



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Social Networks 28 (2006) 247–268

**SOCIAL  
NETWORKS**

[www.elsevier.com/locate/socnet](http://www.elsevier.com/locate/socnet)

# Effects of missing data in social networks

Gueorgi Kossinets\*

*Department of Sociology and Institute for Social and Economic Research and Policy,  
Columbia University, 420 W. 118th Street, 8th Floor, Mail Code 3355,  
New York, NY 10027, USA*

---

## Abstract

We perform sensitivity analyses to assess the impact of missing data on the structural properties of social networks. The social network is conceived of as being generated by a bipartite graph, in which actors are linked together via multiple interaction contexts or affiliations. We discuss three principal missing data mechanisms: network boundary specification (non-inclusion of actors or affiliations), survey non-response, and censoring by vertex degree (fixed choice design), examining their impact on the scientific collaboration network from the Los Alamos E-print Archive as well as random bipartite graphs. The simulation results show that network boundary specification and fixed choice designs can dramatically alter estimates of network-level statistics. The observed clustering and assortativity coefficients are overestimated via omission of affiliations or fixed choice thereof, and underestimated via actor non-response, which results in inflated measurement error. We also find that social networks with multiple interaction contexts may have certain interesting properties due to the presence of overlapping cliques. In particular, assortativity by degree does not necessarily improve network robustness to random omission of nodes as predicted by current theory.

© 2005 Elsevier B.V. All rights reserved.

*PACS:* C34 (truncated and censored models); C52 (model evaluation and testing)

*Keywords:* Missing data; Sensitivity analysis; Graph theory; Collaboration networks; Bipartite graphs

---

\* Tel.: +1 212 854 0367; fax: +1 212 854 7998.

*E-mail address:* [gk297@columbia.edu](mailto:gk297@columbia.edu).

## 1. Introduction

Social network data is often incomplete, which means that some actors or links are missing from the dataset. In a normal social setting, much of the incompleteness arises from the following sources: the so-called boundary specification problem (Laumann et al., 1983); respondent inaccuracy<sup>1</sup> (Bernard et al., 1984; Brewer and Webster, 1999; Marsden, 1990; Butts, 2003); non-response in network surveys (Stork and Richards, 1992; Rumsey, 1993; Robins et al., 2004); or may be introduced via study design (Burt, 1987). Compound missing data mechanisms may be encountered as well. Although missing data is abundant in empirical studies, little research has been conducted on the possible effect of missing links or nodes on the measurable properties of networks at large. In particular, a revision of the original work done primarily in the 1970–1980s (Holland and Leinhard, 1973; Laumann et al., 1983; Bernard et al., 1984) seems appropriate in the light of recent advances that have brought new classes of network models to the attention of the interdisciplinary research community (Amaral et al., 2000; Barabási and Albert, 1999; Newman et al., 2001; Strogatz, 2001; Watts and Strogatz, 1998; Watts, 1999).

This paper aims to highlight the problem of missing data in social network analysis. One approach to deal with it is to develop analytic techniques that capture global statistical tendencies and do not depend on individual interactions (Rapoport and Horvath, 1961). A complementary strategy is to develop remedial techniques that minimize the effect of missing data (Holland and Leinhard, 1973; Robins et al., 2004). Although we do not offer a definitive statistical treatment in this paper, we conduct exploratory analyses and advocate the importance of further work in this direction (cf. Costenbader and Valente, 2003). We use the method of statistical simulation to quantify the uncertainty caused by missing network data and assess sensitivity of graph-level metrics such as average vertex degree, clustering coefficient (Newman et al., 2001), degree correlation coefficient (Newman, 2002), size and mean path length in the largest connected component. Our dataset is the scientific collaboration graph containing authors and papers from the Condensed Matter section of the Los Alamos E-print Archive from 1995 through 1999 (Newman, 2001). We use this example to develop a statistical argument for the general case of social networks with multiple interaction contexts. Owing to the sheer size of the dataset, the numerical estimates have very narrow confidence intervals. The results are compared to the case of random bipartite graphs.

The paper is organized as follows. Section 2 focuses on the sources of missing or false data in social network research. We generalize the boundary specification problem (BSP) for social networks with multiple interaction contexts modeled as bipartite graphs, in which actors are linked via multiple affiliations or collaborations. We discuss the issues of non-response and non-reciprocation in social network studies as well as the degree cutoff bias often introduced by questionnaire design. Section 3 describes relevant network statistics, datasets and simulation algorithms that are used to investigate effects of missing data on network properties. Section 4 presents the results, while Section 5 summarizes the findings.

---

<sup>1</sup> In this paper we do not explicitly model the effect of informant inaccuracy, assuming that either it is consistent with the research framework, or that the network in question was reconstructed from reliable electronic, historical or survey data.

## 2. Sources of missing data in social networks

### 2.1. The boundary specification problem

The boundary specification problem (Laumann et al., 1983) refers to the task of specifying inclusion rules for actors or relations in a network study (Fig. 1). For example, researchers who study intraorganizational networks typically ignore numerous ties that lead outside an organization, reasoning that these ties are irrelevant to the tasks and operations that the organization performs. A classical account is the Bank Wiring Room study (Roethlisberger and Dickson, 1939) which focused on 14 men in the switchboard production section of an electric plant. The sociometric data obtained in that study have been analyzed extensively (Homans, 1950; White et al., 1976) but the effect of interactions outside the wiring room on the workers' behavior and performance at work is unknown and hardly feasible to estimate. The boundary specification problem may be avoided to a certain extent if the community is isolated from the rest of the world as e.g. in Sampson's monastery (Sampson, 1969). By and large, however, network closure is an artifact of research design, i.e. the result of arbitrary definition of network boundaries. Examples include networks based on the formal definition of group membership or positional specification, most commonly defined as occupancy of a ranked position in a formally constituted group, e.g. a country's 100 best known politicians, or 500 top business firms (e.g. Davis and Mizruchi, 1999). When choosing inclusion rules for a network study, a researcher is effectively drawing a non-probability sample from all possible networks of its kind (Laumann et al., 1983). Dynamic changes in the network only exacerbate the problem. An approach advocated by Laumann et al. (1983) is to focus on measurable interactions. The network boundary is then defined by recording who is interacting with whom in a certain context. This approach has been feasible only for small networks until recently, when data on large-scale social interactions become readily available from the records of email communication or virtual communities (Ebel et al., 2002; Guimera et al., 2003; Holme et al., 2004; Newman et al., 2002). It requires an operational specification of the interaction setting or context, and then including all actors who interact within this context.

Since social networks are constructed from actors and relations between actors, the boundary specification problem has two faces to it. In addition to defining a network bound-

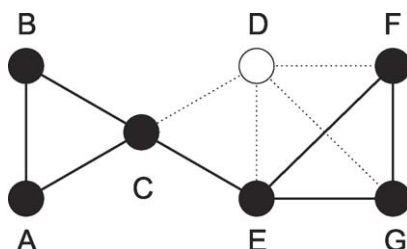


Fig. 1. Illustration of the boundary specification problem. Omission of actors may lead to significant changes in network statistics. In the above example, as a result of exclusion of actor D, the mean network degree  $\bar{z}$  went down 25% from  $3\frac{1}{3}$  to  $2\frac{1}{3}$ .

ary over the set of actors, researchers make decisions on which relations to consider. Here we employ a multicontextual approach based on actors' participation in groups, events or activities. It is the joint network, made by juxtaposition of all relevant kinds of ties between actors, that matters in dynamics of processes based on social influence (White, 1992; White et al., 1976). Each jointly attended event, shared affiliation or interaction context serves as an opportunity to create, maintain, or exercise (manipulate) group and interpersonal ties. The above examples can be represented by a bipartite graph (Wilson, 1982; Wasserman and Faust, 1994), in which one class of vertices represents events, and the second class is actors.<sup>2</sup> If an actor participates in an event, there is an edge drawn between the respective vertices. To focus on the class of actors, we transform the two-mode "affiliation" graph into a one-mode network that captures multiple social relations between the actors (Fig. 2). One-mode projections necessarily consist of many overlapping cliques.<sup>3</sup> Every such clique refers to one or several affiliations or interaction contexts. In the bipartite framework an affiliation tie is added to the network if an actor has participated in the given context. However, correlated contexts are somewhat redundant, in the sense that they contain much the same information about social structure.

