

# Econometric Models of Network Formation

Áureo de Paula<sup>1,2,3</sup>

<sup>1</sup>Department of Economics, University College London, London WC1E 6BT, United Kingdom; email: a.paula@ucl.ac.uk

<sup>2</sup>Centre for Microdata Methods and Practice, London WC1E 7AE, United Kingdom

<sup>3</sup>Institute for Fiscal Studies, London WC1E 7AE, United Kingdom

ANNUAL  
REVIEWS **CONNECT**

[www.annualreviews.org](http://www.annualreviews.org)

- Download figures
- Navigate cited references
- Keyword search
- Explore related articles
- Share via email or social media

Annu. Rev. Econ. 2020. 12:775–99

First published as a Review in Advance on  
May 12, 2020

The *Annual Review of Economics* is online at  
[economics.annualreviews.org](http://economics.annualreviews.org)

<https://doi.org/10.1146/annurev-economics-093019-113859>

Copyright © 2020 by Annual Reviews.  
All rights reserved

JEL codes: C31, C51, D85

## Keywords

network econometrics, dyadic models, strategic network formation models

## Abstract

This article provides a selective review of the recent literature on econometric models of network formation. I start with a brief exposition on basic concepts and tools for the statistical description of networks; then I offer a review of dyadic models, focusing on statistical models on pairs of nodes, and I describe several developments of interest to the econometrics literature. I also present a discussion of nondyadic models in which link formation might be influenced by the presence or absence of additional links, which themselves are subject to similar influences. This argument is related to the statistical literature on conditionally specified models and the econometrics of game theoretical models. I close with a (nonexhaustive) discussion of potential areas for further development.

## 1. INTRODUCTION

Tamás Erdélyi and John Cummings met in their first year of high school in Queens, New York City. Sometime later, John met Doug Colvin, who had moved to the neighborhood from Germany, where his father was stationed with the US Army. Together with a fourth neighborhood friend, Jeffrey Hyman, they formed a band in 1974 and took the name “The Ramones.”

Regardless of one’s opinion on the music by Johnny (John), Joey (Jeffrey), Dee Dee (Doug), and Tommy (Tamás) Ramone, their encounter was a matter of opportunity, talent, and personal affinity.<sup>1</sup> Similar factors inform the formation of networks in many other settings, from friendships in economic and/or social contexts to the establishment of business relationships across firms and of trading ties between countries. These links influence and are facilitated by other outcomes that depend on how a particular group is connected. For example, smoking behavior among adolescents influences and is affected by whom the adolescents are related to. Production decisions are similarly influenced by the set of clients and suppliers for a firm. Especially with the advent and availability of data on networks from surveys and administrative sources, it is increasingly feasible and potentially relevant to provide a systematic assessment of drivers and correlates for the formation of such relationships.

In this article I review recent developments in the nascent literature on the econometrics of networks that provide a means to assess those drivers and correlates. Since an established literature exists in other fields, spanning social sciences, physics, and statistics, I focus here on tools related to the field of economics. Much of this literature builds and gets inspiration from previous developments in other fields, but some of it is quintessentially related to research in economics. This is the case, for instance, for estimable models of strategic network formation.

I start with a brief exposition on basic concepts and tools for the statistical description of networks. I then offer a review of dyadic models and describe several developments of interest to the econometrics literature and connections with related literatures (e.g., panel data). This is followed by a discussion of nondyadic models in which link formation might be influenced by the presence or absence of additional links, which themselves are subject to similar influences. This argument is related to the statistical literature on conditionally specified models and the econometrics of game theoretical models. As in those literatures, potential simultaneity issues are central to the identification and estimation of the quantities of interest, and I explore some of those here as well. I close with a (nonexhaustive) discussion of potential areas for further development.

This review complements other available overviews within econometrics by Chandrasekhar (2015), Graham (2015, 2020), and de Paula (2017). In particular, it follows closely (but expands upon) a previous survey by de Paula (2017) and corresponding material by Graham & de Paula (2020).

## 2. PRELIMINARY FOUNDATIONS

Networks are typically represented by graphs. A graph  $g$  is a pair of sets  $(\mathcal{N}_g, \mathcal{E}_g)$  of nodes (or vertices)  $\mathcal{N}_g$  and edges (or links or ties)  $\mathcal{E}_g$ . I will denote the number of elements in these sets by  $|\mathcal{N}_g| \equiv N$  and  $|\mathcal{E}_g|$ , respectively. Vertices here represent economic agents like individuals, households, or firms. An edge represents a link or connection between two nodes in  $\mathcal{N}_g$ . A graph is undirected when  $\mathcal{E}_g$  is the set of unordered pairs with elements in  $\mathcal{N}_g$ , e.g.,  $\{i, j\}$  with  $i, j \in \mathcal{N}_g$ . (I abstract away from self-links here.) An example is (reciprocal) informal risk-sharing networks

<sup>1</sup>The off-stage personal relationships were not harmonious for long. Nevertheless, Dee Dee Ramone wrote, “I can see now how it was only natural that I would gravitate toward Tommy, Joey, and Johnny Ramone. . . . They were the obvious creeps of the neighbourhood” (cited in Gilmore 2016).

based on kinship or friendship (e.g., Fafchamps & Lund 2003). To accommodate directional relationships, edges are best modeled as ordered pairs, e.g.,  $(i, j) \in \mathcal{N}_g \times \mathcal{N}_g$ . These graphs, known as directed graphs (or digraphs), are more adequate for handling relationships that do not require reciprocity or for which direction carries a particular meaning, as in a supplier-client relationship in a production network (e.g., Atalay et al. 2011). Generalizations allow for weighted links, perhaps representing distances between two individuals or the intensity of a particular relationship. Such weights can be represented as a mapping from the space of pairs (unordered or ordered) into the real line.

A common representation of a graph is through its  $|\mathcal{N}_g| \times |\mathcal{N}_g|$  adjacency matrix  $W$ , where each line represents a different node. The components of  $W$  mark whether an edge between nodes  $i$  and  $j$  (or from  $i$  to  $j$  in a digraph) is present or not, and possibly its weight (in weighted graphs). The adjacency matrix allows one to translate combinatorial operations into linear algebraic ones and can be quite useful in several settings. For an adjacency matrix  $W$  to a simple unweighted graph (i.e., no self-links and at most one link between any pair of nodes), the  $ij$  element of matrix  $W^k$ ,  $k \in \{1, \dots, N-1\}$ , for instance, produces the number of paths of length  $k$  between  $i$  and  $j$ . Two graphs are said to be isomorphic if their adjacency matrices can be obtained from each other, for example through multiplication by a permutation matrix. This corresponds to a relabeling of the vertices in the corresponding graphs.

