

# UCLA

## UCLA Previously Published Works

### Title

Modeling social networks from sampled data

### Permalink

<https://escholarship.org/uc/item/32b0r2sg>

### Journal

The Annals of Applied Statistics, 4(1)

### ISSN

1932-6157

### Authors

Handcock, Mark S  
Gile, Krista J

### Publication Date

2010

### DOI

10.1214/08-aos221

Peer reviewed

## MODELING SOCIAL NETWORKS FROM SAMPLED DATA<sup>1</sup>

BY MARK S. HANDCOCK AND KRISTA J. GILE

*University of California–Los Angeles and Nuffield College*

Network models are widely used to represent relational information among interacting units and the structural implications of these relations. Recently, social network studies have focused a great deal of attention on random graph models of networks whose nodes represent individual social actors and whose edges represent a specified relationship between the actors.

Most inference for social network models assumes that the presence or absence of all possible links is observed, that the information is completely reliable, and that there are no measurement (e.g., recording) errors. This is clearly not true in practice, as much network data is collected through sample surveys. In addition even if a census of a population is attempted, individuals and links between individuals are missed (i.e., do not appear in the recorded data).

In this paper we develop the conceptual and computational theory for inference based on sampled network information. We first review forms of network sampling designs used in practice. We consider inference from the likelihood framework, and develop a typology of network data that reflects their treatment within this frame. We then develop inference for social network models based on information from adaptive network designs.

We motivate and illustrate these ideas by analyzing the effect of link-tracing sampling designs on a collaboration network.

**1. Introduction.** Networks are a useful device to represent “relational data,” that is, data with properties beyond the attributes of the individuals (nodes) involved. Relational data arise in many fields and network models are a natural approach to representing the patterns of the relations between nodes. Networks can be used to describe such diverse ideas as the behavior of epidemics, the interconnectedness of corporate boards, and networks

---

Received December 2007; revised September 2008.

<sup>1</sup>Supported by NIH Grants R01 DA012831 and R01 HD041877, NSF Grant MMS-0851555 and ONR Grant N00014-08-1-1015.

*Key words and phrases.* Exponential family random graph model,  $p^*$  model, Markov chain Monte Carlo, design-based inference.

<p>This is an electronic reprint of the original article published by the <a href="#">Institute of Mathematical Statistics</a> in <i>The Annals of Applied Statistics</i>, 2010, Vol. 4, No. 1, 5–25. This reprint differs from the original in pagination and typographic detail.</p>
--

of genetic regulatory interactions. In social network applications, the nodes in a graph typically represent individuals, and the ties (edges) represent a specified relationship between individuals. Nodes can also be used to represent larger social units (groups, families, organizations), objects (airports, servers, locations) or abstract entities (concepts, texts, tasks, random variables). We consider here stochastic models for such graphs. These models attempt to represent the stochastic mechanisms that produce relational ties, and the complex dependencies thus induced.

Social network data typically consist of a set of  $n$  actors and a relational tie random variable,  $Y_{ij}$ , measured on each possible ordered pair of actors,  $(i, j)$ ,  $i, j = 1, \dots, n, i \neq j$ . In the most simple cases,  $Y_{ij}$  is a dichotomous variable, indicating the presence or absence of some relation of interest, such as friendship, collaboration, transmission of information or disease, etc. The data are often represented by an  $n \times n$  sociomatrix  $Y$ , with diagonal elements, representing self-ties, treated as structural zeros. In the case of binary relations, the data can also be thought of as a graph in which the nodes are actors and the edge set is  $\{(i, j) : Y_{ij} = 1\}$ . For many networks the relations are undirected in the sense that  $Y_{ij} = Y_{ji}, i, j = 1, \dots, n$ .

In the application in this paper we consider a network formed from the collaborative working relations between  $n = 36$  partners in a New England law firm [Lazega (2001)]. We focus on the undirected relation where a tie is said to exist between two partners if and only if both indicate that they collaborate with the other. The scientific objective is to explain the observed structural pattern of collaborative ties as a function of nodal and relational attributes. The relational data is supplemented by four actor attributes: seniority (the rank number of chronological entry into the firm divided by 36), practice (there are two possible values, litigation = 0 and corporate law = 1), gender (3 of the 36 lawyers are female) and office (there are three different offices in three different cities each of different size).

For large or hard-to-find populations of actors it is difficult to obtain information on all actors and all relational ties. As a result, various survey sampling strategies and methods are applied. Some of these methods make use of network information revealed by earlier stages of sampling to guide later sampling. These adaptive designs allow for more efficient sampling than conventional sampling designs. We consider such designs in Section 2.

Most of the work presented here considers the network over the set of actors to be the realization of a stochastic process. We seek to model that process. An alternative is to view the network as a fixed structure about which we wish to make inference based on partial observation.

In this paper we develop a theoretical framework for inference from network data that are partially-observed due to sampling. This work extends the fundamental work of Thompson and Frank (2000). For purposes of presentation, we focus on the relational data itself and suppress reference to

covariates of the nodes. This more general situation is dealt with in [Handcock and Gile \(2007\)](#).

In [Section 2](#) we present a conceptual framework for network sampling. We extend this framework in [Section 3](#) to focus on inference from sampled network data. We first consider the limitations of design-based inference in this setting, then focus on likelihood-based inference. [Section 4](#) presents the rich Exponential Family Random Graph Model (ERGM) family of models that has been applied to complete network data. [Section 5](#) presents a study of the effect of sampling from a known complete network of law firm collaborations. Finally, in [Section 6](#), we discuss the overall ramifications for the modeling of social networks with sampled data and note some extensions.

**2. Network sampling design.** In this section we consider the conceptual and computational theory of network sampling.

There is a substantial literature on network sampling designs. Our development here follows [Thompson and Seber \(1996\)](#) and [Thompson and Frank \(2000\)](#). Let  $\mathcal{Y}$  denote the set of possible networks on the  $n$  actors. Note that in most network samples, the unit of sampling is the actor or node, while the unit of analysis is typically the dyad. Let  $D$  be the  $n \times n$  random binary matrix indicating if the corresponding element of  $Y$  was sampled or not. The value of the  $i, j$ th element is 0 if the  $(i, j)$  ordered pair was not sampled and 1 if the element was sampled. Denote the sample space of  $D$  by  $\mathcal{D}$ . We shall refer to the probability distribution of  $D$  as the *sampling design*. The sampling design is often related to the structure of the graph and a parameter  $\psi \in \Psi$ , so we posit a model for it. Specifically, let  $P(D = d|Y = y; \psi)$  denote the probability of selecting sample  $d$  given a network  $y$  and parameter  $\psi$ .

Under many sampling designs the set of sampled dyads is determined by the set of sampled nodes. Let  $S$  represent a binary random  $n$ -vector indicating a subset of the nodes, where the  $i$ th element is 1 if the  $i$ th node is part of the set, and is 0 otherwise. We often consider situations where  $D$  is determined by some  $S$  which is itself a result of a sample design denoted by  $P(S|Y, \psi)$ . For example, consider an undirected network where the set of observed dyads are those that are incident on at least one of the sampled nodes. In this case  $D = S \circ 1 + 1 \circ S - S \circ S$ , where  $1$  is the binary  $n$ -vector of 1s. A primary example of this is where people are sampled and surveyed to determine all their edges.