The network approach has traditionally sought to separate different relational contexts for the sake of analytical tractability. A textbook definition of a social network (Wasserman and Faust, 1994) assumes a discrete set of actors linked together by a discrete set of relations. One-mode networks have been studied extensively in the recent years with a number of important analytic results obtained (Albert et al., 2000; Barabási and Albert, 1999; Callaway et al., 2000; Cohen et al., 2000, 2001; Newman et al., 2001; Watts and Strogatz, 1998). However, this line of research has focused on simple models for the network (e.g. randomly mixed with respect to vertex degree), which are unlikely to hold in most real situations where both structural and attribute-based processes are important (Girvan and Newman, 2002; Watts et al., 2002; White, 1992). We therefore propose that the multicontextual model of a social network (generated by a bipartite graph) has certain advantages over the models based on simple random graphs. Formulated in a suitable manner, it is analytically tractable (Newman et al., 2001; Watts et al., 2002) and by definition takes care of certain properties observed in empirical social networks that are not easily reproducible with simple random graphs (such as high clustering).

## 2.2. Non-response effects

An important problem in network survey research is that of *survey non-response*. In a standard sampling situation such as drawing a representative sample from some population, special techniques are available to correct parameter estimates for imperfect response rates (Little and Rubin, 2002). Unfortunately, no such definitive treatment is available for social network analysis, although effects of non-response on some network properties have been described previously (Stork and Richards, 1992; Rumsey, 1993). We generally follow

<sup>2</sup> Given the conceptual similarity of affiliation networks, social event attendance and multiple interaction contexts, in the discussion that follows we will take the liberty of using the terms "events", "contexts" or "affiliations" interchangeably, unless specifically mentioned otherwise.

<sup>3</sup> Note that a dyad is a clique of size 2.

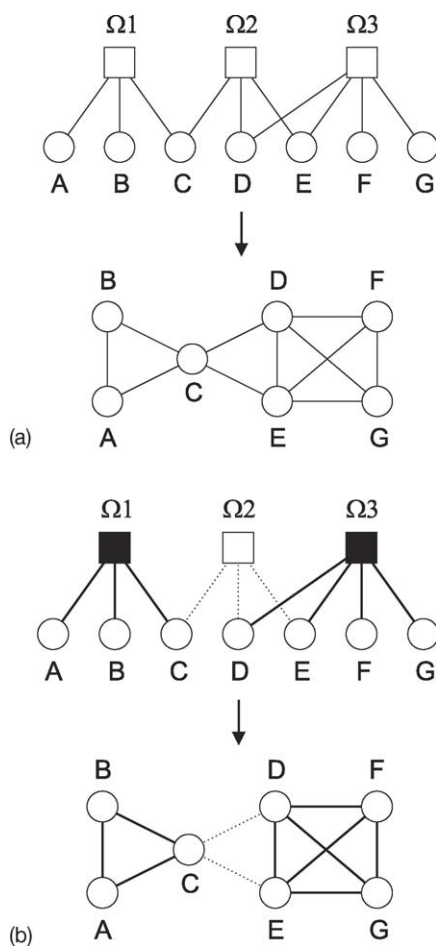


Fig. 2. (a) Explanation of the unipartite projection. Given a bipartite (or ‘two-mode’) affiliation graph, a new network is defined on the set of actors, where two actors are connected if they belong to one or more contexts together in the association graph. In the above example, there are seven actors (A–G) and three groups ( $\Omega 1$ – $\Omega 3$ ). Observe three overlapping cliques in the one-mode projection (ABC, CDE, and DEFG) corresponding to the three interaction contexts. It is possible to differentiate between different levels of intensity of links in the unipartite projection by assigning a weight to each context and calculating a summary weight for each connected pair of actors. However, for the points we wish to make here it is sufficient to use the simple undirected graph representation; that is, to be able to tell if any two actors are connected or not, neglecting the ‘strength’ of connection. (b) Boundary specification problem for relations. Suppose that we fail to include interaction context  $\Omega 2$  in the above example. That may have a drastic effect on the observed properties of the one-mode network, e.g. it may become disconnected, etc.

this exploratory line of research in that we discuss how network structure is affected by simple non-response scenarios and propose some ways to ameliorate the problem. While single-mode networks with non-respondents have been shown to be amenable to statistical treatment (Robins et al., 2004), non-response in networks with multiple interaction contexts

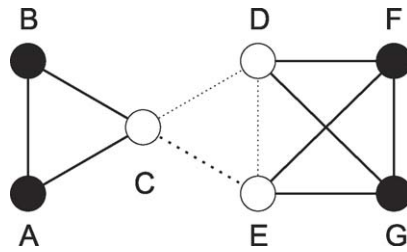


Fig. 3. Non-response in network surveys. Suppose that actors C, D and E did not report their links. However, nominations made by actors A, B, F and G help reconstruct the structure of interactions to a large extent, with a decrease in average degree less than 15%. Compare with the Boundary Specification example (Fig. 1), in which a single missing node caused a 25% deviation in the mean degree.

(modeled as bipartite graphs) may have a number of specific implications. In a survey of an affiliation network, actors are asked to report groups to which they belong. Suppose that we have no other sources of information about affiliations. If any one actor fails to respond, all his affiliations are lost and the resulting missing data pattern becomes equivalent to the boundary specification problem for actors which we model as stochastic omission of some fraction of actors from the network. If however the survey asks actors to name peers with whom they interact (that is, ignoring the multiplexity of ties), then the non-response effect can be balanced out by reciprocal nominations (Stork and Richards, 1992). Suppose actor C did not fill in the network questionnaire (Fig. 3). Yet those of C's interactants who participated in the survey (A and B) must have reported their interactions with C. Intuitively, one would expect that if the number of non-respondents is small relative to the size of the network, and the researcher does not require all nominations to be reciprocated (as a crude validity check), then the amount of missing data caused by non-response should be small if not negligible.

### 2.3. Fixed choice designs

Sometimes social network data may be biased as a result of study design.<sup>4</sup> In this paper we analyze the so-called *fixed choice effect* (Holland and Leinhard, 1973). Consider a friendship network in which actors have anywhere between 1 and 10 friends each. Often network researchers ask respondents to make nominations only up to some fixed number. One would like to know whether and how the network constructed in that particular way is different from the “true” friendship network.

Fixed choice designs introduce right-censoring by vertex degree (Holland and Leinhard, 1973). This missing data mechanism is often present in network surveys. Suppose that actor A belongs to  $k$  groups whereby he is connected to  $x$  other actors (Fig. 4a). In the unipartite case, the actor is requested to nominate up to  $X$  persons from his list of  $x$  interactants, e.g. “ $X$  best friends” (Fig. 4b). If the cutoff is greater than or equal to the actual number

<sup>4</sup> Unless the design explicitly makes use of inherent biases as e.g. in respondent-driven sampling (Salganik and Heckathorn, 2004).