## 2.1. Vertex Features

Various measures can then be defined to characterize a particular vertex in the graph, to relate two or more vertices on a graph, or to represent a global feature of the graph at hand. In what follows I focus on simple, unweighted graphs for ease of exposition. An important characteristic for a particular vertex  $i$ , for example, is the set of neighbors incident with that vertex in a graph  $g$ , denoted by  $N_i(g)$ . The cardinality of this set is known as the degree of that node, and one can then talk about the relative frequency of degrees in a given graph as a whole. (In directed graphs, one can further distinguish between in-degrees and out-degrees, relating respectively to inward and outward edges from and to a given node.) A dense graph, for instance, is one in which nodes display a lot of connections, and a common measure of density is the average degree divided by  $|\mathcal{N}_g| - 1$ , which is the maximum number of possible links available to any given node. It is common to define the (geodesic) distance between two nodes as the shortest path between them. A graph is then said to be connected if the distance between any two vertices is finite (i.e., there is at least one path between any two nodes).

One can also define various measures to characterize the typical subnetwork structure around a given vertex. For brevity, I only mention a basic taxonomy of such measures, as specific definitions are available in most introductory texts on the subject (see, for example, the excellent overview in Jackson 2009). A network aspect of particular interest in social settings is the degree of clustering in the system, intuitively summarized by the propensity of two neighbors to a given node to be themselves directly linked. Theoretically, for example, it may be easier for clustered individuals to coordinate on certain collective actions, since clustering may facilitate common knowledge (Chwe 2000), and different clustering metrics are available to quantify this feature in a network.

Another feature of potential interest in economic and social networks is the degree of centrality of a given vertex, and various measures of centrality are also available. Those aim at characterizing how important a given node is in comparison to the remaining nodes in  $g$ . Aside from how connected a given vertex is (degree centrality) or how far on average a vertex is from any other vertex in the network (closeness centrality), one can also compute the betweenness centrality, illustrating how crucial a given node is in connecting individuals. Another family of popular centrality

measures summarize a node's centrality in reference to its neighbors' centrality (more on this later). The simplest of these measures is the eigenvector centrality (a.k.a. Gould's index of accessibility), corresponding to the dominant eigenvector of the adjacency matrix (Gould 1967, Bonacich 1972). If an individual is central as a consequence of being connected with central individuals, this means that the (eigenvector) centrality vector  $\mathbf{c}$  is such that  $\mathbf{c} = W\mathbf{c}$ . It is thus the eigenvector related to a unit eigenvalue. For a normalized adjacency matrix where all rows add up to one, this is guaranteed to exist and be unique if the network is (strongly) connected by the Perron-Frobenius theorem. Variations on this centrality measure include, for example, Google's PageRank index (Brin & Page 1998). Among the most popular metrics in this family were those proposed by Katz (1953) and Bonacich (1987). The Katz centrality of a node  $i$  can be motivated by ascribing a value of  $\tilde{\beta}^k > 0$  to each connection reached by a walk of length  $k$ . If one adds up the weights for each individual, one has a centrality measure for each individual given by the components of the vector  $\tilde{\beta}W\mathbf{1} + \tilde{\beta}^2W^2\mathbf{1} + \tilde{\beta}^3W^3\mathbf{1} + \dots$ . If  $\tilde{\beta}$  is below the reciprocal of  $W$ 's largest eigenvalue, we can write the above as  $\tilde{\beta}(\mathbf{I} - \tilde{\beta}W)^{-1}W\mathbf{1}$ , where  $\tilde{\beta}$  is a small-enough positive number. The Bonacich centrality generalizes this formula to a two-parameter index defined by the vector  $\alpha(\mathbf{I} - \tilde{\beta}W)^{-1}W\mathbf{1}$ . Such measures turn out to play an important role in the analysis of games and dissemination on networks (see, e.g., Ballester et al. 2006 and the survey in Zenou 2016).

## 2.2. Random Graphs

Letting  $\mathcal{G}$  be a particular set of graphs, one can define a probability model on this sample space. These models can be (and usually are) indexed by features common to the graphs in  $\mathcal{G}$ , like the number of vertices and/or other motifs.

Regardless of the generative model one has in mind, the initial examination of a network typically involves the characterization of some of the features alluded to above, such as the link pattern among node pairs, triads, tetrads, and  $k$ -tuples in general, usually referred to in the first two cases as a dyad or triad census. Isomorphic graphs on the  $k$ -tuple of nodes are usually classified within the same equivalence class. For example, in a dyad from a digraph, a pair with one link will be classified equivalently whether the link is from  $i$  to  $j$  or from  $j$  to  $i$ . Thus, if the network is directed, there are three patterns of subgraphs on pairs: no link, one unreciprocated link from one of the nodes to the other one, or two reciprocated links. In an undirected network, two patterns exist for a pair of nodes: Either those two vertices are connected or not. In triads, an undirected network will feature four patterns: zero, one, two, or three links. A directed network will produce a richer set, with 16 triad patterns. The number of patterns grows rapidly with  $k$  and, in a directed network, "there are so many tetrad types (218) and pentad types (9608) that a  $k$ -subgraph census for  $k \geq 4$  is often more cumbersome than the original sociomatrix" (Holland & Leinhardt 1976, p. 7).<sup>2</sup>

One of the early network models imposes a uniform probability on the class of graphs with a given number of nodes,  $|\mathcal{N}_g|$ , and a particular number of edges,  $|\mathcal{E}_g|$ , for  $g \in \mathcal{G}$  (see Erdős & Rényi 1959, 1960). Another basic, canonical random graph model is one in which the edges between any two nodes follow an independent Bernoulli distribution with equal probability  $p$ . For a large-enough number of nodes and a sufficiently small probability of link formation  $p$ , the degree distribution approaches a Poisson distribution, and the model is consequently known as the Poisson random graph model. This class of models appears in work by Gilbert (1959) and Erdős & Rényi (1960) and has since been studied extensively. While simple to characterize, it fails to

<sup>2</sup>Holland & Leinhardt (1976) provide expressions for the first two moments of such  $k$ -subgraph censuses, and Bickel et al. (2011) and Bhattacharya & Bickel (2015) present further asymptotic results under additional assumptions on the link formation probabilities.

reproduce important dependencies observed in social and economic networks. One category of models that aims at a better representation of the regularities usually encountered in social systems involves models in which nodes are incorporated into the graph sequentially and form ties more or less randomly. One such model is the preferential attachment model (Barabási & Albert 1999), in which the establishment of new links is more likely for higher-degree existing nodes, producing degree distributions with Pareto tails (as well as other regularities usually observed) (see the presentation in Jackson 2009 or Kolaczyk 2009 for a more thorough exposition).