We introduce further notation to allow us to refer to the observed and unobserved portions of the relational structures. Denote the observed part of the complete graph  $Y$  by  $Y_{\text{obs}} = \{Y_{ij} : D_{ij} = 1\}$  and the unobserved part by  $Y_{\text{mis}} = \{Y_{ij} : D_{ij} = 0\}$ . The full *observed data* is then  $\{Y_{\text{obs}}, D\}$ , in contrast to the *complete data*:  $\{Y_{\text{obs}}, Y_{\text{mis}}, D\}$ . We will write the complete graph  $Y = \{Y_{\text{obs}}, Y_{\text{mis}}\}$ . In addition, we make the convention that undefined numbers act as identity elements in addition and multiplication. So a number  $x$  plus

or multiplied by an undefined number  $y$  is  $x$ , and hence  $Y = Y_{\text{obs}} + Y_{\text{mis}}$ . For a given network  $y \in \mathcal{Y}$ , denote the corresponding data as  $\{y_{\text{obs}}, d\}$  and the other elements by their lower-case versions  $y = y_{\text{obs}} + y_{\text{mis}}$ . Finally denote  $\mathcal{Y}(y_{\text{obs}}) = \{v : y_{\text{obs}} + v \in \mathcal{Y}\}$ , that is the set of possible unobserved elements which together with  $y_{\text{obs}}$  result in valid network. The set  $y_{\text{obs}} + \mathcal{Y}(y_{\text{obs}})$  is then the restriction of  $\mathcal{Y}$  to  $y_{\text{obs}}$ .

A sampling design is *conventional* if it does not use information collected during the survey to direct subsequent sampling of individuals (e.g., network census and ego-centric designs). Specifically, a design is conventional if  $P(D = d | Y = y; \psi) = P(D = d | \psi) \forall y \in \mathcal{Y}$ . A simple example of a conventional sampling design for networks is an *ego-centric design*, consisting of a simple random sampling of a subset of the actors, followed by complete observation of the dyads originating from those actors. A complete census of the network is another. More complex examples include designs using probability sampling of pairs and auxiliary variables. Alternatively, we call a sampling design *adaptive* if it uses information collected during the survey to direct subsequent sampling, but the sampling design depends only on the observed data. Specifically, a design is adaptive if:  $P(D = d | Y = y; \psi) = P(D = d | Y_{\text{obs}} = y_{\text{obs}}, \psi) \forall y \in y_{\text{obs}} + \mathcal{Y}(y_{\text{obs}})$ . Hence a design can be adaptive for a given  $y_{\text{obs}}$  (rather than all possible observed data), although most common such designs are adaptive for all possible data observed under them. Conventional designs can be considered to be special cases of adaptive designs.

Note that adaptive sampling designs satisfy

$$(2.1) \quad P(D = d | Y_{\text{obs}}, Y_{\text{mis}}, \psi) = P(D = d | Y_{\text{obs}}, \psi),$$

a condition called “missing at random” by [Rubin \(1976\)](#) in the context of missing data. Note that this is a bit misleading—it does not say that the propensity to be observed is unrelated to the unobserved portions of the network, but that this relationship can be explained by the data that are observed. The observed part of the data are often vital to equality (2.1). Hence adaptive designs are essentially those for which the unobserved dyads are missing at random.

Denote by  $[a]$  the vector-valued function that is 1 if the corresponding element of the vector  $a$  is logically true, and 0 otherwise. Let  $a \times b$  be the elementwise product of the column vector  $a$  and the column vector  $b$  and  $a \cdot b$  be the scalar product  $\sum_j a_j b_j$ . Let  $a \circ b$  be the outer product matrix with  $ij$ th element  $a_i b_j$ . If  $y$  is a matrix and  $b$  a vector let  $y \cdot b$  be the column vector with  $i$ th element  $\sum_j y_{ji} b_j$ .

2.1. *Some adaptive designs for undirected networks.* We now consider several examples of adaptive designs for undirected networks.

2.1.1. *Example: Ego-centric design.* Consider a simple *ego-centric design*:

1. Select individuals at random, each with probability  $\psi$ .
2. Observe all dyads involving the selected individuals (i.e., dyads with at least one of the selected individuals as one of the pair of actors).

The sampling design can be determined for this case. First note that

$$P(D_{ij} = 1|Y, \psi) = 1 - (1 - \psi)^2 \quad \forall i \neq j.$$

This, however, does not give the joint distribution of  $D$ . Let  $S$  be the binary  $n$ -vector where 1 and 0 indicate that the corresponding individual has been selected, or not, respectively. Within this design,  $S$  is determined by  $D$  (i.e.,  $S = [D1 = (n-1)1]$ ). Then  $P(S = s|Y, \psi) = \psi^{1 \cdot s} (1 - \psi)^{n-1 \cdot s}$ ,  $s \in \{0, 1\}^n$ . If the  $i$ th element of  $S$  is 1 then all elements in the  $i$ th row and column of  $D$  are 1.  $D_{ij} = 0$  if and only if both the  $i$ th and  $j$ th elements of  $S$  are both 0. Hence the probability distribution of  $D$  is

$$P(D = d|Y, \psi) = \psi^{1 \cdot s} (1 - \psi)^{n-1 \cdot s}$$

for

$$d = 1 \circ s + s \circ 1 - s \circ s, \quad s \in \{0, 1\}^n.$$

Note that the distribution does not depend on  $Y$ , and is therefore conventional.

2.1.2. *Example: One-wave link-tracing design.* We refer to any sample in which subsequent nodes are enrolled based on their observed relations with other sampled nodes as a *link-tracing design*. Consider the one-wave link-tracing design specified as follows:

1. Select individuals at random, each with probability  $\psi$ .
2. Observe all dyads involving the selected individuals.
3. Identify all individuals reported to have at least one relation with the initial sample, and select them with probability 1.
4. Observe all dyads involving the newly selected individuals.

Let  $S_0$  denote the indicator vector for the initial sample and  $S_1$  the indicator for the added individuals not in the initial sample. Then the whole sample of individuals is  $S = S_0 + S_1$ . As in the undirected ego-centric design,  $D = 1 \circ S + S \circ 1 - S \circ S$ . Note that  $S_1 = [Y S_0 \times (1 - S_0) > 0]$  is derivable from  $S_0$  and  $Y$ . Hence

$$P(D = d|Y, \psi) = \sum_{s_0: s_0 + [Y s_0 \times (1 - s_0) > 0] = s} \psi^{1 \cdot s_0} (1 - \psi)^{n-1 \cdot s_0}$$

for

$$d = 1 \circ s + s \circ 1 - s \circ s, \quad s \in \{0, 1\}^n.$$

2.1.3. *Example: Multi-wave link-tracing design.* Consider a *multi-wave link-tracing design* in which the complete set of partners of the  $k$ th wave are enrolled, that is, the link-tracing process described above is carried out  $k$  times. If  $k$  is fixed in advance this is called  *$k$ -wave link-tracing*.