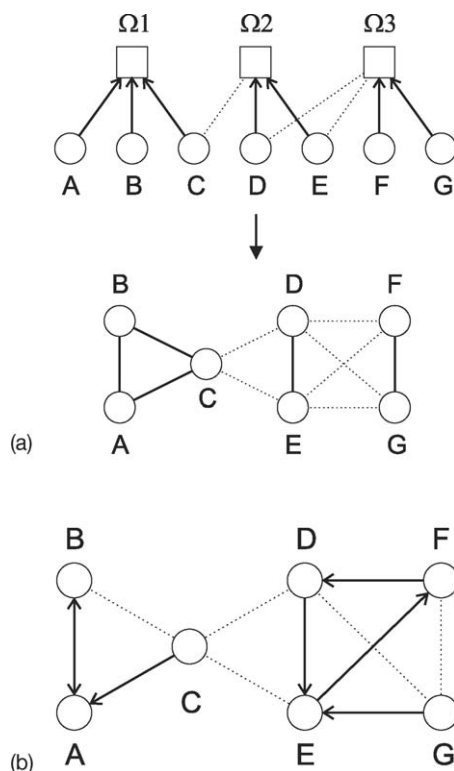


Fig. 4. Illustration of a fixed choice design. (a) Bipartite case: each actor nominates up to a fixed number  $K$  from his affiliations. Nominations are shown as arrows. (b) One-mode case: each actor nominates up to a fixed number  $X$  from his list of acquaintances. In the hypothetical example pictured above  $K = X = 1$ . Note that there is only one reciprocated nomination (between actors A and B).

of friends ( $X \geq x$ ), we assume that all  $x$  links between A and his friends are included in the dataset. If  $X < x$ , actor A must omit  $x - X$  links, but some of those might still be reported by A's friends who are requested to make their nominations likewise. Thus some ties from the original network will be reported by both interactants (reciprocated nominations), some by only one partner (non-reciprocated nominations), and yet some will not be reported (censored links). It is left to the discretion of the researcher whether to include non-reciprocated links which may be qualitatively different from reciprocated ones (e.g., good friends versus casual acquaintances). Fixed choice nominations can easily lead to a non-random missing data pattern. For instance, popular individuals who have more contacts may be more likely to be nominated by their contacts (Feld, 1991; Newman, 2003a). The effect may be different depending on whether the network is mixed disassortatively or assortatively by degree (Newman, 2002; Vázquez and Moreno, 2003): in the first case, vertices with high degrees tend to be matched with vertices with less connections and therefore more censored connections are likely to be restored using reciprocal nominations. This is an example of how the network structure may interact with missing data mechanisms.

### 3. Data and statistics of interest

#### 3.1. Network-level statistics

As we wish to investigate how topological properties of the network are affected by the presence of missing vertices or edges, we measure the following graph-level properties of the unipartite projection onto actors: mean vertex degree  $z$  (average number of interactants per actor), which characterizes network connectivity; clustering  $C$ , loosely interpreted as the probability that any two vertices with a mutual neighbor are themselves connected<sup>5</sup>; assortativity  $r$ , which is the Pearson correlation coefficient of the degrees at endpoints of an edge (Newman, 2002); fractional size of the largest connected component  $S$ ; and average path length (mean geodesic distance) between all pairs of vertices in the largest component  $\ell$ . We accept that the effect of missing data on parameter  $Q$  is tolerable if the relative error  $\varepsilon = \frac{|q-q_0|}{q_0} \leq 10\%$ , where  $q$  is an estimate from a model with missing data and  $q_0$  the respective “true” value calculated from all available data.<sup>6</sup>

#### 3.2. Data

Following previous work, we treat collaboration and affiliation graphs as examples of multicontextual social networks (Davis and Mizruchi, 1999; Mizruchi, 1996; Newman, 2001). We illustrate the problem of missing data in networks using the example of the scientific collaboration graph containing authors and papers from the Condensed Matter section (“cond-mat”) of the Los Alamos E-print Archive from 1995 through 1999 (Newman, 2001) as well as random bipartite graphs. The properties of the dataset are summarized in Table 1.

We compare the collaboration graph to an ensemble of 100 random bipartite graphs with the same number of vertices and edges, i.e. fixing the number of actors  $N = 16726$ , number of groups  $M = 22016$ , mean degree  $\mu = 3.50$  for actors and  $\nu = 2.66$  for groups<sup>7</sup> (Fig. 5b). The degree sequence is not fixed and so is well approximated by the Poisson distribution (Bollobás, 2001; Newman et al., 2001). In the Condensed Matter collaboration network, both the distribution of the number of authors per paper and the distribution of papers per author are considerably skewed to the left relative to the random model (Fig. 5a). The distribution of vertex degree in the one-mode coauthor network (i.e. the number of co-authors) resembles a power-law with exponential cutoff near  $k = 100$  (Fig 5a, dots) while the same distribution in a random graph exhibits the characteristic bimodal shape (Newman et al., 2001) with a clear cutoff in the tail (Fig. 5b). In the unipartite projection of a random bipartite graph there are many vertices with a medium connectivity while very few

<sup>5</sup> There are several ways to measure clustering (Watts, 1999; Newman et al., 2001; Maslov et al., 2002). We adopt the following definition of clustering coefficient:  $C = 3N_{\Delta}/N_3$ , where  $N_{\Delta}$  is the number of triangles in the graph and  $N_3$  is the number of connected triples of vertices. This definition is more representative of average clustering in cases when vertex degree distribution is skewed (Newman et al., 2001).

<sup>6</sup> This measure is sensible only for variables with zero as a natural lower bound, so we do not calculate it for assortativity.

<sup>7</sup> Actually, we need to fix only three parameters since  $\mu N = \nu M$ .



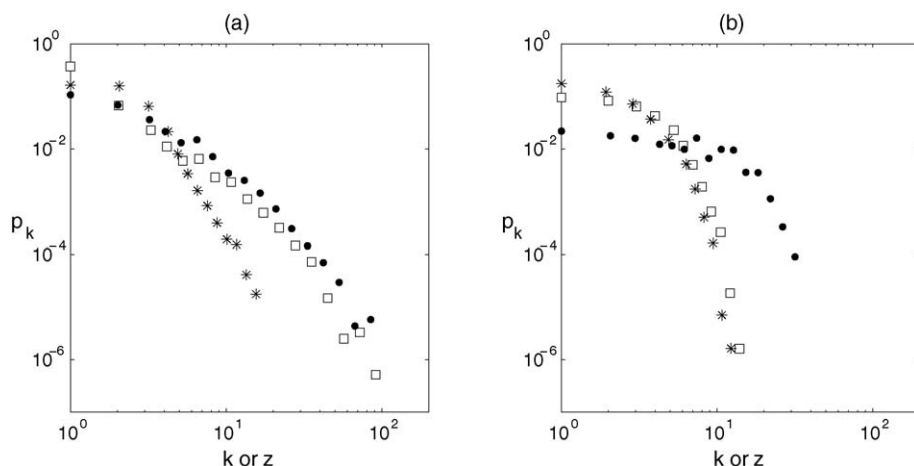


Fig. 5. Distributions of vertex degree in the Condensed Matter collaboration graph (a) and in the comparison random network (b). Squares: number of papers per author; stars: number of authors per paper; dots: number of collaborators per author. The data have been logarithmically binned.

vertices with a very large number of coauthors. The values of mean degree in the one-mode projection are  $z = 5.69$  for the cond-mat graph and  $z = 9.31$  for its random counterpart, which indicates a strongly non-random allocation of authors over papers in the Condensed Matter collaboration network. In both cases  $z \gg 1$ , which implies the existence of the giant connected component (Bollobás, 2001).

As seen from Table 1 the bipartite form of the Condensed Matter collaboration graph is disassortative ( $r_B = -0.18$ ) whereas its one-mode projection is assortative ( $r_U = 0.18$ ).

Table 1  
Properties of the network dataset

Quantity	Notation	Cond-mat	Random <sup>a</sup>
Number of authors	$N$	16726	16726
Number of papers	$M$	22016	22016
Mean papers per author	$\mu$	3.50	3.50
Mean authors per paper	$\nu$	2.66	2.66
Assortativity (degree correlation)	$r_B$	-0.18	-0.054 (4)
Unipartite projection (collaborators)			
Mean degree	$z$	5.69	9.31 (3)
Degree variance	$V$	41.2	33.9 (6)
Clustering	$C$	0.36	0.223 (1)
Assortativity	$r_U$	0.18	0.071 (5)
Number of components	$N_C$	1188	652 (18)
Size of largest component	$S_L$	13861	16064 (18)
Mean path in largest component	$\ell_L$	6.63	4.728 (8)

<sup>a</sup> A random bipartite graph of the same size and mean degree as the original network. Numbers in parentheses are standard deviations on the least significant figures calculated in an ensemble of 100 such graphs.

This implies that authors who work in smaller collaborations publish more papers on average; also, physicists with many collaborators tend to work with those of the same ilk; and similarly, physicists with a few coauthors who are, incidentally, most prolific ones, tend to collaborate with each other.<sup>8</sup> In addition to providing curious insights into the mode of scientific production in Condensed Matter Physics, assortativity has important implications for network robustness (Boguñá et al., 2003; Newman, 2002; Vázquez and Moreno, 2003).