Another alternative is to rely on more general (static) random graph models that explicitly acknowledge the probabilistic dependencies in link formation. Frank & Strauss (1986), for example, focus on random graphs in which two (random) edges that do not share a vertex are conditionally independent given the other remaining (random) edges. This model reflects the intuition that ties are not independent of each other, but their dependency arises only through those ties that are directly involved in the connections in question. This, and a homogeneity assumption (i.e., the assumption that all graphs that are the same up to a permutation of vertices have the same probability), delivers that

$$\mathbb{P}(G = g) \propto \exp \left( \alpha_0 t + \sum_{k=1} \alpha_k s_k \right),$$

where  $t$  is the number of triangles (completely connected triples of vertices) and  $s_k$  is the number of  $k$ -stars (tuples of  $k + 1$  vertices where one of the vertices has degree  $k$  and the remaining ones have degree one). (Notice that the Erdős-Rényi model is a specific case of the above model, where  $\alpha_0 = \alpha_2 = \dots = \alpha_k = 0$ .) This structure suggests a class of probabilistic models that reproduce the exponential functional form above even in cases in which the exact properties used by Frank & Strauss (1986) do not hold. Those models are such that  $\mathbb{P}(G = g) \propto \exp(\sum_{k=1}^p \alpha_k S_k(g))$ , where  $S_k(g)$ ,  $k = 1, \dots, p$  enumerate features of the graph  $g$ . These would be characteristics like the number of edges, the number of triangles, and possibly many others. These models are known as exponential random graph models (ERGMs) (or  $p^*$  models in the social sciences literature) (see Robins et al. 2007) and can be extended beyond undirected random graphs. The models above form an exponential family of distribution over (random) graphs, and exponential distributions (e.g., Bernoulli, Poisson) have well-known probabilistic and statistical properties. For example, the vector  $(S_1(g), \dots, S_p(g))$  constitutes a  $p$ -dimensional sufficient statistic for the parameters  $(\alpha_1, \dots, \alpha_p)$ .

Some recent articles in econometrics are closely related to the exponential random graph model above. More generally, all the models above (and many others) are presented in detail elsewhere (e.g., Bollobás 2001, Jackson 2009, Kolaczyk 2009), and I will selectively discuss features and difficulties as they arise in the literature reviewed here.

### 3. DYADIC MODELS

As their name suggests, dyadic models offer a statistical framework centered on node pairs. In the Erdős-Rényi model with  $|\mathcal{N}_g| = N$  vertices, the class of probability distributions over possible networks,  $\mathcal{P}$ , is indexed by the probability  $p \in (0, 1)$  that a link is formed (independently) between any two nodes. Zheng et al. (2006), for example, use a heterogeneous version of this simple random graph model to obtain estimates for the total size of hard-to-count populations. Observed heterogeneity has also been incorporated via dyadic models that expand on this model, just as a probit or logit model generalizes a simple Bernoulli statistical model,

$$W_{ij} = \mathbf{1}(X_{ij}^\top \beta + \epsilon_{ij} \geq 0),$$

which can be used on directed or undirected settings. The covariates,  $X_{ij}$ , involve individual-specific variables as well as variables defined for the pair (like geographic distance) or individual variables aggregated into pair-specific quantities.<sup>3</sup> For example, the metric  $f(X_i, X_j) = |X_i - X_j|$  between sender and receiver is usually of interest and gathers how relevant similarity (or homophily) in those variables is for link formation. (The coefficient on  $X_{ij}$  is called the homophily parameter by some authors in the statistics literature.) Other pairwise functions possibly used include the inner product or projection functions, though others can also be devised. Whereas estimates can be obtained by usual methods (e.g., likelihood, method of moments), inference in this context typically pays special attention to potential correlations across pairs. If there are individual factors that are common to all observations related to a given node, dyads are not independent and  $\mathbb{E}(\epsilon_{ij}\epsilon_{kl})$  is potentially nonzero when  $i = k, j = l, i = l, \text{ or } j = k$ . In a linear regression model, Fafchamps & Gubert (2007), for example, suggest the following robust variance estimator,

$$\frac{1}{N-K}(\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{i,j,k,l} \frac{m_{ijkl}}{2N} X_{ij} \hat{\epsilon}_{ij} \hat{\epsilon}_{kl} X_{kl}^\top \right) (\mathbf{X}^\top \mathbf{X})^{-1},$$

where the matrix  $\mathbf{X}$  stacks the dyadic covariates and  $m_{ijkl}$  is one when  $i = k, j = l, i = l, \text{ or } j = k$  and zero otherwise. They also discuss standard errors for the logit (see Fafchamps & Gubert 2007, p. 330 n. 7), and additional examination is provided by Cameron & Miller (2014).<sup>4</sup>

It is possible to extend the Erdős-Rényi model to incorporate other features. Consider, for example, the well-known  $p_1$  dyadic model offered by Holland & Leinhardt (1981) for directed networks. Their model, which also belongs to the exponential family, postulates that

$$\mathbb{P}(W_{ij} = W_{ji} = 1) \propto \exp(\alpha^{\text{rec}} + 2\alpha + \alpha_i^{\text{out}} + \alpha_i^{\text{in}} + \alpha_j^{\text{out}} + \alpha_j^{\text{in}}) \quad 1.$$

and

$$\mathbb{P}(W_{ij} = 1, W_{ji} = 0) \propto \exp(\alpha + \alpha_i^{\text{out}} + \alpha_j^{\text{in}}). \quad 2.$$