Let  $S_0$  denote the indicator for the initial sample,  $S_1$  the indicator for the added individuals in the first wave not in the initial sample,  $\dots$ ,  $S_k$  the indicator for the added individuals in wave  $k$  not in the prior samples. Then the whole sample of individuals is  $S = S_0 + S_1 + \dots + S_k$ . As in the ego-centric design  $D = 1 \circ S + S \circ 1 - S \circ S$ . Note that  $S_m = [Y S_{m-1} \times (1 - \sum_{t=0}^{m-1} S_t) > 0]$ ,  $m = 1, \dots, k$  is derivable from  $S_0$  and  $Y$ . Then

$$P(D = d|Y, \psi) = \sum_{s_0: s_0 + s_1 + \dots + s_k = s} \psi^{1 \cdot s_0} (1 - \psi)^{n - 1 \cdot s_0}$$

for  $d = 1 \circ s + s \circ 1 - s \circ s$ ,  $s \in \{0, 1\}^n$ . Here  $S_m = [Y S_{m-1} \times (1 - \sum_{t=0}^{m-1} S_t) > 0] = [Y_{\text{obs}} S_{m-1} \times (1 - \sum_{t=0}^{m-1} S_t) > 0]$ ,  $m = 1, \dots, k$  so that the individuals selected in the successive waves only depend on the observed part of the graph, and not on the unobserved portions of the graph. Clearly, this is also true for one-wave link-tracing as a simple case of  $k$ -wave link-tracing. Note that it may be possible that  $S_m = \emptyset$  for some  $m < k$ , so that subsequent waves do not increase the sample size (i.e.,  $S_k = \emptyset$ ). A variant of the  $k$ -wave link-tracing design is the *saturated link-tracing* design, in which sampling continues until wave  $m$ , such that  $S_m = \emptyset$ . We interpret  $k$  as the bound on the number of waves sampled imposed by the sampling design. Since saturated link-tracing does not restrict the number of waves sampled, we represent it by setting  $k = \infty$ .

2.2. *Some adaptive designs for directed networks.* We can also consider variants of these adaptive designs for directed networks.

2.2.1. *Example: Ego-centric design.* Consider a simple *ego-centric design*:

1. Select individuals at random, each with probability  $\psi$ .
2. Observe all directed dyads originating at the selected individuals.

As before, the sampling design can be determined for this case. Since a directed dyad is observed only if its tail node is sampled,

$$P(D_{ij} = 1|Y, \psi) = \psi \quad \forall i \neq j$$

and  $D = S_0 \circ 1$ . Hence the probability distribution of  $D$  is

$$P(D = d|Y, \psi) = \psi^{1 \cdot s} (1 - \psi)^{n - 1 \cdot s}$$

for  $d = s \circ 1$ ,  $s \in \{0, 1\}^n$  and the distribution does not depend on  $Y$ . As in the undirected case, this design is therefore conventional.

2.2.2. *Example: One-wave link-tracing design.* Consider a one-wave link-tracing design on a directed network specified as follows:

1. Select individuals at random, each with probability  $\psi$ .
2. Observe all directed dyads originating at the selected individuals.
3. Identify all individuals receiving an arc from a member of the initial sample, and select them with probability 1.
4. Observe all directed dyads originating at the newly selected individuals.

Let  $S_0$  denote the indicator vector for the initial sample and  $S_1$  the indicator for the added individuals not in the initial sample. Then the whole sample of individuals is  $S = S_0 + S_1$ . As in the ego-centric design  $D = S \circ 1$  and

$$P(D = d|Y, \psi) = \sum_{s_0: s_0 + [Y s_0 \times (1 - s_0)] > 0] = s} \psi^{1 \cdot s_0} (1 - \psi)^{n - 1 \cdot s_0}$$

for  $d = s \circ 1, s \in \{0, 1\}^n$ .

2.2.3. *Example: Multi-wave link-tracing design.* Consider a directed version of the multi-wave link-tracing design in which the complete set of out-partners of the  $k$ th wave are enrolled. The whole sample of individuals is  $S = S_0 + S_1 + \dots + S_k$ . And  $S_m = [Y \cdot S_{m-1} \times (1 - \sum_{t=0}^{m-1} S_t) > 0]$ ,  $m = 1, \dots, k$  is derivable from  $S_0$  and  $Y$ . Then

$$P(D = d|Y, \psi) = \sum_{s_0: s_0 + s_1 + \dots + s_k = s} \psi^{1 \cdot s_0} (1 - \psi)^{n - 1 \cdot s_0}$$

for  $d = s \circ 1, s \in \{0, 1\}^n$ , where we note that  $S_m = [Y \cdot S_{m-1} \times (1 - \sum_{t=0}^{m-1} S_t) > 0] = [Y_{\text{obs}} \cdot S_{m-1} \times (1 - \sum_{t=0}^{m-1} S_t) > 0]$ ,  $m = 1, \dots, k$  so that the individuals selected in successive waves of depend only on the previously observed part of the graph, and not on the unobserved portions. The saturated link-tracing design is represented by  $k = \infty$ .

**3. Inferential frameworks.** In this section we consider two frameworks for inference based on sampled data. In the *design-based* framework  $y$  represents the fixed population and interest focuses on characterizing  $y$  based on partial observation. The random variation considered is due to the sampling design alone. A key advantage of this approach is that it does not require a model for the data themselves, although a model may also be used to guide design-based inference [Särndal, Swensson and Wretman (1992)]. Under the *model-based* framework,  $Y$  is stochastic and is a realization from a stochastic process depending on a parameter  $\eta$ . Here interest focuses on  $\eta$  which characterizes the mechanism that produced the complete network  $Y$ . We find severe limitations of the design-based framework for data from link-tracing samples, and focus on likelihood inference within the model-based framework.



3.1. *Design-based inference for the network.* In the design-based frame, the unobserved data values, or some functions thereof, are analogous to the parameters of interest in likelihood inference. The population of data values is treated as fixed, and all uncertainty in the estimates is due to the sampling design, which is typically assumed to be fully known (not just up to the parameter  $\psi$ ).

Inference typically focuses on identifying design-unbiased estimators for quantities of interest measured on the complete network. In an undirected network analysis setting, for example, we can consider estimating  $\tau = \sum_{i < j} y_{ij}$ , the number of edges in the network. Note that  $y$  is a partially-observed matrix of constants in this setting. Then  $\hat{\tau}$  is design-unbiased for  $\tau$  if

$$\mathbb{E}_D[\hat{\tau}|\psi, y] = \tau,$$

where the expectation is taken over realizations of the sampling process. Specifically,

$$\mathbb{E}_D[\hat{\tau}(Y_{\text{obs}}, D)|\psi, y] = \sum_{d \in \mathcal{D}} \hat{\tau}(y_{\text{obs}}(d), d)P(D = d|\psi, y),$$

where  $\hat{\tau}(y_{\text{obs}}(d), d)$  is the estimator expressed as a function of the observed network information. Similarly, the variance of the estimator is computed with respect to the variation induced by the sampling procedure

$$\mathbb{V}_D[\hat{\tau}(Y_{\text{obs}}, D)|\psi, y] = \sum_{d \in \mathcal{D}} (\hat{\tau}(y_{\text{obs}}(d), d) - \tau)^2 P(D = d|\psi, y).$$

The Horvitz–Thompson estimator is a classic tool of design-based inference, and is based on inverse-probability weighting the sample. In our example, it is

$$\hat{\tau}(Y_{\text{obs}}, D) = \sum_{i < j: D_{ij}=1} \frac{y_{ij}}{\pi_{ij}},$$

where the *dyadic sampling probability*  $\pi_{ij} = P(D_{ij} = 1|\psi, y)$  is the probability of observing dyad  $(i, j)$ .

