}\sum_{s=1}^{S} E\left\{Z_s'\nabla[R_s(p)\theta]\right\}.$$

For each group $s$ and individual $i \le n_s$, let $R_{s,i}(p)$ denote the corresponding row in $\mathbf{R}(p)$, and $\nabla_p R_{s,i}(p)$ be the $(K+1)$-by-2 Jacobian matrix of $R_{s,i}(p)$ with respect to $p$. Let $\nabla_p[R_s(p)\theta]$ denote an $n_s$-by-2 matrix with each row $i$ being $\theta'\nabla_p R_{s,i}(p)$, and let $\nabla_p[\mathbf{R}(p)\theta]$

be an $N$-by-2 matrix that stacks them for $s \leq S$. Then,

$$
\begin{aligned}
\tfrac{1}{\sqrt{S}}\mathbf{Z}'\left[Y - \mathbf{R}(\widehat{p})\theta\right] &= \tfrac{1}{\sqrt{S}}\mathbf{Z}'\left[Y - \mathbf{R}(p)\theta\right] - \left(\tfrac{1}{S}\mathbf{Z}'\nabla_p\left[\mathbf{R}(p)\theta\right]\right)\sqrt{S}(\widehat{p} - p) + o_p(1) \\
&= \tfrac{1}{\sqrt{S}}\sum_s Z_s'\left[y_s - R_s(p)\theta\right] - F_0\left(\tfrac{1}{\sqrt{S}}\sum_s \tau_s\right) + o_p(1) \\
&= \tfrac{1}{\sqrt{S}}\sum_s \underbrace{Z_s'v_s - F_0\tau_s}_{\equiv \zeta_s} + o_p(1). \quad\quad (17)
\end{aligned}
$$

The first equality follows form a Taylor approximation around the actual misclassification rates $p = (p_0, p_1)'$; the second from $\tfrac{1}{S}\mathbf{Z}'\nabla_p\left[\mathbf{R}(p)\theta\right] \overset{p}{\longrightarrow} \lim_{S\to\infty} S^{-1}\sum_s E\left\{Z_s'\nabla_p\left[R_s(p)\theta\right]\right\}$ and from the asymptotic linear representation of the estimator $\widehat{p} = (\hat{p}_0, \hat{p}_1)$; the third from $y_s = R_s(p)\theta + v_s$. This proves the limiting distribution of $\sqrt{S}(\widehat{\theta} - \theta)$ in Proposition 4.

# B. Proofs and Further Details for Section 5

## B1. Proofs in Section 5.1

*Proof of Proposition 5.* Under (A3),

$$
\begin{aligned}
E(Gy|X, G) &= E[GM(X\beta + GX\gamma + \varepsilon)|X, G] = GM\left(X\beta + GX\gamma\right), \\
E(Wy|X, G) &= E[WME\left(X\beta + GX\gamma + \varepsilon|X, G, H\right)|X, G] = E(W|G, X)M(X\beta + GX\gamma).
\end{aligned}
$$

Under (A1) and (A2), $E(W|G, X) = G$. It then follows that $E(\eta|X, G) = 0$.

Next, note

$$
\begin{aligned}
E\left[\zeta(X)'WWy|G, X\right] &= \zeta(X)'E(W^2|G, X)M(X\beta + GX\gamma); \\
E[\zeta(X)'WWX|G, X] &= \zeta(X)'E(W^2|G, X)X; \\
E\left[\zeta(X)'WGy|G, X\right] &= \zeta(X)'E(W|G, X)GM(X\beta + GX\gamma); \\
E[\zeta(X)'WGX|G, X] &= \zeta(X)'E(W|G, X)GX.
\end{aligned}
$$

As shown in (1), under (A4), $E(W^2|G, X) = G^2$. Because $E(W|G, X) = G$ under (A1) and (A2), this implies $E\left[\zeta(X)'W\eta\right] = 0$. Also, $E[\zeta(X)'H\eta] = (1 - p_0 - p_1)E[\zeta(X)'W\eta] +$

9

$E[\zeta(X)'p_0(\iota\iota' - I)\eta] = 0$, where the second equality holds because $E(\eta|X, G) = 0$. $\qquad\square$

## B2. The Setting of a single large network

In the main text, we focus on cases where the sample consists of many small, fixed-sized groups, where no links exist between members of different groups.