Here, the parameter  $\alpha_i^{\text{out}}$  encodes the tendency of node  $i$  to send out links irrespective of the target (i.e., its gregariousness), and  $\alpha_j^{\text{in}}$  captures node  $j$ 's tendency to receive links regardless of the sender's identity (i.e., its attractiveness). Since they regulate individual link probabilities, these parameters drive heterogeneity in individual (in and out) degrees across agents in the group. They also influence how dense or sparse the network tends to be: The more negative those parameters are, the less likely a link is. The parameter  $\alpha^{\text{rec}}$  registers the tendency for directed links to be reciprocated: Large, positive values of  $\alpha^{\text{rec}}$  will increase the likelihood of symmetric adjacency matrices. [Note that when  $\alpha_i^{\text{out}} = \alpha_i^{\text{in}} = \alpha^{\text{rec}} = 0$  for all  $i$ , links form independently with probability given by  $\exp(\alpha)/(1 + \exp(\alpha))$ , and the model would correspond to a logit.] Given an observed network, the authors suggest estimating the model above by maximum likelihood and offer an iterative scaling algorithm to ease computation, given the possibly large number of parameters to be estimated.

This model has been generalized and expanded upon. Hoff (2005), for example, considers an augmented model where multiplicative interactions between individual unobserved factors are added to the probability specification above (i.e.,  $\mathbf{z}_i \times \mathbf{z}_j$ , where  $\mathbf{z}_i$  is a vector of  $i$ -specific factors),

<sup>3</sup>In a recent article, Comola & Fafchamps (2017) extend the model to accommodate misreporting/discordance of elicited links by the individuals involved, which is not uncommon in survey-based network data.

<sup>4</sup>Related works, focusing on linear models, are those by Aronow et al. (2015) and Tabord-Meehan (2019).

and those plus the additive gregariousness and attractiveness features defined previously (i.e.,  $\alpha_i^{\text{out}}$  and  $\alpha_i^{\text{in}}$ ) are modeled as random effects. It is also possible to add covariates, as done by Hoff (2005), who estimates his model using Bayesian methods.

Dzemeski (2019) focuses on such a model, where a link from node  $i$  to node  $j$  obeys

$$W_{ij} = \mathbf{1}(X_{ij}^\top \beta + \alpha_i^{\text{out}} + \alpha_j^{\text{in}} + \epsilon_{ij} \geq 0),$$

where  $\epsilon_{ij}$  is a standard normal random variable. As proposed by Holland & Leinhardt (1981), the parameters  $\alpha_i^{\text{out}}$  and  $\alpha_j^{\text{in}}$  appear as fixed effects, hence allowing for an arbitrary correlation between those and with any observed characteristic in the model. Reciprocity in link formation is captured by a nonzero correlation  $\rho$  between  $\epsilon_{ij}$  and  $\epsilon_{ji}$ , which are jointly normal.<sup>5</sup> Otherwise, errors are independent across pairs. As pointed out by Graham (2017) in an undirected network framework, the presence of individual effects encoded in  $\alpha_i^{\text{out}}$  and  $\alpha_j^{\text{in}}$  may obscure the identification and estimation of the homophily parameter  $\beta$ , as a gregarious individual (i.e., high  $\alpha_i^{\text{out}}$ ) will tend to send links indiscriminately, even if there is a homophilous tendency—that is, high, negative  $\beta$  on the distance between the individual and another node on a particular feature  $X$ . It is important thus to control for those, should degree heterogeneity reflect such sources of heterogeneity.

Because each of the  $N$  individuals has at their disposal  $N - 1$  potential liaisons, the setting is akin to a panel where both dimensions ( $N$  and  $T$ ) are roughly comparable. Dzemeski (2019) thus uses tools from the (large- $N$ , large- $T$ ) panel data literature (see Fernández-Val & Weidner 2016, 2018) to establish large sample properties for the estimator. When the network is sufficiently dense (i.e., the probability of link formation between any two pairs does not vanish asymptotically), the maximum likelihood estimator for  $\beta$  is such that

$$\hat{\beta} \stackrel{a}{\sim} \mathcal{N}(\beta_0 + B_N, V_N),$$

where the asymptotic bias  $B_N$  disappears for large  $N$  and the variance decreases with  $N$ .<sup>6</sup> [Since there are  $N(N - 1)$  ordered pairs, the convergence rate is the usual parametric rate.] Because the number of individual parameters is proportional to  $N$ , an incidental parameter problem yields an asymptotic bias (see Fernández-Val & Weidner 2018). While the estimator is consistent, the asymptotic bias matters for inference. Dzemeski (2019) provides a test of the model based on the prevalence of transitive triads (i.e., vertex triples where links are transitive), which are not used in estimating it, and an application to the microfinance-related networks collected and analyzed by Banerjee et al. (2014).<sup>7</sup>

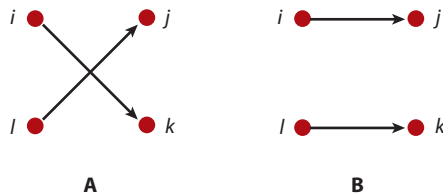
Yan et al. (2019) provide a related analysis for a similar framework, where the errors  $\epsilon_{ij}$  are independent (and  $\rho$  is thus zero) logistic random variables. Under these conditions, the likelihood function for the model is given by

$$\exp \left( \sum_{i,j} w_{ij} X_{ij}^\top \beta + \alpha^{\text{in}\top} \mathbf{d}^{\text{in}} + \alpha^{\text{out}\top} \mathbf{d}^{\text{out}} - \mathcal{C}(\beta, \alpha^{\text{in}}, \alpha^{\text{out}}) \right), \quad 3.$$

<sup>5</sup>One alternative is to follow Holland & Leinhardt (1981) and encode reciprocity in a parameter  $\alpha^{\text{rec}}$ . This would lead to a model akin to  $W_{ij} = \mathbf{1}(X_{ij}^\top \beta + \alpha_i^{\text{out}} + \alpha_j^{\text{in}} + \alpha^{\text{rec}} W_{ji} + \epsilon_{ij} \geq 0)$ . This model is nonetheless compatible with multiple joint distributions for  $(W_{ij}, W_{ji})$  and is thus incomplete. In this case, the model is akin to a game between  $i$  and  $j$  and point-identification is not guaranteed (see de Paula 2013). This is not the case in the original model proposed by Holland & Leinhardt (1981) (see Equations 1 and 2), as the probability distribution over  $(W_{ij}, W_{ji})$  is well defined there.