Consider an estimator of  $\tau$  based on relations observed through the egocentric design of Section 2.1.1. Then

$$\pi_{ij} = 1 - (1 - \psi)^2 \quad \forall i, j.$$

The classic Horvitz–Thompson estimator  $\hat{\tau}$  of  $\tau$  then weights each observation by the inverse of its sampling probability

$$\hat{\tau} = \sum_{i < j: D_{ij}=1} \frac{y_{ij}}{\pi_{ij}} = \frac{1}{1 - (1 - \psi)^2} \sum_{i < j: D_{ij}=1} y_{ij}.$$

Then

$$\mathbb{V}(\hat{\tau}) = \sum_{i < j} \sum_{k < l} \{[1 - (1 - \psi)^2]^{-2} \pi_{ij,kl} - 1\} y_{ij} y_{kl},$$

where  $\pi_{ij,kl} = P(S_{0i} + S_{0j} > 0, S_{0k} + S_{0l} > 0)$  or

$$\pi_{ij,kl} = \begin{cases} \pi_{ij}, & i = k, j = l, \\ \pi_{ij} \pi_{kl}, & i \notin \{k, l\} \text{ and } j \notin \{k, l\}, \\ \psi^3 - 3\psi^2, & \text{otherwise.} \end{cases}$$

Among the many available estimators for the variance of the Horvitz–Thompson estimator is the Horvitz–Thompson variance estimator:

$$\hat{\mathbb{V}}(\hat{\tau}) = \sum_{i < j: D_{ij}=1} \sum_{k < l: D_{kl}=1} \frac{1}{\pi_{ij,kl}} \{[1 - (1 - \psi)^2]^{-2} \pi_{ij,kl} - 1\} y_{ij} y_{kl}.$$

Note the importance of the unit sampling probabilities in these estimators. This is a hallmark of design-based inference: inference relies on full knowledge of the sampling procedure in order to make unbiased inference without making assumptions about the distribution of the unobserved data. This typically requires knowledge of the sampling probability of each unit in the sample. This procedure is complicated in the network context, in that we require the sampling probabilities of the units of analysis, dyads, which are different from the units of sampling, nodes. In fact, for even single-wave link-tracing samples, the dyadic sampling probabilities are not observable.

To see this, define the *nodal neighborhood of a dyad*  $(i, j)$ ,  $N(i, j)$ , where  $k \in N(i, j) \iff \{S_{0k} = 1 \implies D_{ij} = 1\}$ . Then  $\pi_{ij} = P(\exists k: S_{0k} = 1, k \in N(i, j))$ .

For the one-wave link-tracing design of Section 2.1.2,  $N(i, j) = \{k\}: y_{ik} = 1$  or  $y_{jk} = 1$  or  $k \in \{i, j\}$ . Then if the initial sample  $S_0$  is drawn according to the design in Section 2.1.2,  $\pi_{ij} = 1 - (1 - \psi)^{\|N(i, j)\|}$ . Suppose  $S_{0i} = 1$ , and  $S_{0j} = 0$ . Then dyad  $(i, j)$  is observed, but  $\|N(i, j)\|$  is unknown because it is unknown which  $k$  satisfy  $y_{jk} = 1$ . The link-tracing sampling structures for which nodal and dyadic sampling probabilities are observable are summarized in Table 1. For directed networks, we assume sampled nodes provide information on their out-arcs only, so that  $D$  is not symmetric and  $D_{ij} = 1 \iff S_i = 1$ .

Of the designs considered here, dyadic sampling probabilities are observable only for ego-centric samples, and never for link-tracing designs. Nodal sampling probabilities are also observable for ego-centric sampling, as well as for one-wave and saturated link-tracing designs in undirected networks. Overall, this table presents strong limitations to the applicability of design-based methods requiring the knowledge of sampling probabilities to link-tracing designs. Note that this limitation is not specific to dyad-based network statistics. Estimation of triad-based network statistics such as a triad

TABLE 1

*Observable sampling probabilities under various sampling schemes for directed and undirected networks. Nodal and dyadic sampling probabilities are considered separately. “X” indicates observable sampling probabilities, while a blank indicates unobservable sampling probabilities*

Sampling scheme	Nodal probabilities $\pi_i$		Dyadic probabilities $\pi_{ij}$	
	Undirected	Directed	Undirected	Directed
Ego-centric	X	X	X	X
One-wave	X			
$k$ -wave, $1 < k < \infty$				
saturated	X			

census would be subject to similar limitations. A Horvitz–Thompson style estimator would rely on a weighted sum of observed triads, weighted according to sampling probabilities. Sampling probabilities for triads would be even more complex, as they would typically require sampling of two of the three nodes involved in an undirected case, and at least two of the three nodes in an directed case, depending on the triad census. Both of these sampling probabilities would not be possible to compute for link-tracing samples in which the degrees or in-degrees of some involved nodes are unobserved.

Not surprisingly, most of the work on design-based estimators for link-tracing samples has focused on the cases where sampling probabilities are observable: typically for one-wave or saturated samples used to estimate population means of nodal covariates. Frank (2005) presents a good overview and extensive citations to this literature. See also Thompson and Collins (2002); Snijders (1992). Although examples tend to focus on instances where sampling probabilities are observable, the limited applicability of classical design-based methods in estimating structural network features based on link-tracing samples has not been emphasized in the literature.

In the absence of observable sampling probabilities, design-based inference requires a mechanism for estimating sampling probabilities. This is most often necessary in the context of out-of-design missing data, and addressed with approaches such as propensity scoring [Rosenbaum and Rubin (1983)], which rely on auxiliary information available for the full sampling frame to estimate unknown sampling probabilities. Link-tracing differs from the traditional context of such methods in that the sampling probabilities are unobserved even when the design is executed faithfully, and in that the unknown sampling probabilities result directly from the unobserved variable of interest. In particular, estimating unknown sampling probabilities is equivalent to estimating unobserved relations based on the observed relations. One approach is to augment the sample with sufficient information to allow for determination of the sampling probabilities. However in most

cases, this requires a substantial expansion of the sampling design. Therefore, in practice we must rely on a model relating the observed portions of the network structure to the unobserved portions. Lack of reliance on an assumed outcome model is a great advantage of the design-based framework over the model-based framework. By introducing a model to estimate sampling probabilities based on the outcome of interest, we reintroduce this reliance on model form, negating much of the advantage of the design-based framework. Furthermore, note that the naive use of this approach has an ad-hoc flavor, while still requiring complex observation weights and variance estimators.

In the next section, we describe an alternative more flexible likelihood approach to network inference based on link-tracing samples.

*3.2. Likelihood-based inference.* Consider a parametric model for the random behavior of  $Y$  depending on a parameter  $p$ -vector  $\eta$ :

$$(3.1) \quad P_\eta(Y = y), \quad \eta \in \Xi.$$

In the model-based framework, if  $Y$  is completely observed, inference for  $\eta$  can be based on the likelihood

$$L[\eta|Y = y] \propto P_\eta(Y = y).$$

This situation has been considered in detail in [Hunter and Hancock \(2006\)](#) and the references therein. In the general case, where  $Y$  may be only partially observed, we can consider using the (so-called) *face-value likelihood* based solely on  $Y_{\text{obs}}$ :

$$(3.2) \quad L[\eta|Y_{\text{obs}} = y_{\text{obs}}] \propto \sum_{v \in \mathcal{Y}(y_{\text{obs}})} P_\eta(Y = y_{\text{obs}} + v).$$

This ignores the additional information about  $\eta$  available in  $D$ . Inference for  $\eta$  and  $\psi$  should be based on all the available observed data, including the sampling design information. This likelihood is any function of  $\eta$  and  $\psi$  proportional to  $P(D, Y_{\text{obs}}|\eta, \psi)$ :

$$\begin{aligned} L[\eta, \psi|Y_{\text{obs}} = y_{\text{obs}}, D = d_{\text{obs}}] & \\ & \propto P(D = d_{\text{obs}}, Y_{\text{obs}} = y_{\text{obs}}|\eta, \psi) \\ & = \sum_{v \in \mathcal{Y}(y_{\text{obs}})} P(D = d_{\text{obs}}|Y = y_{\text{obs}} + v, \psi) P_\eta(Y = y_{\text{obs}} + v). \end{aligned}$$

Thus the correct model is related to the complete data model through the sampling design as well as the observed nodes and dyads.