We now show how the idea of an adjusted 2SLS also applies when there is interdependence between *all* individuals in a sample. Specifically, we consider a setting in which the sample is partitioned into well-defined, *approximate groups*, which we henceforth refer to as "*blocks*". Formally, the individuals in the sample are partitioned into $S$ blocks. Links within each block $s$ are dense (i.e., the probability of forming links between individuals within the same block does *not* diminish as the sample size increases); links between different blocks are sparse, with the rate of formation diminishing as the number of blocks increases.

The sample size is $N \equiv \sum_{s=1}^{S} n_s$. Let $G_N$ and $H_N$ denote the true and noisy measure of $N$-by-$N$ adjacency matrices respectively, which span the $S$ blocks in the sample. Link misclassification exists in $H_N$ in two ways. First, links within each block are randomly misclassified in the sample at rate $p_0$ and $p_1$. Second, sparse cross-block links are *never* reported in the sample. By definition, $H_N$ is block-diagonal, with each diagonal block indexed as $H_{N,s}$ for $s = 1, 2, ..., S$.

To facilitate derivation of the asymptotic properties of our 2SLS estimator, let $\widetilde{G}_N$ be a hypothetical *block-diagonal approximation* of $G_N$, which perfectly reports all within-block links but drops all cross-block links. That is, for all individual $i$,

$$\widetilde{G}_{N,ij} = G_{N,ij} \text{ if } j \in s(i); \ \widetilde{G}_{N,ij} = 0 \text{ if } j \notin s(i),$$

where $s(i)$ indicates the block that $i$ belongs to. By construction, all elements outside the diagonal blocks in $\widetilde{G}_N$ are zeros. We maintain the following assumptions on the measurement errors in $H_N$:

$$(N1) \ E(H_{N,ij}|\widetilde{G}_N, X) = E(H_{N,ij}|\widetilde{G}_{N,ij}, X) \ \forall i \neq j;$$

(N2) $E(H_{N,ij}|\widetilde{G}_{N,ij} = 1, X) = 1 - p_1$, $E(H_{N,ij}|\widetilde{G}_{N,ij} = 0, X) = p_0$ $\forall i$ and $j \neq i$ in $s(i)$.

As before, assume $p_0 + p_1 < 1$. Furthermore, we maintain that the block-specific random arrays, $H_{N,s}$, $\widetilde{G}_{N,s}$, $X_{N,s}$, $\epsilon_{N,s}$ (with $H_{N,s}$, $\widetilde{G}_{N,s}$ being $n_s$-by-$n_s$ matrices), are drawn independently across the blocks. Under these maintained conditions, we can consistently estimate the misclassification rates following the same approach as in Section 4.1 and using linked pairs within diagonal blocks only. For the rest of this section, we take $(p_0, p_1)$ as given, and focus on the asymptotic properties of an adjusted 2SLS that removes misclassification bias by adjusting the diagonal block measures.

Let $W_N$ be a block-diagonal matrix, with each of its $S$ diagonal blocks adjusted as $W_{N,s} \equiv [H_{N,s} - p_0(\iota_{n_s}\iota'_{n_s} - I_{n_s})]/(1 - p_0 - p_1)$. In the Web Appendix, we show that the structural model

$$y_N = \lambda G_N y_N + X_N \beta + \varepsilon_N$$

can be written as

$$y_N = \lambda W_N y_N + X_N \beta + v_N + u_N, \tag{18}$$

where $u_N \equiv (I_N - \lambda W_N)\left(I_N - \lambda \widetilde{G}_N\right)^{-1} \lambda \Delta_N y_N$ with $\Delta_N \equiv G_N - \tilde{G}_N$ and

$$v_N \equiv \varepsilon_N + \lambda \left(\widetilde{G}_N - W_N\right) \widetilde{y}_N \text{ with } \widetilde{y}_N \equiv (I_N - \lambda \widetilde{G}_N)^{-1}(X_N \beta + \varepsilon_N).$$

Note that we decompose composite errors in (18) into $u_N$ and $v_N$, which are both vectorizations of block-specific vectors $u_{N,s}$ and $v_{N,s}$. While $v_{N,s}$ are independent across the blocks, $u_{N,s}$ are correlated across the blocks because of interdependence between $y_{N,s}$ due to sparse links between the blocks in $G_N$. This difference requires us to apply separate tactics to characterize their contribution to the estimation errors.