<sup>6</sup>The estimate for  $\rho$  is obtained in a second step once the maximum likelihood estimate for  $\beta$  is available.

<sup>7</sup>Interestingly, the estimated distributions of gregariousness and attractiveness appear to cluster in a few groups, suggesting group-level heterogeneity.



**Figure 1**

The probability that  $l$  sends a link to  $j$ , as in configuration A, given either configuration A or configuration B, does not depend on individual-specific parameters ( $\alpha$ s).

where  $\mathcal{C}(\beta, \alpha^{\text{in}}, \alpha^{\text{out}}) \equiv \sum_{i \neq j} \ln(1 + \exp(X_{ij}^\top \beta + \alpha_i^{\text{out}} + \alpha_j^{\text{in}}))$  is a normalizing constant, and  $\alpha^{\text{in}}, \alpha^{\text{out}}, \mathbf{d}^{\text{in}}$ , and  $\mathbf{d}^{\text{out}}$  are  $N$ -dimensional vectors stacking the  $\alpha_i^{\text{in}}, \alpha_i^{\text{out}}$  parameters and in-degree and out-degree sequences, respectively. They also rely on similar (though not identical) large- $N$ , large- $T$  panel data manipulations to handle the incidental parameters bias in estimating  $\beta$ , but they pay special attention to the estimation of the individual-specific parameters  $\alpha$ . Their setup relaxes some of the (non-)sparsity constraints posed by Dzemski (2019), who places a lower bound on the likelihood of link formation (see above), allowing the  $\alpha$  parameters to drift in value at a  $\ln N$  rate. As in Dzemski's paper, they also demonstrate the performance of the estimator in simulations and with empirical applications.

In sparser settings, where the errors  $\epsilon_{ij}$  are independent logistic random variables, Charbonneau (2017)—and, for undirected networks, Graham (2017)<sup>8</sup> (see discussion below)—adapts insights from the nonlinear panel data literature on short (i.e., fixed- $T$ ) horizons going back to work by Rasch (1960), to “difference out” the individual parameters in this logistic regression (see, e.g., Arellano & Honoré 2001, section 5.1). The probability mass function for links is the exponential family model in Equation 3. Given the properties of the exponential family, the in-degree and out-degree vectors  $\mathbf{d}^{\text{in}}$  and  $\mathbf{d}^{\text{out}}$  are sufficient statistics for the parameter vectors  $\alpha^{\text{in}}$  and  $\alpha^{\text{out}}$ . This means that the conditional probability for the edges  $W_{ij}$  given those random vectors does not depend on the incidental parameters  $\alpha^{\text{in}}$  and  $\alpha^{\text{out}}$ , thus providing a (conditional) likelihood function. Unfortunately, such conditional likelihood is computationally complex. It is nonetheless still possible to condition the incidental parameters out. Consider, for example, four generic vertices  $i, j, l$ , and  $k$ . The idea is to focus on the probability for one of the two configurations, A or B in **Figure 1**, conditional on either one of the two patterns occurring.

It is not hard to see that the conditional probability that  $l$  sends a link to  $j$  (i.e.,  $\{W_{lj} = 1\}$ ) as in configuration A, given that only one of those links is in place (i.e.,  $\{W_{lj} + W_{lk} = 1\}$ ), equals

$$\mathbb{P}(W_{lj} = 1 | \mathbf{X}, \alpha, W_{lj} + W_{lk} = 1) = \frac{\exp[(X_{lj} - X_{lk})^\top \beta + \alpha_j^{\text{in}} - \alpha_k^{\text{in}}]}{1 + \exp[(X_{lj} - X_{lk})^\top \beta + \alpha_j^{\text{in}} - \alpha_k^{\text{in}}]},$$

where  $\mathbf{X}$  stacks the covariates to all the pairs involved. Note that the expression above does not depend on  $\alpha_k^{\text{out}}$ . Likewise, the likelihood that  $i$  sends a link to  $j$  (i.e.,  $\{W_{ij} = 1\}$ ), as in configuration B, given that only one link is sent to either one of these potential counterparts (i.e.,  $\{W_{ij} + W_{ik} = 1\}$ ), is

$$\mathbb{P}(W_{ij} = 1 | \mathbf{X}, \alpha, W_{ij} + W_{ik} = 1) = \frac{\exp[(X_{ij} - X_{ik})^\top \beta + \alpha_j^{\text{in}} - \alpha_k^{\text{in}}]}{1 + \exp[(X_{ij} - X_{ik})^\top \beta + \alpha_j^{\text{in}} - \alpha_k^{\text{in}}]},$$

<sup>8</sup>This paper was previously circulated with the title “An empirical model of network formation: detecting homophily when agents are heterogeneous.”



which again does not depend on  $\alpha_i^{\text{out}}$ . These equations look like a conventional panel logit model with regressors  $X_j - X_k$  and  $\alpha_j^{\text{in}} - \alpha_k^{\text{in}}$  as a fixed effect. A similar manipulation can then be carried out focusing on the event that  $i$  sends a link to either  $k$  or  $j$  (but not both) (i.e.,  $\{W_{ij} + W_{ik} = 1\}$ ),  $l$  sends a link to either  $k$  or  $j$  (but not both) (i.e.,  $\{W_{lj} + W_{lk} = 1\}$ ), and  $i$  and  $l$  do not send links to the same counterpart (i.e.,  $\{W_{lj} + W_{ij} = 1\}$ ). Given this event, the probability that  $l$  sends a link to  $j$  can be seen to be

$$\begin{aligned} \mathbb{P}(W_{lj} = 1 | \mathbf{X}, \alpha^{\text{in}}, \alpha^{\text{out}}, W_{lj} + W_{lk} = 1, W_{ij} + W_{ik} = 1, W_{lj} + W_{ij} = 1) \\ = \frac{\exp[(X_{lj} - X_{lk}) - (X_{ij} - X_{ik})]^\top \beta}{1 + \exp[(X_{lj} - X_{lk}) - (X_{ij} - X_{ik})]^\top \beta}, \end{aligned}$$

where the individual parameters are no longer present! This strategy offers a (quasi-)likelihood function that can be used to estimate  $\beta$ . Since it does not depend on the incidental parameters, it is not restricted by the constraints on sparsity that the previous results rely on, but it uses fewer data points. Jochmans (2018) obtains the large sample properties for the estimator thus defined. Under regularity assumptions, he shows that  $\hat{\beta}$  is asymptotically normal (without the bias term) and  $\|\hat{\beta} - \beta_0\| = O_p([N(N-1)p_N]^{-1/2})$ , where  $p_N$  is the proportion of quadruples, like those in **Figure 1**, that contribute to the estimation objective function. Intuitively, the more of those there are, the more precise the estimator is. He offers simulations where the (conditional) maximum likelihood estimator defined above is shown to perform well even for sparser networks.