In likelihood inference, the sampling parameter  $\psi$  is a nuisance parameter, and modeling the sampling design along with the data structure adds a great

deal of complexity. It is natural to ask when we might consider the simpler face-value likelihood, (3.2), which ignores the sampling design.

In the context of missing data, Rubin (1976) introduced the concept of *ignorability* to specify when inference based on the face-value likelihood is efficient. We introduce the term *amenability* to represent the notion of ignorability for network sampling strategies within a likelihood framework.

In many situations where models are used, the parameters  $\eta \in \Xi$  and  $\psi \in \Psi$  are *distinct*, in the sense that the joint parameter space of  $(\eta, \psi)$  is  $\Psi \times \Xi$ . If the sampling design is adaptive and the parameters  $\eta$  and  $\psi$  are distinct,

$$\begin{aligned} L[\eta, \psi | Y_{\text{obs}} = y_{\text{obs}}, D = d_{\text{obs}}] \\ &\propto P(D = d_{\text{obs}} | Y_{\text{obs}} = y_{\text{obs}}, \psi) \sum_{v \in \mathcal{Y}(y_{\text{obs}})} P_{\eta}(Y = y_{\text{obs}} + v) \\ &\propto L[\psi | D = d_{\text{obs}}, Y_{\text{obs}} = y_{\text{obs}}] \times L[\eta | Y_{\text{obs}} = y_{\text{obs}}]. \end{aligned}$$

Thus if the sampling design is adaptive and the structural and sampling parameters are distinct, then the sampling design is *ignorable* in the sense that the resulting likelihoods are proportional. When this condition is satisfied likelihood-based inference for  $\eta$ , as proposed here, is unaffected by the (possibly unknown) sampling design. This leads to the following definition and result.

**DEFINITION.** Consider a sampling design governed by parameter  $\psi \in \Psi$  and a stochastic network model  $P_{\eta}(Y = y)$  governed by parameter  $\eta \in \Xi$ . We call the sampling design *amenable to the model* if the sampling design is adaptive and the parameters  $\psi$  and  $\eta$  are distinct.

**RESULT.** Consider networks produced by the stochastic network model  $P_{\eta}(Y = y)$  governed by parameter  $\eta \in \Xi$  which are observed by a sampling design with parameter  $\psi \in \Psi$  amenable to the model. Then the likelihood for  $\eta$  and  $\psi$  is

$$L[\eta, \psi | Y_{\text{obs}} = y_{\text{obs}}, D = d_{\text{obs}}] \propto L[\psi | D = d_{\text{obs}}, Y_{\text{obs}} = y_{\text{obs}}] \times L[\eta | Y_{\text{obs}} = y_{\text{obs}}].$$

Thus likelihood-based inference for  $\eta$  from  $L[\eta, \psi | Y_{\text{obs}}, D]$  will be the same as likelihood-based inference for  $\eta$  based on  $L[\eta | Y_{\text{obs}}]$ .

This result shows for standard designs such as the ego-centric, single wave and multi-wave sampling designs in Section 2, likelihood-based inference can be based on the face-value likelihood  $L[\eta | Y_{\text{obs}}]$ . This was first noted in the foundational paper of Thompson and Frank (2000). Explicitly, this is

$$L[\eta | Y_{\text{obs}} = y_{\text{obs}}] \propto P(Y_{\text{obs}} = y_{\text{obs}} | \eta) = \sum_{v \in \mathcal{Y}(y_{\text{obs}})} P_{\eta}(Y = y_{\text{obs}} + v).$$

Hence we can evaluate the likelihood by just enumerating the full data likelihood over all possible values for the missing data.

We may also wish to make inference about the design parameter  $\psi$ . The likelihood for  $\psi$  based on the observed data is any function of  $\psi$  proportional to  $P(D, Y_{\text{obs}}|\psi)$ . For designs amenable to the model this is

$$\begin{aligned} L[\psi|D = d_{\text{obs}}, Y_{\text{obs}} = y_{\text{obs}}] &\propto P(D = d_{\text{obs}}|Y_{\text{obs}} = y_{\text{obs}}, \psi) \\ &= P(D = d_{\text{obs}}|Y = y_{\text{obs}} + v, \psi) \end{aligned}$$

for any choice of  $v$  in  $\mathcal{Y}(y_{\text{obs}})$ . Hence it can be computed without reference to the network model.

**4. Exponential family models for networks.** The models we consider for the random behavior of  $Y$  rely on a  $p$ -vector  $g(Y)$  of statistics and a parameter vector  $\eta \in R^p$ . The canonical exponential family model is

$$(4.1) \quad P_{\eta}(Y = y) = \exp\{\eta \cdot g(y) - \kappa(\eta)\}, \quad y \in \mathcal{Y}$$

where  $\exp\{\kappa(\eta)\} = \sum_{u \in \mathcal{Y}} \exp\{\eta \cdot g(u)\}$  is the familiar normalizing constant associated with an exponential family of distributions [Barndorff-Nielsen (1978); Lehmann (1983)].

The range of network statistics that might be included in the  $g(y)$  vector is vast—see Wasserman and Faust (1994) for the most comprehensive treatment of these statistics—though we will consider only a few in this article. We allow the vector  $g(y)$  to include covariate information about nodes or edges in the graph in addition to information derived directly from the matrix  $y$  itself.

There has been a great deal of work on models of the form (4.1), to which we refer as exponential family random graph models or ERGMs for short. [We avoid the lengthier EFRGM, for “exponential family random graph models,” both for the sake of brevity and because we consider some models in this article that should technically be called *curved* exponential families Hunter and Handcock (2006).]

The normalizing constant is usually difficult to compute directly for  $\mathcal{Y}$  containing large numbers of networks. Inference for this class of models was considered in the seminal paper by Geyer and Thompson (1992), building on the methods of Frank and Strauss (1986) and the above cited papers. Until recently, inference for social network models has relied on maximum pseudolikelihood estimation [Besag (1974); Frank and Strauss (1986); Strauss and Ikeda (1990); Geyer and Thompson (1992)]. Geyer and Thompson (1992) proposed a stochastic algorithm to approximate maximum likelihood estimates for model (4.1), among other models; this Markov chain Monte Carlo (MCMC) approach forms the basis of the method described in this article. The development of these methods for social

network data has been considered by [Corander, Dahmström and Dahmström \(1998\)](#); [Crouch, Wasserman and Trachtenberg \(1998\)](#); [Snijders \(2002\)](#); [Handcock \(2002\)](#); [Corander, Dahmström and Dahmström \(2002\)](#); [Hunter and Handcock \(2006\)](#).