This decomposition of the composite error is useful for illustrating two main steps for deriving the asymptotic result. Let $Z_N$ denote the matrix of instruments, with $Z_{N,s}$ being its sub-matrix specific to block $s$. Instrument exogeneity requires $E(Z'_{N,s}v_{N,s}) = 0$ for all $s$. Recall the 2SLS estimator that uses $Z_N$ as instruments for $R_N \equiv (W_N y_N, X_N)$ is

$\widehat{\theta} = \left(A_N' B_N^{-1} A_N\right)^{-1} A_N' B_N^{-1} Z_N' y_N$, where $A_N \equiv Z_N' R_N$ and $B_N \equiv Z_N' Z_N$. By definition,

$$\widehat{\theta} - \theta = \left(A_N' B_N^{-1} A_N\right)^{-1} A_N' B_N^{-1} Z_N' (v_N + u_N).$$

The asymptotic property of the estimator thus depends on that of $Z_N' v_N$ and $Z_N' u_N$, which we investigate sequentially.

First, we characterize the order of $Z_N' v_N$, using the fact that $v_{N,s}$ are independent across blocks $s$. To see why such independence holds, recall that $H_{N,s}$, $\widetilde{G}_{N,s}$, $X_{N,s}$, $\epsilon_{N,s}$ are assumed independent across blocks $s$. By construct, $\widetilde{G}_N$, $H_N$, $W_N$ and $(I - \lambda \widetilde{G}_N)^{-1}$ are all block-diagonal. Hence $\widetilde{y}_{N,s} = (I_s - \lambda \widetilde{G}_{N,s})^{-1}(X_{N,s}\beta + \varepsilon_{N,s})$ are independent across $s$.[1] It then follows that $v_{N,s} = \varepsilon_{N,s} + \lambda \left(\widetilde{G}_{N,s} - W_{N,s}\right) \widetilde{y}_{N,s}$, and are independent across $s$.

We maintain exogeneity and independence conditions which are analogous to (A3) and (A4) for the case with small groups in Section 3:

(N3)   $E(\varepsilon_{N,s} | X_{N,s}, G_{N,s}, H_{N,s}) = 0$ for all $s$;

(N4)   Conditional on $(G_N, X_N)$, $H_{N,ij} \perp H_{N,kl}$ for all $(i,j) \neq (k,l)$.

Under these conditions, $E(v_{N,s} | X_{N,s}, G_{N,s}) = 0$. The independence between $v_{N,s}$ mentioned above then allows us to apply the law of large numbers to show that

$$\frac{1}{S} Z_N' v_N = \frac{1}{S} \sum_s Z_{N,s}' v_{N,s} = O_p(S^{-1/2}).$$

Second, the order of $\frac{1}{S} Z_N' u_N$ is bounded above by the expected number of misclassified links across the blocks, which are assumed to be sparse in the following sense:

(S-LOB) $\sum_{i=1}^N \sum_{j \notin s(i)} E\left(|\Delta_{N,ij}|\right) = O(S^\rho)$ for some $\rho < 1$.

This condition is the same as in Lewbel et al. (2023), who provide examples with primitive conditions. Among other things, it requires the links outside these blocks, or approximate

---

[1] We refer to $\widetilde{y}_N$ as a *hypothetical* reduced form, because it is based on the block-diagonal approximation $\widetilde{G}_N$ rather than the actual $G_N$.

groups, to be sparse with diminishing formation rates as $S \to \infty$. Regularity conditions for deriving asymptotic properties are collected in Condition (S-REG) in the Web Appendix. Applying arguments similar to those in Proposition 3.1 and 3.2 of Lewbel et al. (2023), we have the following proposition.