A characteristic feature of assortatively mixed ( $r_U > 0$ ) networks is the so-called core group consisting of interconnected high-degree vertices. The core group provides exponentially many distinct pathways to connect vertices of smaller degrees. From an epidemiology point of view, the core forms a reservoir that is capable of sustaining a disease outbreak even though the overall network density is too low for an epidemic to occur. The good news, however, is that an outbreak in assortatively mixed networks is likely to be confined to a smaller subset of the vertices. Disassortative networks are particularly susceptible to targeted attacks on high-degree vertices due to the fact that the latter provide much of the global network connectivity (Newman, 2003a).

Although a random graph has zero assortativity in the limit of large system size, it may acquire some disassortativity as a finite-size effect, e.g. from the constraint forbidding multiple edges between two vertices (Maslov et al., 2002; Newman, 2003a). In a similar fashion, random bipartite graphs exhibit disassortative mixing if the number of groups differs from the number of actors. This follows from the definition of a bipartite graph (no edges connect vertices of the same class) and the requirement that no actor belongs to the same group twice. The ensemble of random bipartite graphs simulated here exhibit small but significant disassortativity ( $r_B = -0.054 \pm 0.004$ ) while the corresponding one-mode networks are assortatively mixed by degree ( $r_U = -0.071 \pm 0.005$ ). It is important to keep in mind that clustering, assortativity (or generally, the mixing pattern) and degree distribution are not independent. In particular, disassortative mixing in simple graphs may cause a decrease in clustering by suppressing connections between high-degree vertices in favor of vertices of lower degree, thus reducing the number of triads in the network (Maslov et al., 2002; Newman, 2003a).

### 3.3. Algorithms

The outline of the simulation algorithm is as follows: (1) take a real (large enough) social network or an ensemble of random graphs and assume that network data is complete; (2) remove a fraction of entities to simulate different sources of error; (3) measure network properties and compare to the “true” values (from the “complete” network); (4) repeat (2)–(3) a number of times to generate distributions of statistics of interest. Table 2 summarizes our simulation models.

<sup>8</sup> Additional simulations (not shown here) indicate that the presence of heavy-tailed group size distribution in a bipartite graph may cause assortativity in its one-mode projection onto actors. This lead us to suggest that assortativity of the one-mode Physics collaboration graph might be to some extent an artifact of the skewed distribution of collaboration sizes.

Table 2  
Simulation algorithms for sensitivity analysis

Label	Problem	Model <sup>a</sup>
BSPC	Boundary specification problem for contexts	Remove a fraction of contexts at random
BSPA	Boundary specification problem for actors	Remove a fraction of actors at random
NRE	Non-response effect	Remove links within subgraph induced by a specified fraction of actors
FCC	Fixed choice (contexts)	Apply censoring by degree to actors
FCA	Fixed choice (actors)	Create unipartite projection; apply censoring by degree; keep non-reciprocated links
FCR	Fixed choice (actors), reciprocated nominations only	Create unipartite projection; apply censoring by degree; keep only reciprocated links

<sup>a</sup> We measure properties of the unipartite projection in all models.

## 4. Results and discussion

### 4.1. Comparison of boundary specification and non-response effects

The results of the simulations for the Condensed Matter collaboration graph and for comparable random bipartite networks are plotted in Figs. 6, 8–11. The proportion of missing data increases from left to right and at the leftmost point we assume that all information about the network is available. We model the boundary specification problem for contexts (BSPC) by randomly removing vertices of the corresponding class (“papers”) from the network. The boundary specification problem for actors (BSPA) is modeled as random deletion of vertices corresponding to “authors” in the case of collaboration network. Survey non-response is different from BSPA in that in the former vertices are not removed from the network but all edges between randomly assigned “non-respondents” are deleted.

#### 4.1.1. Mean vertex degree

For a random bipartite graph, the mean degree in the unipartite projection onto actors decreases linearly with random removal of actors or groups:  $z = \mu\nu(1 - \theta)$ , where  $\theta$  is a relative number of missing actors or groups, respectively<sup>9</sup> (observe overlapping curves in Fig. 6b). However, in the one-mode collaboration network average degree decreases slower in the simulation of BSPC (Fig. 6a, dots) than in BSPA (squares). This behavior implies non-random allocation of actors (authors) to groups (papers) and leads us to introduce the notion of “redundancy” in group affiliation.

One way to capture the average importance of an interaction context is to measure what we call the *redundancy* of a bipartite graph. We define redundancy as  $\beta = \frac{\mu\nu - z}{\mu\nu} = 1 - \frac{z}{\mu\nu}$ , where  $\mu$  is average number of groups per actor,  $\nu$  is average size of the group, and  $z$  is actual (observed) mean actor degree in the unipartite projection onto the set of actors. In a complete bipartite graph all affiliations but one are redundant in the sense that they connect actors who

<sup>9</sup> Here we have made use of the fact that the mean vertex degree  $z = \mu\nu$  in the unipartite projection of random bipartite graph (for large  $N$  and  $M$ ,  $\mu > 1$  and  $\nu > 1$ ), which is symmetrical with respect to changes in either  $\mu$  or  $\nu$  (Newman et al., 2001).

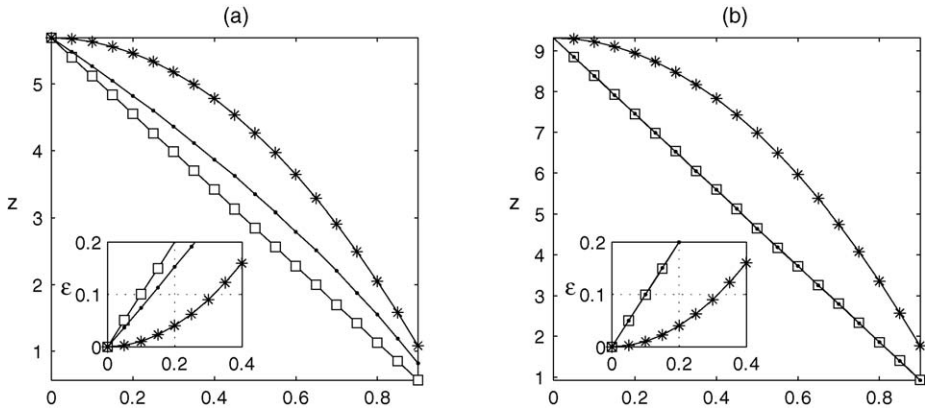


Fig. 6. Sensitivity of mean vertex degree in the unipartite projection  $z$  to different missing data mechanisms: (a) in the Condensed Matter graph; (b) in a bipartite random graph. Dots: boundary specification (non-inclusion) effect for interaction contexts (BSPC); the horizontal axis corresponds to the fraction of papers missing from the database. Squares: non-inclusion effect for actors (BSPA) with the  $x$ -axis corresponding to the fraction of authors missing from the database. Note that in panel (b) dots overlap with squares. Stars: simulation of survey non-response among authors (NRE); vertices are assumed non-responding at random. The  $x$ -axis indicates the fraction of non-respondents. Insets: relative error  $\varepsilon = |z - z_0|/z_0$ , where  $z_0$  is the true value. Each data point is an average over 50 iterations. Lines connecting datapoints are meant as a guide for the eye.

are already connected (Fig. 7a), consequently  $\beta_C = 1 - \frac{N-1}{MN} \rightarrow 1$  as  $M \rightarrow \infty$  ( $M$  is the number of affiliations). At the other extreme are acyclic bipartite graphs (Fig. 7b), in which if any two actors belong to the same affiliation it is the only affiliation they share, therefore  $z = \mu\nu$  and  $\beta_A = 0$ . Consider a bipartite graph such that every connected pair of actors have attended exactly three events together. The mean degree in the actors one-mode network will be  $z = \mu\nu/3$ , and redundancy therefore is  $\beta = 1 - 1/3 = 2/3$ . Redundancy of a random bipartite graph is expected to be close to zero since  $z \approx \mu\nu$ , which becomes exact as the graph size increases (Newman et al., 2001). In general, high redundancy implies that as new interaction contexts emerge, they will likely link already connected actors. Redundancy of the Condensed Matter collaboration graph is  $\beta = 1 - 5.69/(3.50 \times 2.66) \approx 0.38$ , which means that if the collaboration sizes were sharply peaked around the mean, then about forty percent of collaborations could be omitted without any significant change in the structure of unipartite projection. However, this is not exactly the case here (Fig. 6a) because the group size distribution is quite skewed (Fig. 5a). There are certain important collaborations

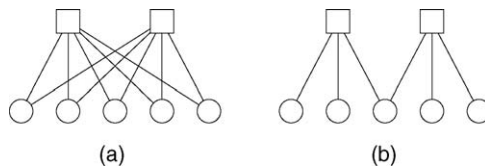


Fig. 7. Examples of (a) complete (maximally redundant) and (b) acyclic (non-redundant) bipartite graphs.