4.1. *Likelihood-based inference for ERGM.* In this section we consider likelihood inference for  $\eta$  in the case where  $Y = Y_{\text{obs}} + Y_{\text{mis}}$  is possibly only partially observed.

As the direct computation of the likelihood is difficult when the number of networks in  $\mathcal{Y}$  is large, we can approximate the likelihood by using the MCMC approach of randomly sampling from the space of possible values of the missing data and taking the mean. Alternatively, consider the conditional distribution of  $Y$  given  $Y_{\text{obs}}$ :

$$P_{\eta}(Y_{\text{mis}} = v | Y_{\text{obs}} = y_{\text{obs}}) = \exp[\eta \cdot g(v + y_{\text{obs}}) - \kappa(\eta | y_{\text{obs}})], \quad v \in \mathcal{Y}(y_{\text{obs}}),$$

where  $\exp[\kappa(\eta | y_{\text{obs}})] = \sum_{u \in \mathcal{Y}(y_{\text{obs}})} \exp[\eta \cdot g(u + y_{\text{obs}})]$ . This formula gives a simple way to sample from the conditional distribution and hence produce multiple imputations of the full data. Specifically, the conditional distribution of  $Y$  given  $Y_{\text{obs}}$  is an ERGM on a constrained space of networks, and hence one can simulate from it using a variant of the standard MCMC for ERGM [[Hunter and Handcock \(2006\)](#); [Handcock et al. \(2003\)](#)] that restricts the proposed networks to the subset of networks that are concordant to the observed data.

Also note that

$$L[\eta | Y_{\text{obs}} = y_{\text{obs}}] \propto \exp[\kappa(\eta | y_{\text{obs}}) - \kappa(\eta)]$$

which can then be estimated by MCMC samples: the first term by a chain on the complete data and the second by a chain conditional on  $y_{\text{obs}}$ . So the sampled data situation is only slightly more difficult than the complete data case.

**5. Two-wave link-tracing samples from a collaboration network.** In this section we investigate the effect of network sampling on estimation by comparing network samples to the situation where we observe the complete network. Specifically, we consider the collaborative working relations between 36 partners in a New England law firm introduced in Section 1. These data have been studied by many authors including [Lazega \(2001\)](#), [Snijders et al. \(2006\)](#) and [Hunter and Handcock \(2006\)](#) (whom we follow).

We consider an ERGM (4.1) with two network statistics for the direct effects of seniority and practice of the form

$$\sum_{1 \leq i, j \leq n} y_{ij} X_i,$$

where  $X_i$  is the seniority or practice of partner  $i$ . We also consider three dyadic homophily attributes based on practice, gender and office. These are included as three network statistics indicating matches between the two partners in the dyad on the given attribute:

$$\sum_{1 \leq i < j \leq n} y_{ij} \mathcal{I}(X_i = X_j),$$

where  $\mathcal{I}(x)$  indicates the truth of the condition  $x$  and  $X_i$  and  $X_j$  are the practice, gender or office attribute of partner  $i$  and  $j$ , respectively. We also include statistics that are purely functions of the relations  $y$ . These are the number of edges (essentially the density) and the geometrically weighted edgewise shared partner statistic (denoted by GWESP), a measure of the transitivity structure in the network [Snijders et al. (2006)]. The model is a slightly reparameterized form of Model 2 in Hunter and Handcock (2006) obtained by replacing the alternating  $k$ -triangle term with the GWESP term. The scale parameter for the GWESP term is fixed at its optimal value (0.7781). See Hunter and Handcock (2006) for details.

As discussed in Hunter and Handcock (2006), this model provides an adequate fit to the data, and we will use it here to assess the effect of sampling on model fit. A summary of the MLE parameters used is given in the *complete data value* column of Table 2. Note that we are taking these parameters as “truth” and considering data produced by sampling from this network.

We construct all possible datasets produced by a two-wave link-tracing design starting from two randomly chosen nodes (the “seeds”). This adaptive design is amenable to the model. As there are 36 partners and the

TABLE 2  
*Bias and Root Mean Squared Error (RMSE) of natural parameter MLE based on two-wave samples as percentages of true parameter values and efficiency losses*

Natural parameter	Complete data value	Bias (%)	RMSE (%)	Efficiency loss (%)
<i>Structural</i>				
Edges	-6.51	0.2	1.2	1.7
GWESP	0.90	0.8	3.7	5.1
<i>Nodal</i>				
Seniority	0.85	0.3	3.1	1.3
Practice	0.41	0.4	5.3	3.5
<i>Homophily</i>				
Practice	0.76	0.8	4.3	2.9
Gender	0.70	0.9	4.7	1.7
Office	1.15	0.7	2.9	2.8



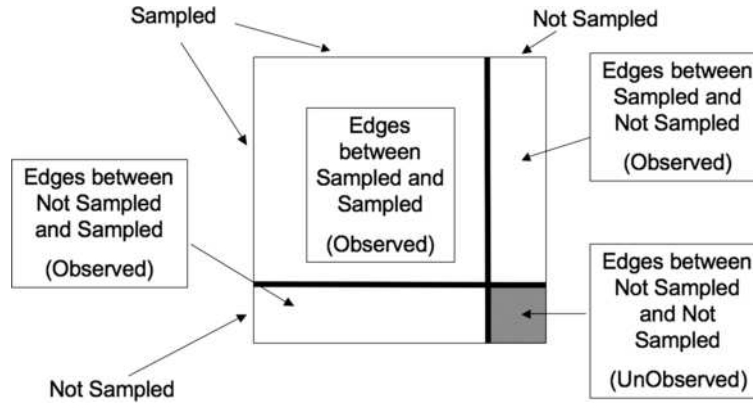


FIG. 1. Schematic depiction of sampled and unobserved arc data when the sampling is over an undirected network.

sample is deterministic given the seeds, there are  $\binom{36}{2} = 630$  possible data sets. The number of actors in each dataset varies from just 2 to all 36 depending on the degree of connectedness of the seeds. The data pattern is shown in Figure 1. Consider a partition of the sampled from the nonsampled and the corresponding  $2 \times 2$  blocking of the sociomatrix, with the four blocks representing dyads from sampled and nonsampled to sampled and nonsampled. The complete data consists of the full sociomatrix. The first three blocks contain the observed data, the dyads involving at least one sampled node, and the last block contains the unobserved data, those between the nonsampled.

For each of these samples we use the methods of Section 4.1 to estimate the parameters. We can then compare them to the MLE for the complete dataset. For these networks, the MLEs are obtained using `statnet` [Handcock et al. (2003)], both for the natural parametrization and for the mean value parameterization [see Handcock (2003)].

The mean value parameters are a function of the natural parameters, specifically the expected values of the sufficient statistics given the values of the natural parameters.

There are two isolates, that is nodes with no relations. If these two are selected as the two seeds, only 69 of the 630 dyads are observed, and no edges are observed. Therefore, the MLE associated with this sample includes (negative) infinite values, on the boundary of the convex hull. For this reason, we exclude this sample from our analyses. Practically, this exclusion is reasonable in that it is unlikely any researcher drawing a link-tracing sample including only two isolated nodes will proceed with analysis of that sample.