**Proposition A** *Suppose (N1), (N2), (N3) and (N4) hold. If Assumptions (S-LOB) and (S-REG) hold, then*

$$\widehat{\theta} - \theta = O_p(S^{-1/2} \vee S^{\rho - 1}).$$

*If in addition $\rho < 1/2$, then*

$$\sqrt{S} \left(\widehat{\theta} - \theta\right) \xrightarrow{d} \mathcal{N}(0, \Omega),$$

*where $\Omega \equiv \left(A_0' B_0^{-1} A_0\right)^{-1} A_0' B_0^{-1} \omega_0 B_0^{-1} A_0 \left(A_0' B_0^{-1} A_0\right)^{-1}$ with $A_0, B_0, \omega_0$ being constant arrays defined in Section B3.*

## B3. Proof of Proposition A in Section B2

In this section we derive the asymptotic property of adjusted 2SLS in the setting of a single, large network that is near-block diagonal. Our objective is to show that, when the order of magnitude of the misclassification errors outside the diagonal blocks, or approximate groups, are small enough in the sense of (S-LOB), a 2SLS that only adjusts the link measure within each block while ignoring sparse, off-diagonal links is a root-n, consistent, asymptotically normal estimator for social effects.

To focus on this main goal, we take the misclassification rates $(p_0, p_1)$ as given and fixed in the adjustment. (A proof that also accounts for estimation errors in the initial estimates of $(p_0, p_1)$ would follow from steps similar to Proposition 4 in Section 4.2, but do not add any insight for the main goal.) Also, for conciseness, we only investigate the case with a single, unsymmetrized measure as in Section 3.3.1; parallel results for the case of multiple, symmetrized measure follow from analogous arguments and are omitted for brevity.

We begin by deriving the noisy, feasible structural form in (17). First off, note that the

reduced form of $y_N$ is:

$$
\begin{aligned}
y_N &= (I_N - \lambda G_N)^{-1}(X_N\beta + \varepsilon_N) \\
&= (I_N - \lambda \widetilde{G}_N)^{-1}(X_N\beta + \varepsilon_N) - \left[(I_N - \lambda \widetilde{G}_N)^{-1} - (I_N - \lambda G_N)^{-1}\right](X_N\beta + \varepsilon_N) \quad (19) \\
&= \underbrace{(I_N - \lambda \widetilde{G}_N)^{-1}(X_N\beta + \varepsilon_N)}_{\equiv \widetilde{y}_N} + (I_N - \lambda \widetilde{G}_N)^{-1}\lambda \underbrace{(G_N - \widetilde{G}_N)}_{\equiv \Delta_N}\underbrace{(I_N - \lambda G_N)^{-1}(X_N\beta + \varepsilon_N)}_{=y_N}.
\end{aligned}
$$

where the third equality follows from the fact that $\mathcal{A}^{-1} - \mathcal{B}^{-1} = \mathcal{A}^{-1}(\mathcal{B} - \mathcal{A})\mathcal{B}^{-1}$ for invertible matrices $\mathcal{A}$, $\mathcal{B}$. Next, write (14) as

$$
\begin{aligned}
y_N &= W_N y_N + X_N\beta + \varepsilon_N + \lambda\left(\widetilde{G}_N - W_N\right)y_N + \lambda\Delta_N y_N \\
&= W_N y_N + X_N\beta + \underbrace{\varepsilon_N + \lambda\left(\widetilde{G}_N - W_N\right)\widetilde{y}_N}_{\equiv v_N} + \underbrace{\lambda^2\left(\widetilde{G}_N - W_N\right)(I_N - \lambda \widetilde{G}_N)^{-1}\Delta_N y_N + \lambda\Delta_N y_N}_{\equiv u_N},
\end{aligned}
$$

where the second equality holds because we substitute $y_N$ in $\lambda\left(\widetilde{G}_N - W_N\right)y_N$ using the r.h.s. of (19). Furthermore, we can write

$$
u_N = \left[\lambda\left(\widetilde{G}_N - W_N\right)(I_N - \lambda \widetilde{G}_N)^{-1} + I_N\right]\lambda\Delta_N y_N = (I_N - \lambda W_N)\left(I_N - \lambda \widetilde{G}_N\right)^{-1}\lambda\Delta_N y_N.
$$

This establishes equation (17) in the text.