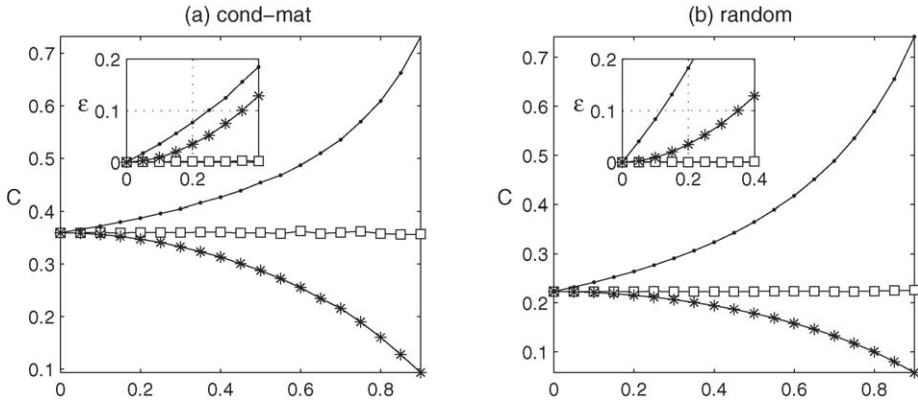


Fig. 8. Sensitivity of clustering  $C$  in the unipartite projection: omission of interaction contexts (dots); omission of actors (squares); survey non-response (stars).

that serve as “hubs” that stitch together local groups of coauthors, which may increase the sensitivity of this network to BSPC. Also recall that the degree correlation coefficient in the original bipartite network is  $r_B = -0.18$ , implying that on average authors who work in smaller collaborations tend to be more productive (this fact may reflect the nature of the dataset and its limited time frame; see Newman, 2001).

As could be expected, due to counting in non-reciprocated nominations, the non-response effect is somewhat less severe than BSP and may be tolerated for response rates of 70% and better where the relative error is less than 10% (Fig. 6, insets).

#### 4.1.2. Clustering

Random omission of actors (Fig. 8, squares) appears to have no effect on clustering in the unipartite projection. This result could be expected since all clustering is engendered via joint membership in groups, whose pattern is unaffected by random deletion of actors. It is intuitively plausible that interaction contexts are responsible for the resulting clustering and mixing pattern in the bipartite model of a social network. Fig. 8 (dots) implies that omission of contexts (BSPC) results in increased clustering. As has been mentioned above, each interaction context or group in a bipartite graph corresponds to a clique in the one-mode network of actors. If redundancy of the bipartite graph is sufficiently high, these cliques tend to overlap. As more interaction contexts are removed, cliques in the one-mode network disconnect from each other thus effectively reducing the number of connected triples of vertices  $N_3$  while keeping the number of triads  $N_\Delta$  high. This causes the clustering coefficient  $C = 3N_\Delta/N_3$  to grow.

On the contrary, non-response (Fig. 8, stars) results in lower clustering. Since missing links under non-response are the ones that connect non-responding nodes and otherwise network connectivity is not affected, this mechanism opens up triples faster than producing dyads or isolates, and therefore the clustering coefficient is decreasing.

The relative deterioration rate (Fig. 8b, inset) depends on the “true” value of clustering. For one-mode networks generated from random graphs with Poisson degree distributions,

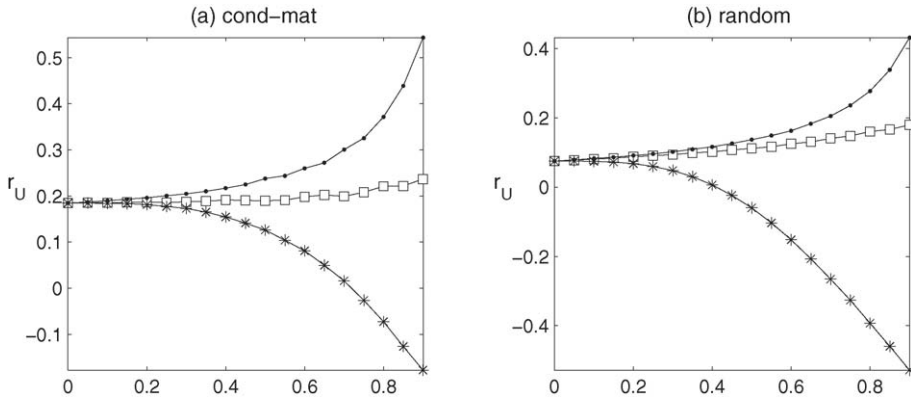


Fig. 9. Sensitivity of degree assortativity coefficient  $r_U$  in the unipartite projection: omission of interaction contexts (dots); omission of actors (squares); survey non-response (stars).

clustering coefficient changes as  $C(\theta) = 1/(1 + \mu(1 - \theta))$  in the case of BSPC, and  $C(\theta)$  is fairly close to  $\theta/(1 + \mu(1 - \theta))$  under non-response, where  $\theta$  denotes the fraction of missing groups or non-responding vertices, respectively. The first result follows trivially from the formula  $C = 1/(1 + \mu)$ , derived by Newman et al. (2001); the second is our conjecture based on simulations. It seems plausible that BSPC and non-response may compensate each other under some fortunate circumstances, yet separately they drastically affect the estimate of clustering coefficient and inflate the measurement error. Ironically, eliminating one source of error but not the other could severely impair the estimate of clustering in the network!

#### 4.1.3. Assortativity

The simulation results plotted in Fig. 9 show that, as in the case of clustering, BSPC increases degree-to-degree correlation in the unipartite projection while non-response causes it to diminish, and ultimately leads to a disassortative mixing pattern. We should emphasize these facts as they increase the uncertainty about the estimates of clustering and assortativity in networks with unknown missing data patterns.

It has been shown that unipartite networks that are assortatively mixed by degree are more robust to removal of vertices than disassortative or neutral networks (Newman, 2003b). Several social networks, including the one-mode collaboration graph analyzed in this paper have been found to be assortatively mixed. In such networks, the assortative core can form a reservoir that will sustain the disease even in the absence of epidemic in the network at large (Section 3.2). Observe, however, that one tends to overestimate the mixing coefficient in networks with multiple interaction contexts as a consequence of the boundary specification problem for contexts (Fig. 9, dots) and, to a lesser extent, BSP for actors. Therefore complete social networks may actually possess less assortativity than they appear to have, provided that researchers take measures to minimize non-response. This finding may turn out to be an important factor in cost-benefit analyses of disease prevention strategies that are based on empirical network data.

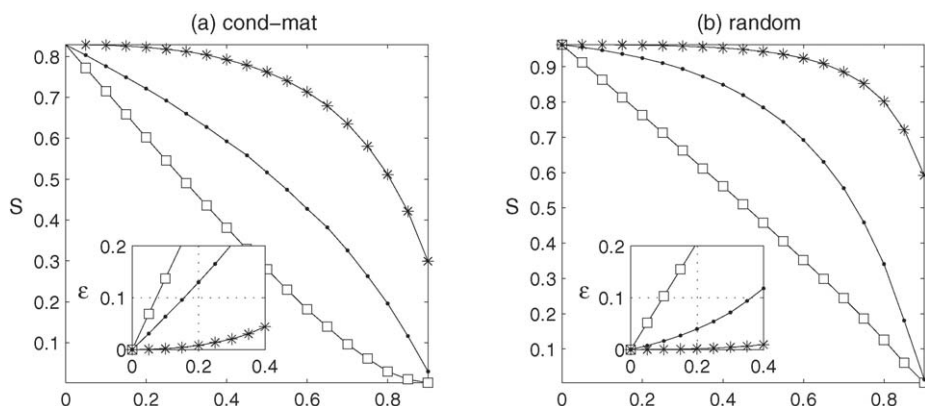


Fig. 10. Relative size of the largest connected component in the unipartite projection: omission of interaction contexts (solid dots); omission of actors (squares); survey non-response (stars).