One way to assess the effect of the link-tracing design is to compare the estimates from the sampled data to that of the complete data. As a measure

of the difference between the estimates in the metric of the model, we use the Kullback–Leibler divergence from the model implied by the complete data estimate to that of the sampled data estimate. Recall that the Kullback–Leibler divergence of a distribution with probability mass function  $p$  from the distribution with probability mass function  $q$  is

$$E_q[\log(q) - \log(p)].$$

Let  $\eta$  and  $\xi$  be alternative parameters for the model (4.1). The Kullback–Leibler divergence,  $KL(\xi, \eta)$ , of the ERGM with parameter  $\eta$  from the ERGM with parameter  $\xi$  is

$$\begin{aligned} E_\xi \left[ \log \left( \frac{P_\xi(Y = y)}{P_\eta(Y = y)} \right) \right] &= \sum_{y \in \mathcal{Y}} \log \left( \frac{P_\xi(Y = y)}{P_\eta(Y = y)} \right) P_\xi(Y = y) \\ &= \sum_{y \in \mathcal{Y}} (\xi - \eta) \cdot y P_\xi(Y = y) + \kappa(\eta) - \kappa(\xi) \\ &= (\xi - \eta) \cdot E_\xi[g(Y)] + \kappa(\eta) - \kappa(\xi). \end{aligned}$$

If  $\xi$  is the complete data MLE then  $E_\xi[g(Y)] = g(Y_{\text{obs}})$  are the observed statistics (given in the *complete data value* column of Table 3). The divergence can be easily computed using the MCMC algorithms of Section 4.1.

Figure 2 plots the Kullback–Leibler divergence of the MLEs based on the 629 samples from the complete data MLE. The Kullback–Leibler divergence of the two smallest samples, including only 5 nodes (165 dyads), are about 14 and have not been plotted to reduce the vertical scale. The horizontal axis is the number of observed dyads in the sample. The plot indicates

TABLE 3  
*Bias and Root Mean Squared Error (RMSE) of mean value parameter MLE based on two-wave samples as percentages of true parameter values and efficiencies*

Natural parameter	Complete data value	Bias (%)	RMSE (%)	Efficiency loss (%)
<i>Structural</i>				
Edges	115.00	0.4	2.0	1.8
GWESP	190.31	0.4	2.8	1.9
<i>Nodal</i>				
Seniority	130.19	0.3	1.8	1.4
Practice	129.00	0.2	2.6	3.4
<i>Homophily</i>				
Practice	72.00	0.1	2.0	1.7
Gender	99.00	0.5	2.1	1.8
Office	85.00	0.7	2.7	3.0

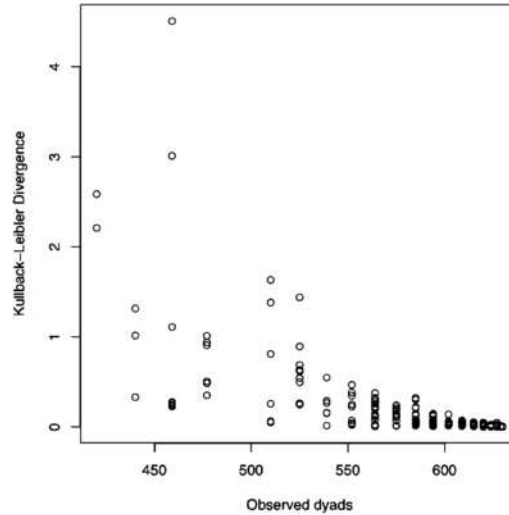


FIG. 2. *Kullback–Leibler divergence of the MLEs based on the samples compared to the complete data MLE. As the number of dyads sampled increases, the information content of the samples approaches that of the complete data. The information loss for the majority of samples is modest.*

how the information in the data about the complete data MLE approaches that of the complete data as the number of sampled dyads approaches the full number. The key feature of this figure is the *variation* in information content among samples of the same size especially for the smaller sample sizes. Different seeds lead to samples that tell us different things about the model even when the numbers of partners surveyed is the same.

For more specific information on the individual estimates, we can compute the bias of the estimates based on the samples as the mean difference between the parameter estimates from the samples and that of the complete network. The root mean squared error (RMSE) is the square-root of the mean of the squared difference between the parameter estimates from each sample and the complete data estimates. The efficiency loss of the sampled estimate is the ratio of the mean squared error and the variance of the sampling distribution of the estimate based on the full data. This standardizes the error in the sampled estimates by the variation in the complete data estimates. We also complete a similar comparison of the estimates under the alternative mean value parametrization [Handcock (2003)].

The properties of the natural parameter estimates are summarized in Table 2. The bias and root mean squared error are presented in percentages of the complete data parameter estimates.

The bias is very small and the RMSE is modest. The efficiency loss is 2%–3% on average. Note that these population-average figures obscure the variation in loss over individual samples apparent in Figure 2.

Table 3 is the mean value parameterization analog of Table 2. As these are on the same measurement scale as the statistics they are easier to interpret. Again we see that the estimates are approximately unbiased and the RMSE and efficiency losses are small.

**6. Discussion.** In this paper we give a concise and systematic statistical framework for dealing with partially observed network data resulting from a designed sample. The framework includes, but is not restricted to, adaptive network sampling designs. We present a definition of a network design which is amenable to a given model and a result on likelihood-based inference under such designs.

An important simple result of this framework is that sampled networks are not “biased” but can be representative if analyzed correctly. Many authors have confused the ideas of simple random sampling of the dyads with representative designs. The results of this paper indicate that simple random sampling is not necessary for valid inference. In fact, the most commonly used designs can be easily taken into account. Hence, despite their form, inference from adaptive network samples is tractable.

It is illustrative to compare our approach to that of [Stumpf, Wiuf and May \(2005\)](#). These authors highlight the difference between the structure of a network and that of a sub-network induced by Bernoulli sampling of its nodes. The framework in this paper allows valid inference for the properties of the network based on its partial observation. This is because we fit a broad class of models compatible with an arbitrary set of network statistics (e.g., ERGM) for the complete network and use a method of inference that does not rely on equality between the structure of the full and sub-networks. As illustrated by the work of [Stumpf, Wiuf and May \(2005\)](#), treating the observed portion as if it were the full network may lead to invalid inference about characteristics of the full network such as the degree distribution.

We have also shown that likelihood-based inference from an adaptive network sample can be conducted using a complete network model. We have shown that such inference is both principled and practical. The likelihood framework naturally accommodates standard sampling designs. Note that in a design-based frame, principled inference would require a great deal of effort to precisely characterize the sampling designs. The result that link-tracing designs are adaptive and can be analyzed with complex likelihood based methods is very valuable in practice as these designs have previously not been analyzed with general exponential family random graph (or similar) models. The only prior work appears to be that of [Thompson and Frank \(2000\)](#) who applied a less complex model class.

In our application we show that an adaptive network sampling of a collaboration network can lead to effective estimates of the model parameters in the vast majority of cases. We find that the MLEs from the samples have only modest bias (compared to the complete data estimate) and an error that only increases slowly with the number of unobserved dyads. We also show that the information content of the sample (with respect to the model), varies greatly even for samples of the same size. For conventional samples of i.i.d. random variables, the Fisher information is simply proportional to the sample size. In the network setting with dependence terms, however, the Fisher information will depend on the specific set of nodes and dyads sampled. For example, the information component corresponding to the GWESP term in the example will be larger for samples in which more pairs of nodes joined by edges are sampled, as GWESP applies only to pairs of nodes joined by edges. If no such dyads were sampled, there would be no information in the sample about the propensity for nodes joined by edges to have relations in common.