Next, we introduce the regularity conditions for establishing the asymptotic properties in Proposition 6. Suppose $I_N - \lambda G_N$ and $I_N - \lambda \widetilde{G}_N$ are invertible almost surely, and denote $M_N \equiv (I_N - \lambda G_N)^{-1}$, $\widetilde{M}_N \equiv (I_N - \lambda \widetilde{G}_N)^{-1}$. Let $\widetilde{R}_{N,s} \equiv (W_{N,s}\widetilde{M}_{N,s}X_{N,s}, X_{N,s})$.

(**S-REG**) (i) For all $i$, $\sup_i \left[\sum_j |M_{N,ij}|\right] < \infty$; $\sup_j E\left(|X_{N,j}\beta| + |\varepsilon_{N,j}|| \Delta_N\right) < \infty$; $\sup_j \left|\left(X_N' H_N W_N \widetilde{M}_N\right)_{ij}\right| < \infty$ and $\sup_j \left|\left(X_N' W_N \widetilde{M}_N\right)_{ij}\right| < \infty$ almost surely.
(ii) $(H_{N,s}, \widetilde{G}_{N,s}, X_{N,s}, \epsilon_{N,s})$ are independent across blocks $s = 1, 2, ..., S$.
(iii) There exist $\delta > 0$ s.t. for all $s$, $E\left[||Z_{N,s}'\widetilde{R}_{N,s}||^{1+\delta}\right]$, $E|\left[||Z_{N,s}'W_{N,s}\widetilde{M}_{N,s}\varepsilon_{N,s}||^{1+\delta}\right]$, and $E\left(||Z_{N,s}'Z_{N,s}||^{1+\delta}\right)$ are uniformly bounded.
(iv) For some $\delta > 0$, $E\left||Z_{N,s}'v_{N,s}\right||^{2+\delta} < \Delta < \infty$ and $S^{-1}\sum_{s=1}^S Var(Z_{N,s}'v_{N,s}) > \delta' > 0$ for $S$ sufficiently large.

14

(v) $\sup_j \left| \left[ (I_N - \lambda W_N) \, \widetilde{M}_N \right]_{ij} \right| < \infty$ for all $i$ almost surely.

(vi) $\lim_{S \to \infty} \frac{1}{S} \sum_s E \left( Z'_{N,s} Z_{N,s} \right)$ and $\lim_{S \to \infty} \frac{1}{S} \sum_s E \left( Z'_{N,s} \widetilde{R}_{N,s} \right)$ exist and are non-singular.

Assumption (S-REG) collects regularity conditions needed for deriving the asymptotic properties of $\widehat{\theta} - \theta$. Part (ii) implies that exogenous variables are drawn independently across the blocks. Part (i) and (v) introduce bound conditions on exogenous arrays in the model. These allow us to relate differences between $y_N$ and its near-block diagonal approximation $\widetilde{y}_N$ to the order of difference between $G_N$ and $\widetilde{G}_N$. Parts (iii) and (iv) are boundedness conditions on population moments that ensure a law of large numbers and a central limit theorem apply to components of the estimator.

**Lemma A1.** *Let $a_N$, $b_N$ be random vectors in $\mathbb{R}^N$. Suppose there exist constants $C_1, C_2 < \infty$ such that $\Pr\{\sup_{i \leq N} E(|a_i| \| \Delta_N) \leq C_1\} = 1$ and $\Pr\{\sup_{j \leq N} E(|b_j| \| \Delta_N) \leq C_2\} = 1$. Then Assumption S-LOB implies $\frac{1}{S} a'_N \Delta_N b_N = O_p(S^{\rho-1})$.*