Graham (2017) investigates a similar model, also with observed covariates and logistic idiosyncratic errors, but in an undirected network of the form

$$W_{ij} = W_{ji} = \mathbf{1}(X_{ij}^\top \beta + \alpha_i + \alpha_j + \epsilon_{ij} \geq 0). \quad 4.$$

There, the links are established with probability proportional to  $\exp(\beta^\top X_{ij} + \alpha_i + \alpha_j)$  and can be interpreted as the outcome from a pairwise stable arrangement with transfers (see discussion below). Here, the distinction between sender gregariousness and receiver attractiveness for a given node is moot since the link is not directional, but the individual parameters  $\alpha_i$  can be seen as the propensity by node  $i$  to establish connections. As in the articles discussed above, such parameters are also treated as fixed effects. In the absence of covariates, this model was christened the  $\beta$ -model by Chatterjee et al. (2011), who analyze its features. Yan & Xu (2013) establish the large sample properties for its maximum likelihood estimator (for sufficiently dense networks). In the presence of covariates, Graham (2017) analyzes the large sample properties for the maximum likelihood estimator for  $\beta$  (and  $\alpha$ ) as in the articles above, but he also offers a conditional maximum likelihood estimator for  $\beta$  constructed using sufficient statistics for  $\alpha$ , which allows him to “condition those parameters out” and to circumvent the incidental parameters problem in estimating  $\beta$ . Similarly to the directed case, the degree sequence here is such a sufficient statistic.<sup>9</sup> As is the case there, the enumeration of possible degree sequences becomes computationally complex even at moderately sized networks, and Graham also relies on tetrads to offer a conditional maximum likelihood estimator, the tetrad logit, that does not depend on the incidental parameters. The tetrad logit estimator is asymptotically normal as  $N$  grows and is able to tackle sparser models, albeit at a slower convergence rate. While for dense models (where network density converges to a constant as  $N$  grows) the convergence rate is parametric ( $\|\hat{\beta} - \beta_0\| = O_p([N(N-1)/2]^{-1/2})$ ), it is

<sup>9</sup>This property is also used by Pelican & Graham (2019) to test for the presence of externalities in the model, i.e., to verify whether a link between two nodes depends on edges between those nodes and other vertices, conditional on covariates and individual effects.

$O_p([N(N-1)/2]^{-1/4})$  for sparse network sequences where the density converges to zero at a rate proportional to  $N$ .

Recent efforts have focused on extending the models above in several directions. Toth (2018) and Gao (2020), for example, study semiparametric versions of the undirected network model in Equation 4 without a distributional assumption on  $\epsilon_{ij}$ . Toth (2018) examines the identification and estimation of this variation of the model, adapting ideas related to Han's (1987) (semiparametric) maximum rank correlation estimator. Gao (2020) studies the semiparametric identification of the model but also discusses nonseparable models where the index  $X_{ij}^\top \beta + \alpha_i + \alpha_j + \epsilon_{ij}$  is replaced by more general forms like  $\phi(X_{ij}, \alpha_i, \alpha_j) + \epsilon_{ij}$ . Other extensions veer a bit farther from the model in Equation 4: Shi & Chen (2016), for example, examine the estimation of a double-hurdle model<sup>10</sup> in which

$$W_{ij} = W_{ji} = \mathbf{1}(X_{ij}^\top \beta + \alpha_i + \epsilon_{ij} \geq 0) \times \mathbf{1}(X_{ji}^\top \beta + \alpha_j + \epsilon_{ji} \geq 0).$$

The condition above can be interpreted as a pairwise stability requirement on a model where individual utility from links depends additively on every direct connection and there are no externalities (see discussion below). Here the link is observed whenever it is beneficial to both parties involved, and not if it is detrimental to at least one of them. A related model (without individual effects) is considered by Comola & Fafchamps (2014) to study whether (directionally) elicited links are unilaterally or bilaterally formed.

#### 4. BEYOND DYADS

While dyadic models provide an important angle on the study of link formation, it is plausible that a connection between two nodes depends on liaisons with other nodes in the group. In a directed network, one such specification would be

$$W_{ij} = \mathbf{1} \left( \gamma \sum_{k \neq j} W_{jk} + \epsilon_{ij} \geq 0 \right), \quad 5.$$

where for simplicity I abstract from covariates  $X_{ij}$  and individual fixed effects  $\alpha$ . The individual link formation may also depend on the remaining connections in the network in more general ways. Here, node  $i$  sends an edge to  $j$ , taking into account those to whom  $j$  links (including whether  $j$  is connected to  $i$ ). One relevant feature of the model above is the econometric endogeneity issue that derives from having  $\sum_{k \neq j} W_{jk}$  as a covariate. If each of the elements in this sum is simultaneously determined according to an expression like Equation 5, this covariate will be related to  $\epsilon_{ij}$ . This would require caution in using partial information methods (focused on dyadic likelihoods, for instance) to estimate the parameters above. It might still nonetheless allow for the use of full information methods, if a joint distribution for all the links is consistent with the conditionally specified model in Equation 5. Importantly, though, such models are well known to be incomplete: They may not produce a unique joint distribution for  $W_{ij}$ ,  $i, j = 1, \dots, N$  (see footnote 5). If there are only two nodes, for example, we obtain  $W_{12} = \mathbf{1}(\gamma W_{21} + \epsilon_{12} \geq 0)$  and  $W_{21} = \mathbf{1}(\gamma W_{12} + \epsilon_{21} \geq 0)$ . If  $\gamma$  is positive, so that there is a tendency to reciprocate links, there are realizations of  $\epsilon_{12}$  and  $\epsilon_{21}$  for which the model predicts both  $(W_{12}, W_{21}) = (0, 0)$  and  $(W_{12}, W_{21}) = (1, 1)$  as solutions to this system of equations. In the statistical literature, this corresponds to the conditionally specified