In practice the sample is a result of a combination of the sampling design and an *out-of-design* mechanism. The sampling design is the part of the observation process under the control of the surveyor. When adaptive designs are executed faithfully, the unknown dyads are assumed to be intentionally unobserved, or missing by design. Note that the definition of control may be extended to nonamenable sampling designs, for example by allowing the design to depend on unknown factors, such as the unrecorded values of variables used for stratification. The out-of-design mechanism is the nonintentional nonobservation of network information (e.g., due to the failure to report links, incomplete measurement of links and attrition from longitudinal surveys). This is also referred to, in general, as the *non-response mechanism*. We consider the joint effect of sampling and missing data in a companion paper [Handcock and Gile (2007)].

**Acknowledgments.** The authors would like to thank the members of the UW Network Modeling Group (Martina Morris, P. I.), Stephen Fienberg and the reviewers for their helpful input.

## SUPPLEMENTARY MATERIAL

### Supplement: Software used in the simulation study

(DOI: [10.1214/08-AOAS221SUPP](https://doi.org/10.1214/08-AOAS221SUPP); .zip). The code used to perform this study is written in the R statistical language [R Development Core Team (2007)] and is based on `statnet`, an open-source software suite for network modeling [Handcock et al. (2003)]. We provide the code and documentation for it with links to the `statnet` website.

## REFERENCES

- BARNDORFF-NIELSEN, O. E. (1978). *Information and Exponential Families in Statistical Theory*. Wiley, New York. [MR0489333](#)
- BESAG, J. (1974). Spatial interaction and the statistical analysis of lattice systems (with discussion). *J. Roy. Statist. Soc. Ser. B* **36** 192–236. [MR0373208](#)
- CORANDER, J., DAHMSTRÖM, K. and DAHMSTRÖM, P. (1998). Maximum likelihood estimation for Markov graphs. Research report, Dept. Statistics, Univ. Stockholm.
- CORANDER, J., DAHMSTRÖM, K. and DAHMSTRÖM, P. (2002). Maximum likelihood estimation for exponential random graph models. In *Contributions to Social Network Analysis, Information Theory, and Other Topics in Statistics; A Festschrift in Honour of Ove Frank* (J. Hagberg, ed.) 1–17. Dept. Statistics, Univ. Stockholm.
- CROUCH, B., WASSERMAN, S. and TRACHTENBERG, F. (1998). Markov chain Monte Carlo maximum likelihood estimation for  $p^*$  social network models. In *The XVIII International Sunbelt Social Network Conference, Sitges, Spain*.
- FRANK, O. (2005). Network Sampling and Model Fitting. In *Models and Methods in Social Network Analysis* (J. S. P. Carrington and S. S. Wasserman, eds.) 31–56. Cambridge Univ. Press, Cambridge.
- FRANK, O. and STRAUSS, D. (1986). Markov Graphs. *J. Amer. Statist. Assoc.* **81** 832–842. [MR0860518](#)
- GEYER, C. J. and THOMPSON, E. A. (1992). Constrained Monte Carlo maximum likelihood calculations (with discussion). *J. Roy. Statist. Soc. Ser. B* **54** 657–699. [MR1185217](#)
- HANDCOCK, M. S. (2002). Degeneracy and inference for social network models. In *The Sunbelt XXII International Social Network Conference, New Orleans, LA*.
- HANDCOCK, M. S. (2003). Assessing degeneracy in statistical models of social networks. Working paper 39, Center for Statistics and the Social Sciences, Univ. Washington. Available at <http://www.csss.washington.edu/Papers>.
- HANDCOCK, M. S. and GILE, K. J. (2007). Modeling social networks with sampled or missing data. Working paper 75, Center for Statistics and the Social Sciences, Univ. Washington. Available at <http://www.csss.washington.edu/Papers>.
- HANDCOCK, M. S. and GILE, K. J. (2010). Supplement to “Modeling social networks from sampled data.” DOI: [10.1214/08-AOAS221SUPP](https://doi.org/10.1214/08-AOAS221SUPP).
- HANDCOCK, M. S., HUNTER, D. R., BUTTS, C. T., GOODREAU, S. M. and MORRIS, M. (2003). statnet: Software tools for the statistical modeling of network data statnet project <http://statnet.org/>, Seattle, WA. R package version 2.0. Available at <http://CRAN.R-project.org/package=statnet>.
- HUNTER, D. R. and HANDCOCK, M. S. (2006). Inference in curved exponential family models for networks. *J. Comput. Graph. Statist.* **15** 565–583. [MR2291264](#)
- LAZEGA, E. (2001). *The Collegial Phenomenon: The Social Mechanisms of Cooperation Among Peers in a Corporate Law Partnership*. Oxford Univ. Press, Oxford.
- LEHMANN, E. L. (1983). *Theory of Point Estimation*. Wiley, New York, NY. [MR0702834](#)
- R DEVELOPMENT CORE TEAM (2007). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, Version 2.6.1. Available at <http://www.R-project.org/>.
- ROSENBAUM, P. R. and RUBIN, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* **70** 41–55. [MR0742974](#)
- RUBIN, D. B. (1976). Inference and missing data. *Biometrika* **63** 581–592. [MR0455196](#)
- SÄRNDAL, C.-E., SWENSSON, B. and WRETMAN, J. (1992). *Model Assisted Survey Sampling*. Springer, New York. [MR1140409](#)
- SNIJDERS, T. A. B. (1992). Estimation on the basis of snowball samples: How to weight. *Bulletin Methodologie Sociologique* **36** 59–70.

- SNIJDERS, T. A. B. (2002). Markov chain Monte Carlo estimation of exponential random graph models. *Journal of Social Structure* **3** 1–41.
- SNIJDERS, T. A. B., PATTISON, P., ROBINS, G. L. and HANDCOCK, M. S. (2006). New specifications for exponential random graph models. *Sociological Methodology* **36** 99–153.
- STRAUSS, D. and IKEDA, M. (1990). Pseudolikelihood estimation for social networks. *J. Amer. Statist. Assoc.* **85** 204–212. [MR1137368](#)
- STUMPF, M. P. H., WIUF, C. and MAY, R. M. (2005). Subnets of scale-free networks are not scale-free: Sampling properties of networks. *Proc. Natl. Acad. Sci. USA* **102** 4221–4224.
- THOMPSON, S. K. and COLLINS, L. M. (2002). Adaptive sampling in research on risk-related behaviors. *Drug and Alcohol Dependence* **68** S57–S67.
- THOMPSON, S. K. and FRANK, O. (2000). Model-based estimation with link-tracing sampling designs. *Survey Methodology* **26** 87–98.
- THOMPSON, S. K. and SEBER, G. A. F. (1996). *Adaptive Sampling*. Wiley, New York. [MR1390995](#)
- WASSERMAN, S. and FAUST, K. (1994). *Social Network Analysis: Methods and Applications*. Cambridge Univ. Press.

DEPARTMENT OF STATISTICS  
UNIVERSITY OF CALIFORNIA  
LOS ANGELES, CALIFORNIA 90095-1554  
USA  
E-MAIL: [handcock@stat.washington.edu](mailto:handcock@stat.washington.edu)

NUFFIELD COLLEGE  
UNIVERSITY OF OXFORD  
NEW ROAD  
OXFORD OX1 1NF  
UNITED KINGDOM  
E-MAIL: [krista.gile@nuffield.ox.ac.uk](mailto:krista.gile@nuffield.ox.ac.uk)