#### 4.1.4. Size of the largest connected component

As can be seen from Fig. 10, the collaboration network is quite robust to survey non-response (stars): good estimates can be obtained with response rates of 70% and better (50% for random graphs with similar parameters). On the other hand, omission of actors (squares) leads to immediate and severe deterioration of the network connectivity. The effect of missing interaction contexts (dots) is somewhere in-between. Non-inclusion of actors (as well as actor non-response with required reciprocation, for that matter) is analogous to the so-called “node failures” analyzed in several recent studies of computer networks (Albert et al., 2000; Callaway et al., 2000; Cohen et al., 2000, 2001; Vázquez and Moreno, 2003). This line of literature has focused on the effects that random failures or intentional attacks on Internet routers might have on the global connectivity properties of the Internet, such as the size of the largest connected component. In particular, it has been shown that for random breakdowns, networks whose degree distribution is approximated by a power-law remain essentially connected even for very large breakdown rates (Cohen et al., 2000). It has been also demonstrated under quite general assumptions that disassortativity increases network fragility as it works against the process of formation of the giant component; on the other hand, assortative correlations make graph robust to random damage (Vázquez and Moreno, 2003). However, our simulation results do not fully agree with these notions. The one-mode coauthorship network is assortatively mixed and has a heavy-tailed degree distribution, while the projection of a random bipartite graph has near zero assortativity and quickly decaying degree distribution (Fig. 5a and b, respectively, dots). Yet under BSPA the size of the largest component decreases faster in the one-mode collaboration network (compare Fig. 10a and b, squares).

To separate possible effects of mixing pattern and degree distribution, we have run simulations with bipartite networks obtained by randomly rewiring the collaboration graph. These networks have the same degree sequences as the original bipartite graph but zero assortativity coefficient. The rewired networks behave very similarly to random graphs with Poisson degree distribution. An important difference, however, is that random removal of

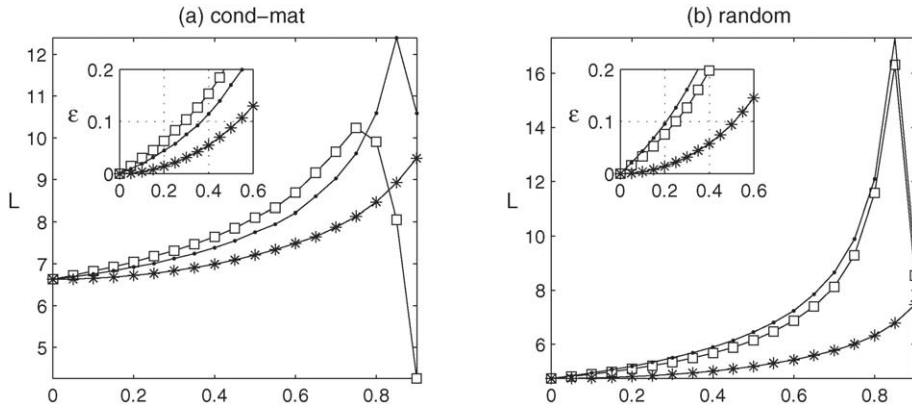


Fig. 11. Mean path length in the largest component of the unipartite projection: omission of interaction contexts (dots); omission of actors (squares); survey non-response (stars). Note the drop in path length corresponding to the loss of connectivity as the network becomes fragmented and the largest component becomes increasingly small.

actors initially leads to a faster decrease in the size of the giant component  $S_L$ , but for large removal rates  $S_L$  approaches zero size continuously in a rewired network (not shown here), while both random graph and the original collaboration network exhibit a discontinuity (easily seen in the plot of average path length, Fig. 11). We conclude that a rewired version of the collaboration graph is more resilient to BSPA than the original, despite its lack of assortativity. Hence, assortativity alone does not necessarily imply network robustness, contrary to previous assertions, and may have substantially different implications for networks engendered via joint membership in groups or interaction contexts. The compound effect of the mixing pattern and degree sequences in such networks therefore deserves a further investigation.

#### 4.1.5. Mean path length in the largest connected component

As may be seen from Fig. 11, BSPA and BSPC have a similar effect on the average path length. Path length exhibits a discontinuity when mean vertex degree becomes less than unity. Because of the skewed degree distribution of the Condensed Matter collaboration network BSPA has a stronger impact on mean degree than BSPC, and consequently, the phase transition (breakdown of the largest component into many small ones) occurs at  $\theta \approx 0.75$  for BSPA and  $\theta \approx 0.9$  for BSPC. The effects of missing data mechanisms on the mean path length may be tolerated (i.e. relative error not exceeding 10%) for amounts of missing data up to 20% in case of BSPA or BSPC, and for response rates of 50% and better in case of actor non-response.

#### 4.2. Degree censoring (fixed choice effect)

We consider the impact of fixed-choice questionnaire design (right-censoring by vertex degree) on network properties in the following three cases: (1) we record up to  $K$  interaction contexts out of average  $\mu$  for every actor; (2) each actor nominates up to  $X$  out of average



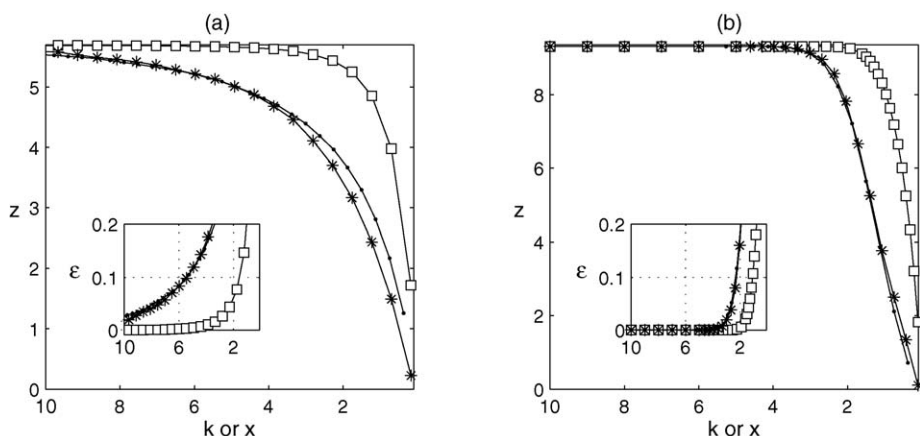


Fig. 12. Fixed choice effect on the mean degree of the unipartite projection  $z$  in the Condensed Matter collaboration graph (a) and a comparable random graph (b). Dots: censoring collaborations. The question asked of each author would be to “nominate” up to  $K$  papers coauthored by him. The horizontal axis represents the relative degree cutoff  $k = K/\mu$ , where  $\mu = 3.5$  is the mean number of affiliations per actor. Note that the amount of missing data increases as we lower the threshold value. For example,  $k = 5$  means that the actual cutoff is  $K = 5\mu$ , five times the mean actor degree in the bipartite network. Squares: censoring coauthors, no reciprocation required. The question asked of each author would be to nominate up to  $X$  coauthors. The horizontal axis represents relative degree cutoff  $x = X/z$  in units of  $z$ , the mean number of collaborators per author, where (a)  $z = 5.69$  in the Physics collaboration graph and (b)  $z = 9.31$  in a random network. Stars: only reciprocated nominations, relative cutoff  $x = X/z$  in units of  $z$ . Insets: relative error  $\epsilon = |z - z_0|/z_0$ , where  $z_0$  is the true value. Each data point is an average over 50 iterations. Lines connecting datapoints are a guide for the eye only.