<sup>10</sup>Traditionally, multiple hurdle models refer to variations of work by Cragg (1971). In the binary outcome context, the model is related to the partially observable model presented by Poirier (1980).

model being incompatible (see Arnold & Press 1989). In econometrics, such models are referred to as incomplete, as they lack a selection protocol between  $(W_{12}, W_{21}) = (0, 0)$  and  $(W_{12}, W_{21}) = (1, 1)$  in the example above, and are commonly found in the econometric study of strategic interactions with discrete actions (see de Paula 2013).

A modeling alternative that avoids these issues relies on temporal restrictions and postulates, for example, that

$$W_{ij,t} = \mathbf{1} \left( \rho W_{ij,t-1} + \gamma \sum_{k \neq j} W_{jk,t-1} + \epsilon_{ij,t} \geq 0 \right).$$

Here, the edge  $W_{ij,t}$  between  $i$  and  $j$  in period  $t$  depends on lagged network features  $W_{ij,t-1}$  and  $\sum_{k \neq j} W_{jk,t-1}$ , which are predetermined with respect to the innovation  $\epsilon_{ij,t}$ . When empirically adequate, this specification circumvents the simultaneity issue and allows for the addition of other covariates  $X_{ij,t}$ , while requiring repeated observations on the network.<sup>11</sup> It is also possible to include individual parameters  $\alpha$ , though proper attention should be given to the dynamic panel data nature of the setting in this case (see, e.g., Arellano & Honoré 2001, section 8, for the conventional panel data treatment). Such a model (with individual effects) is indeed contemplated by Graham (2016), in whose work similar conditioning arguments as those by Charbonneau (2017) or Graham (2017), though conditional on different events, are used to examine (semiparametric) identification and estimation on undirected networks (see also Han et al. 2019).

Another strategy to bypass the simultaneity issues referred to above is to model not only dyads, but also triads, tetrads, and more in general  $k$ -tuples, directly. Chandrasekhar & Jackson (2016) propose one such framework, which they call the subgraph generation model (SUGM). The model specifies a set of  $K$  subgraph classes, potentially involving more than two nodes each and including probabilities for each of these two nodes. If  $K = 2$ , for example, one of those subgraph classes could be taken to be the possible graphs between two nodes and the second to be the possible networks among three vertices. One then forms subgraphs at random from each of these classes among all the tuples in the group, and the final network is obtained by taking the union of all the edges thus sampled. Some of these edges may be redundant, as a link between two nodes could have appeared as a draw from both the two-node subgraph class and the three-node one. The subgraphs may also get meshed together: A complete subnetwork among nodes  $i, j$ , and  $k$  could arise as a genuine subgraph on that triad or as three independent complete graphs on the three pairs involving those three vertices. Disentangling the count of subgraphs in the model that are genuinely formed or occur by happenstance from the composition of other subgraphs can be done by noting that the counts of possible subgraphs are a mixture of both genuinely and incidentally formed subgraphs. This renders a system of equations that can then be solved for the parameters of interest and potentially matched to moments for the desired parameters.

#### 4.1. Strategic Formation Models

An alternative is to directly handle the challenges highlighted above. Since econometric models with interacting agents are natural platforms where such difficulties arise, much of the work speaking to such issues corresponds to estimable game theoretic frameworks. (While the economic narrative is not as salient in the statistical models described previously, many of those can also be interpreted through the lens of behavioral models.)

<sup>11</sup>The reader is referred to the related discussion of temporal lags in social interactions models and their identification by Manski (1993, section 4). Bramoullé et al. (2020) provide a review of that literature in this volume.

The first step in framing the network formation as a game involves a specification for the set of players, their actions, and their payoffs. In the present setting, the players are represented by the group of nodes  $\mathcal{N}_g$  to be eventually connected in equilibrium, corresponding for example to individuals, households, or firms. Their actions are in turn related to the formation of links. When the model is aimed at a directed network characterization, actions encode whether a node sends a link to another node. In undirected environments, on the other hand, the connections might need to be agreed upon. For a given graph  $g \equiv (\mathcal{N}_g, \mathcal{E}_g)$  designating how players are connected, a node  $i$  is assigned a network-dependent payoff. Whereas the payoff specification will depend on the context, a common parameterization (on an undirected graph) (see, e.g., de Paula et al. 2018) is given by

$$U_i(g) \equiv \sum_{j \neq i} W_{ij} \times (u + \epsilon_{ij}) + |\cup_{j: W_{ij}=1} N_j(g) - N_i(g) - \{i\}|v + \sum_j \sum_{k > j} W_{ij} W_{ik} W_{jk} \omega, \quad 6.$$

where  $W_{ij} = 1$  if  $i$  and  $j$  are connected,  $N_i(g)$  denotes the set of nodes directly connected to node  $i$  in the graph, and  $|\cdot|$  is the cardinality of a given set. The first term on the right-hand side of Equation 6 registers the payoff from direct connections and involves the parameter  $u$  and idiosyncratic variables  $\epsilon_{ij}$ , unobserved by the econometrician. The second term encodes the utility obtained from indirect connections:  $|\cup_{j: W_{ij}=1} N_j(g) - N_i(g) - \{i\}|$  is the number of individuals connected to direct counterparts of  $i$  but not directly connected to  $i$ . Finally, the last term summarizes any benefits accruing from two direct connections also being connected and induces incentives for clustering, a commonly observed phenomenon. While I omit observable covariates in the expression above, the parameters  $(u, v, \omega)$  can also be made to depend on those.

It is also important to clearly establish the information structure. In the literature so far (and in accordance with the statistical models discussed above), it is often assumed that information is complete, so agents are informed about others' (observed and unobserved) payoffs and incentives perfectly. Incomplete information, albeit possibly more plausible and epistemically more adequate in certain contexts, has less often been analyzed in the literature.