*Proof of Lemma A1.* From Assumption S-LOB, $\sum_i \sum_j E |\Delta_{N,ij}| = O(S^\rho)$ for some $\rho < 1$. By construction,

$$
\begin{aligned}
E \left( \left| \tfrac{1}{S} a'_N \Delta b_N \right| \right) &\leq \tfrac{1}{S} E \left[ \sup_{i,j} E \left( |a_i b_j| \mid \Delta_N \right) \cdot \left( \sum_i \sum_j |\Delta_{N,ij}| \right) \right] \\
&\leq \tfrac{1}{S} E \left[ C_1 C_2 \left( \sum_i \sum_j |\Delta_{N,ij}| \right) \right] = O(S^{\rho-1}).
\end{aligned}
$$

It then follows that $\frac{1}{S} a'_N \Delta_N b_N = O_p(S^{\rho-1})$. $\qquad\square$

**Lemma A2.** *Under the conditions in (S-REG)-(i), there exists a constant $C^* < \infty$ such that $\Pr\{\sup_{i \leq N} E(|y_i| \| \Delta_N) \leq C^*\} = 1$ for all $N$.*

*Proof of Lemma A2.* Let $M_N \equiv (I_N - \lambda G_N)^{-1}$. For any matrix $\mathcal{A}$, let $\mathcal{A}_{(i)}$ denote its $i$-th row; and $\mathcal{A}_{ij}$ denote its $(i,j)$-th component. It then follows from the reduced form that

$$
\begin{aligned}
\sup_{i \leq N} E(|y_{N,i}| \mid \Delta_N) &= \sup_i E \left( \left| \sum_j M_{N,ij} \left( X_{N,(j)} \beta + \varepsilon_j \right) \right| \, \Big\| \Delta_N \right) \\
&\leq \sup_i \left[ \sum_j |M_{N,ij}| \right] \times \sup_j E \left( |X_{N,(j)} \beta| + |\varepsilon_{N,j}| \| \Delta_N \right).
\end{aligned}
$$

15

Hence, there exists some constant $C^* < \infty$ with $\Pr\{\sup_i E(|y_i| \| \Delta_N) \le C^*\} = 1$. □

**Lemma A3.** *Under the conditions in (S-REG), $\frac{1}{S} R'_N Z_N = A_0 + o_p(1)$, $\frac{1}{S} Z'_N Z_N = B_0 + o_p(1)$, and $\frac{1}{S} Z'_N v_N = O_p(S^{-1/2})$.*

*Proof of Lemma A3.* By definition, $\frac{1}{S} Z'_N Z_N = \frac{1}{S} \sum_{s=1}^S Z'_{N,s} Z_{N,s}$, with $Z_{N,s}$ independent across $s$ due to (S-REG)-(ii). Then by (S-REG)-(iii) and the law of large numbers for independent and heterogeneously distributed observations (e.g., Corollary 3.9 in White (2001)), $\frac{1}{S} Z'_N Z_N = B_0 + o_p(1)$ where $B_0 \equiv \lim_{S\to\infty} \frac{1}{S} \sum_s E\left(Z'_{N,s} Z_{N,s}\right)$. Next, note by construction and (19),

$$
\frac{1}{S} Z'_N R_N = \frac{1}{S} \begin{pmatrix} X'_N H_N W_N \widetilde{y}_N & X'_N H_N X_N \\ X'_N W_N \widetilde{y}_N & X'_N X_N \end{pmatrix} + \frac{1}{S} \lambda \begin{pmatrix} X'_N H_N W_N \widetilde{M}_N \Delta_N y_N & 0 \\ X'_N W_N \widetilde{M}_N \Delta_N y_N & 0 \end{pmatrix}. \tag{20}
$$

By (S-REG)-(i) and Lemma A2, $y_N$ satisfies the condition on $b_N$ in Lemma A1. It then follows from Lemma A1 that the *second* term on the right-hand side of (20) is $O_p(S^{\rho-1})$. Besides, the *first* term on the r.h.s. of (20) is

$$
\frac{1}{S} \sum_{s=1}^S Z'_{N,s} \widetilde{R}_{N,s} + \frac{1}{S} \sum_{s=1}^S \left(Z'_{N,s} W_{N,s} \widetilde{M}_{N,s} \varepsilon_{N,s}, \mathbf{0}\right). \tag{21}
$$

By (N3), $E\left(Z'_{N,s} W_{N,s} \widetilde{M}_{N,s} \varepsilon_{N,s}\right) = 0$. It then follows from (S-REG)-(iii) that the expression in (21) is $A_0 + o_p(1)$, with $A_0 \equiv \lim_{S\to\infty} \frac{1}{S} \sum_s E\left(Z'_{N,s} \widetilde{R}_{N,s}\right)$.