$z$  interaction partners; the link is present if either one or both members of a dyad report it; (3) same as previous, but every dyadic link must be reported by both partners. Varying the cutoff values  $K$  and  $X$ , we have explored how these missing data mechanisms affect the unipartite social network under assumption of random nominations. Sensitivity curves for the mean vertex degree are shown in Fig. 12. The results for other statistics discussed in the previous sections are qualitatively similar to the corresponding BSP/non-response effects up to the direction of error (see Tables 3 and 4 for details).

It appears that degree censoring has a much more severe effect on the Condensed Matter collaboration graph (left plot) than on a random bipartite network with the same parameters  $N$ ,  $M$  and  $\mu$  (right plot). In a random graph, a fixed choice of  $K = k\mu$  interaction contexts (collaborations) or reciprocated nomination of  $X = xz$  partners practically does not affect mean degree  $z$  as long as relative cutoffs  $k > 3$  or  $x > 3$ . In the collaboration graph, however, mean degree departs from its true value as soon as the relative cutoff  $k$  or  $x$  becomes less than 15. As a consequence, this impairs estimates of such network properties as the number of components, size of the largest component and geodesics length (not shown). The effects of degree censoring on network properties are quantified in Table 4, where we report approximate minimal cutoff values such that parameter estimates are within  $\pm 10\%$  around their respective true values. It is noteworthy that fixed choice errors are virtually non-existent in random graphs for relative cutoff values  $k$  or  $x \gtrsim 2$ . On the contrary, the real collaboration network appears to be very sensitive to degree bound effects.

Table 3

Approximate tolerable fractional amount of missing data<sup>a</sup> and direction of deviation<sup>b</sup> for boundary specification and non-response effects

Property of one-mode network	Symbol	BSPC <sup>c</sup>	BSPA <sup>d</sup>	NRE <sup>e</sup>
Mean degree	$z$	0.14 (0.1) <sup>f</sup> ↓	0.1 (0.1) ↓	0.3 (0.3) ↓
Clustering	$C$	0.25 (0.1) ↑	n.a. <sup>g</sup>	0.35 (0.35) ↓
Size of largest component	$S_L$	0.15 (0.35) ↓	0.08 (0.1) ↓	n.a.
Mean path in largest component	$\ell_L$	0.4 (0.2) ↑	0.3 (0.25) ↑	0.5 ↑

<sup>a</sup> Missing data is tolerable if it causes relative error not exceeding 10%, i.e.  $\varepsilon = \left| \frac{q-q_0}{q_0} \right| \leq 0.1$ , where  $q$  is an estimate from a model with missing data and  $q_0$  is the value calculated from complete data.

<sup>b</sup> We use ↑ or ↓ to indicate the direction of departure of the estimate from the true value (up or down, respectively) for a small amount of missing data such that the network is kept above the percolation threshold, i.e. mean vertex degree  $z > 1$ .

<sup>c</sup> Boundary specification for interaction contexts or affiliations.

<sup>d</sup> Boundary specification for actors (missing actors).

<sup>e</sup> Non-response, reciprocated nominations are not required.

<sup>f</sup> Numbers in parentheses are results for an ensemble of 100 random bipartite graphs with the same number of vertices and edges.

<sup>g</sup> Very slow change: less than 10% error for 50% of missing data.

While there may be a number of different mechanisms at work, it is likely that this difference in behavior is a joint effect of the non-random mixing and skewed degree distributions observed in the Condensed Matter collaboration graph. Censoring by degree has little effect on the random graph because its degree variance is quite small, i.e. it is rather sharply peaked around the mean. Therefore, when we cut edges in excess to, say,  $2\mu$  or  $2z$  in a random graph, the number of actually removed links is negligible. On the other hand, the distribution of papers by authors and the distribution of the number of collaborators in the one-mode network both have a heavy tail (Fig. 5), i.e. there is a considerable fraction of vertices with degrees greater than twice the average value. If the one-mode network is mixed assortatively by degree as in the case of the Condensed Matter graph, then degree censoring

Table 4

Approximate minimal tolerable cutoffs<sup>a</sup> and direction of deviation<sup>b</sup> for degree censoring effects

Property (projection)	Symbol	FCC <sup>c</sup>	FCA <sup>d</sup>	FCR <sup>e</sup>
Mean degree	$z$	$5.5\mu$ (2.5) <sup>f</sup> ↓	$1.5z$ (1) ↓	$5.5z$ (2.5) ↓
Clustering	$C$	$8\mu$ (2.5) ↑	$1.5z$ (1)	$6z$ (1.6)
Size of largest component	$S_L$	$3.5\mu$ (1.2) ↓	$1z$ (0.2) ↓	$2z$ (0.7) ↓
Mean path in largest component	$\ell_L$	$6.5\mu$ (2) ↑	$1.8z$ (0.9) ↑	$5z$ (2) ↑

<sup>a</sup> The degree cutoff is tolerable if the relative error caused by censoring  $\varepsilon = \left| \frac{q-q_0}{q_0} \right| \leq 10\%$ , where  $q$  is an estimate from a model with missing data and  $q_0$  is the value calculated from complete data.

<sup>b</sup> We use ↑ or ↓, where applicable, to indicate the direction of departure of the estimate from the true value (up or down, respectively) for a small amount of missing data such that the network is kept above the percolation threshold, i.e. mean vertex degree  $z > 1$ .

<sup>c</sup> Fixed choice of interaction contexts.

<sup>d</sup> Fixed choice of actors, reciprocation not required.

<sup>e</sup> Fixed choice of actors, only reciprocated nominations.

<sup>f</sup> Numbers in parentheses are results for an ensemble of 100 random bipartite graphs with the same number of vertices and edges.

will likely eliminate most connections within the network core and quickly break down the giant component. Additional computer experiments (not shown) with a randomly rewired version of the cond-mat network, which has the same degree distribution but zero mixing, support this explanation. Whereas skewed actor degree distribution alone may have a limited impact on the robustness of network statistics with respect to the fixed choice effects, when present together with assortative mixing, it makes the network increasingly more sensitive. We would like to stress that one-mode projections of bipartite graphs, assortativity may arise as a structural artifact of a skewed group size distribution (see footnote 8), rather than being a substantive property of some network process. Hence it is important when doing empirical research that possible fixed choice effects be carefully examined if there are reasons to think that the network under study has been engendered by a multicontextual affiliation graph.

## 5. Conclusions

In this paper we have compared a number of missing data effects in social networks with multiple interaction contexts. Social interactions are modeled as a bipartite graph, consisting of the set of actors and the set of interaction contexts or affiliations. The conventional single-mode network of actors is a unipartite projection of the bipartite graph onto the set of actors. We have measured structural properties of this projection while varying the amount of missing data in the generating bipartite graph by omitting actors, interaction contexts, or individual interactions. As examples of multicontextual social networks we analyzed the Los Alamos Condensed Matter collaboration graph and an ensemble of random bipartite graphs with similar parameters.

The findings reported in this paper are based on a case study and simulated random graphs and therefore may not apply to all social networks. Moreover, we have modeled all missing data mechanisms as random, which is a big simplification. With all due limitations, however, several results of significance follow from our studies. Boundary specification can dramatically alter estimates of network-level statistics, in particular, the assortativity coefficient and mean degree, even if context redundancy is large. In a fixed choice survey design, the errors introduced by missing data are relatively small up to certain degree cutoff values, which depend on the vertex degree distribution and mixing pattern; the apparently worst case being networks with highly skewed degree distributions, which may produce unreliable statistics, especially in the presence of assortative mixing.

We find that assortativity coefficient is overestimated via omission of interaction contexts (affiliations) or fixed choice of affiliations. On the other hand, actor non-response or fixed choice of collaborators leads to an underestimated mixing coefficient and may even cause an assortatively mixed network to appear as disassortative. In a similar fashion, the observed clustering coefficient increases via omission of interaction contexts or fixed choice thereof, and decreases with actor non-response. The clustering coefficient is unaffected by random omission of actors since all clustering in the bipartite model of social networks is engendered via interaction contexts (group affiliation). The divergent effect of the two missing data mechanisms results in inflated measurement error. It is ironic that by eliminating one source of error (e.g., non-response) but not the other (boundary

specification effect) one might actually end up with worse estimates of clustering or assortativity. Finally, the confounding effect of mixing pattern and degree distribution on network robustness under random omission of actors is found to be different from what is assumed in the current literature. We have found that under certain circumstances the largest component in a network assortatively mixed by vertex degree is less robust to random deletion of vertices than in a comparable neutral network. We attribute this peculiar behavior to the detailed structural composition of the networks that we have focused on; namely, the presence of multiple overlapping cliques in the one-mode network as a result of unipartite projection.