Next, note that by definition,

$$
\frac{1}{S} Z'_N v_N = \frac{1}{S} \sum_{s=1}^S Z'_{N,s} \varepsilon_{N,s} + \lambda \frac{1}{S} \sum_{s=1}^S Z'_{N,s} \left(\widetilde{G}_{N,s} - W_{N,s}\right) \widetilde{y}_{N,s}. \tag{22}
$$

By construction, $Z_{N,s}$, $\varepsilon_{N,s}$, $\widetilde{G}_{N,s}$ and $H_{N,s}$ are independent across blocks $s = 1, 2, ..., S$. Also, recall that $\widetilde{y}_{N,s}$ is defined as $\widetilde{y}_{N,s} \equiv (I_s - \lambda \widetilde{G}_{N,s})^{-1}(X_{N,s}\beta + \varepsilon_{N,s})$, Hence $\widetilde{y}_{N,s}$ is also independent across the blocks. Assumption (N3) implies $E(Z'_{N,s} \varepsilon_{N,s}) = 0$; Assumptions (N1) and (N2) imply

$$
E\left(W_{N,s} | \widetilde{G}_{N,s}, X_{N,s}\right) = \widetilde{G}_{N,s}.
$$

16

Furthermore, the same argument as in the proof of Proposition 2 shows that under (N1), (N2), (N3) and (N4)

$$E\left(H_{N,s}W_{N,s}|\widetilde{G}_{N,s}, X_{N,s}\right) = E\left(H_{N,s}\widetilde{G}_{N,s}\Big|\widetilde{G}_{N,s}, X_{N,s}\right),$$

so that

$$E\left[Z'_{N,s}\left(\widetilde{G}_{N,s} - W_{N,s}\right)\widetilde{y}_{N,s}\right] = 0.$$

It then follows from (S-REG)-(iv) and the Central Limit Theorem that $\frac{1}{S}Z'_N v_N = O_p(S^{-1/2})$. □

*Proof of Proposition A.* As shown in Lemma A3, $\frac{1}{S}R'_N Z_N = A_0 + o_p(1)$, $\frac{1}{S}Z'_N Z_N = B_0 + o_p(1)$, and $\frac{1}{S}Z'_N v_N = O_p(S^{-1/2})$ under (N1)-(N4), (S-LOB) and (S-REG). Furthermore, with (S-REG)-(v), Lemma A1 and Lemma A2 imply that $\frac{1}{S}Z'_N u_N = O_p(S^{\rho-1})$. When $\rho < 1/2$, we have

$$\frac{1}{\sqrt{S}}Z'_N(u_N + v_N) \xrightarrow{d} \frac{1}{\sqrt{S}}Z'_N v_N \xrightarrow{d} \mathcal{N}(0, \omega_0),$$

where $\omega_0 = \lim_{S\to\infty} \frac{1}{S}\sum_s E\left(Z'_{N,s} v_{N,s} v'_{N,s} Z_{N,s}\right)$. Hence,

$$\sqrt{S}(\widehat{\theta} - \theta) = \left(A'_0 B_0^{-1} A_0\right)^{-1} A'_0 B_0^{-1}\frac{1}{\sqrt{S}}Z'_N v_N + o_p(1)$$
$$\xrightarrow{d} \mathcal{N}(0, \left(A'_0 B_0^{-1} A_0\right)^{-1} A'_0 B_0^{-1}\omega_0 B_0^{-1} A_0 \left(A'_0 B_0^{-1} A_0\right)^{-1}).$$

□

# References

Lewbel, A., X. Qu, and X. Tang (2023). Ignoring measurement errors in social networks. *The Econometrics Journal*, utad028. pages 12, 13

White, H. (2001). *Asymptotic theory for econometricians.* Academic press. pages 16