In practice it may be difficult to estimate the effects of missing data and to identify and separate its sources. Therefore one should take measures against multiple possible missing data effects. We emphasize the importance of further research to better understand patterns and consequences of missing data in social networks and to provide statistical guidance to researchers in the field.

## Acknowledgements

The author thanks Peter Dodds, Andrew Gelman, Nobuyuki Hanaki, Catherine Hecht, Alexander Peterhansl, Duncan Watts, Harrison White and anonymous reviewers for useful comments, and Mark Newman for providing the E-print Archive collaboration data.

## References

- Albert, R., Jeong, H., Barabási, A.L., 2000. Error and attack tolerance of complex networks. *Nature* 406, 378–382.
- Amaral, L.A.N., Scala, A., Barthélemy, M., Stanley, H.E., 2000. Classes of small-world networks. *Proceeding of the National Academy of Sciences of the USA* 97, 11149–11152.
- Barabási, A.L., Albert, R., 1999. Emergence of scaling in random networks. *Science* 286, 509–512.
- Bernard, H.R., Killworth, P., Kronenfeld, D., Sailer, L., 1984. The problem of informant accuracy: the validity of retrospective data. *Annual Review of Anthropology* 13, 495–517.
- Boguñá, M., Pastor-Satorras, R., Vespignani, A., 2003. Epidemic spreading in complex networks with degree correlations. In: Rubi, J.M. (Ed.), *Proceedings of the XVIII Sitges Conference Statistical Mechanics of Complex Networks*. Springer-Verlag, Berlin.
- Bollobás, B., 2001. *Random Graphs*, 2nd ed. Cambridge University Press, Cambridge.
- Brewer, D.D., Webster, C.M., 1999. Forgetting of friends and its effects on measuring friendship networks. *Social Networks* 21, 361–373.
- Burt, R.S., 1987. A note on missing social network data in the General Social Survey. *Social Networks* 9, 63–73.
- Butts, C.T., 2003. Network inference, error, and informant (in)accuracy: a Bayesian approach. *Social Networks* 25, 103–140.
- Callaway, D.S., Newman, M.E.J., Strogatz, S.H., Watts, D.J., 2000. Network robustness and fragility: percolation on random graphs. *Physical Review Letters* 85, 5468–5471.
- Cohen, R., Erez, K., ben Avraham, D., Havlin, S., 2000. Resilience of the internet to random breakdowns. *Physical Review Letters* 85, 4626–4628.
- Cohen, R., Erez, K., ben Avraham, D., Havlin, S., 2001. Breakdown of the internet under intentional attack. *Physical Review Letters* 86, 3682–3685.
- Costenbader, E., Valente, T.W., 2003. The stability of centrality measures when networks are sampled. *Social Networks* 25, 283–307.

- Davis, G.F., Mizruchi, M.S., 1999. The money center cannot hold: commercial banks in the US system of corporate governance. *Administrative Science Quarterly* 44, 215–239.
- Ebel, H., Mielsch, L.I., Bornholdt, S., 2002. Scale-free topology of e-mail networks. *Physical Review E* 66, 035103.
- Feld, S.L., 1991. Why your friends have more friends than you do. *American Journal of Sociology* 96, 1464–1477.
- Girvan, M., Newman, M.E.J., 2002. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences of the USA* 99, 7821–7826.
- Guimera, R., Danon, L., Diaz-Guilera, A., Giralt, F., Arenas, A., 2003. Self-similar community structure in organisations. *Physical Review E* 68, 065103.
- Holland, P.W., Leinhard, S., 1973. Structural implications of measurement error in sociometry. *Journal of Mathematical Sociology* 3, 85–111.
- Holme, P., Edling, C.R., Liljeros, F., 2004. Structure and time-evolution of an Internet dating community. *Social Networks* 26, 155–174.
- Homans, G.C., 1950. *The Human Group*. Harcourt, Brace and World, New York.
- Laumann, E.O., Marsden, P.V., Prensky, D., 1983. The boundary specification problem in network analysis. In: Burt, R.S., Minor, M.J. (Eds.), *Applied Network Analysis*. Sage Publications, London, pp. 18–34.
- Little, R.J.A., Rubin, D.B., 2002. *Statistical Analysis with Missing Data*, 2nd ed. Wiley–Interscience, Hoboken, NJ.
- Marsden, P.V., 1990. Network data and measurement. *Annual Review of Sociology* 16, 435–463.
- Maslov, S., Sneppen, K., Zaliznyak, A., 2002. Detection of topological properties in complex networks: correlation profile of the Internet. <http://arxiv.org/abs/cond-mat/0205379>.
- Mizruchi, M.S., 1996. What do interlocks do? An analysis, critique, and assessment of research on interlocking directorates. *Annual Review of Sociology* 22, 271–298.
- Newman, M.E.J., 2001. The structure of scientific collaboration networks. *Proceedings of the National Academy of Sciences of the USA* 98, 404–409.
- Newman, M.E.J., 2002. Assortative mixing in networks. *Physical Review Letters* 89, 208701.
- Newman, M.E.J., 2003. Ego-centered networks and the ripple effect. *Social Networks* 25, 83–95.
- Newman, M.E.J., 2003. Mixing patterns in networks. *Physical Review E* 67, 026126.
- Newman, M.E.J., Forrest, S., Balthrop, J., 2002. Email networks and the spread of computer viruses. *Physical Review E* 66, 035101.
- Newman, M.E.J., Strogatz, S.H., Watts, D.J., 2001. Random graphs with arbitrary degree distributions and their applications. *Physical Review E* 64, 026118.
- Rapoport, A., Horvath, W.J., 1961. A study of a large sociogram. *Behavioral Science* 6, 279–291.
- Robins, G., Pattison, P., Woolcock, J., 2004. Missing data in networks: exponential random graph ( $p^*$ ) models for networks with non-respondents. *Social Networks* 26, 257–283.
- Roethlisberger, F.J., Dickson, W.J., 1939. *Management and the Worker*. Harvard University Press, Cambridge, MA.
- Rumsey, D.J., 1993. Nonresponse models for social network stochastic processes (Markov chains). Ph.D. Thesis. The Ohio State University.
- Salganik, M.J., Heckathorn, D.D., 2004. Sampling and estimation in hidden populations using respondent-driven sampling. *Sociological Methodology* 34, 193–239.
- Sampson, S., 1969. *Crisis in a cloister*. Ph.D. Thesis. Cornell University.
- Stork, D., Richards, W.D., 1992. Nonrespondents in communication network studies: problems and possibilities. *Group and Organization Management* 17 (2), 193–209.
- Strogatz, S.H., 2001. Exploring complex networks. *Nature* 410, 268–276.
- Vázquez, A., Moreno, Y., 2003. Resilience to damage of graphs with degree correlations. *Physical Review E* 67, 015101.
- Wasserman, S., Faust, K., 1994. *Social Network Analysis: Methods and Applications*. Cambridge University Press, Cambridge.
- Watts, D.J., 1999. *Small Worlds: The Dynamics of Networks Between Order and Randomness*. Princeton University Press, Princeton, NJ.
- Watts, D.J., Dodds, P.S., Newman, M.E.J., 2002. Identity and search in social networks. *Science* 296, 1302–1305.

- Watts, D.J., Strogatz, S.H., 1998. Collective dynamics of 'small-world' networks. *Nature* 393, 440–442.
- White, H.C., 1992. *Identity and Control: A Structural Theory of Social Action*. Princeton University Press, Princeton, NJ.
- White, H.C., Boorman, S.A., Breiger, R.L., 1976. Social structure from multiple networks. I. Blockmodels of roles and positions. *American Journal of Sociology* 81, 730–780.
- Wilson, T.P., 1982. Relational networks: an extension of sociometric concepts. *Social Networks* 4, 105–116.