Another feature of the environment that requires attention relates to transferability. This refers to the possibility for agents to transfer payoffs among themselves, not only monetarily but also through other means. When available, this possibility allows nodes, for example, to bid for their preferred counterparts by accepting lower payoffs themselves. At two opposite ends of the spectrum are nontransferable utility (NTU) models, when there is no technology enabling agents to decrease their utility to benefit a potential partner, and transferable utility (TU) models, which allow transfers of utility at a constant exchange rate, and the total gain from the matching (surplus) is what matters for the stability of the relationship.<sup>12</sup> Which one is adopted depends, again, on the context at hand.

In closing the model, one then relies on a solution concept prescribing how individual behaviors are aggregated to generate an equilibrium network. Whereas traditional concepts in game theory (e.g., Nash equilibrium) can be envisioned and adapted to this case, the theoretical literature has offered additional notions to better capture the peculiarities of certain network formation contexts. When modeling undirected networks, for example, Jackson & Wolinsky (1996) propose pairwise stability as an alternative solution concept. A network  $g$  is pairwise stable according to Jackson & Wolinsky (1996) if any link present in  $g$  is mutually beneficial and any absent link is detrimental

<sup>12</sup>Another possibility is the intermediate scenario with imperfect transferable utility (ITU), where transfers are allowed but at an exchange rate between individual utilities that is not constant and is possibly endogenous to the economic environment. While this would also be categorized as transferable utility, the conventional terminology focuses on the constant exchange rate case (see Chiappori 2019).

to at least one of the parties involved. More formally, we have

$$\forall ij \in g, U_i(g) \geq U_i(g_{-ij}) \text{ and } U_j(g) \geq U_j(g_{-ij})$$

and

$$\forall ij \notin g, U_i(g) > U_i(g_{+ij}) \text{ or } U_j(g) > U_j(g_{+ij}),$$

where  $ij \in (\notin)g$  signifies that the link between  $i$  and  $j$  pertains (or not) to the set of edges in  $g$ .<sup>13</sup> The network  $g_{-ij}$  is  $g$  without the link between  $i$  and  $j$ , and  $g_{+ij}$  is the network  $g$  with the link between  $i$  and  $j$ . The theoretical literature has contemplated several variations to this stability concept (e.g., pairwise Nash stability, strong stability) (see the discussion in Jackson 2009), and Bloch & Jackson (2006) adapt this solution concept to an NTU environment.

For a given parameter vector and observable covariates, the model translates the distribution of unobservable variables (e.g.,  $\epsilon_{ij}$ ) into a probability distribution over the equilibrium sets for the game described above. Since this probability distribution is indexed by the parameter vector, this delivers a statistical model on which one can in principle perform estimation and inference. One difficulty in doing this is that there might be more than one stable network for a given parameter value and given realizations of observable and unobservable random variables. This is related to the statistical difficulties highlighted earlier in this section and is potentially problematic for identification (i.e., the reverse mapping between observed distributions and parameters), computation, and inference.

Another possibility is to rely on iterative procedures in which the presence or absence of links is evaluated as individuals or pairs take turns in a random meeting protocol, as in the preferential attachment framework mentioned earlier, in which agents form links sequentially and the establishment of new links is more likely with higher-degree existing nodes (see Barabási & Albert 1999).<sup>14</sup> Before discussing models relying on the equilibrium notions above, I discuss alternatives relying on iterative protocols.

**4.1.1. Iterative network formation.** Mele (2017), Badev (2018), and Christakis et al. (2020) are notable examples of strategic network formation models based on a stochastic meeting protocol whereby individuals or pairs sequentially revise their links. While such protocols are not meant to be directly fit to the data, the random meeting sequences and unobservable errors guiding the decisions to establish or interrupt connections lead to a potentially estimable distribution over networks. Mele (2017) studies a directed network, and Badev (2018) expands the analysis there to the joint determination of links and behaviors.<sup>15</sup> Christakis et al. (2020), on the other hand, examine an undirected network. The last two works present empirical applications to links among adolescents using the Add Health data, and Mele (2020) applies the methodology developed in previous work (Mele 2017) to examine segregation in high schools, also using Add Health data.

<sup>13</sup>This is related to, but different from, the stability concept typically used in marriage market models.

<sup>14</sup>It is worth noting that such stochastic revision processes are not unrelated to the (noniterative) equilibrium notions mentioned earlier and discussed below. For example, in an NTU setting, Jackson & Watts (2002) demonstrate that a process where pairs meet sequentially and are offered to (myopically) form or maintain links that are mutually beneficial and to dissolve links that are not beneficial to at least one of the parties involved converges to a pairwise stable network or a cycle.

<sup>15</sup>The paper by Mele (2017) was previously circulated as “A structural model of segregation in social networks” (in 2015). Previous versions of the paper by Badev (2018) were circulated as “Discrete games with endogenous networks: theory and policy” (in 2013) and “Discrete games in endogenous networks: equilibria and policy” (in 2017).

Mele (2017) models a directed network where  $W_{ij} = 1$  if individual  $i$  offers a link to individual  $j$  and  $W_{ij} = 0$  otherwise. The utility function for individual  $i$  is given by

$$U_i(g) \equiv \sum_{j \neq i} W_{ij} u_{ij}^\theta + \sum_{j \neq i} W_{ij} W_{ji} m_{ij}^\theta + \sum_{j \neq i} W_{ij} \sum_{k \neq i, j} W_{jk} v_{ik}^\theta + \sum_{j \neq i} W_{ij} \sum_{k \neq i, j} W_{ki} v_{kj}^\theta,$$

where  $u_{ij}^\theta \equiv u(X_i, X_j; \theta)$  represents the utility from directly linking to other individuals, and the first term involving  $v_{ij}^\theta \equiv v(X_i, X_j; \theta)$  encodes the utility from indirectly linking to a friend's friend. These parameters then play the same role as  $u$  and  $v$  in the utility function previously introduced. Since this is a directed network,  $m_{ij}^\theta \equiv m(X_i, X_j; \theta)$  marks the utility from a mutual, reciprocated link. (In undirected networks, links are reciprocated by definition and this term does not show up in the utility function presented earlier in this review.) The second term involving  $v_{ij}^\theta$  internalizes some of the impact a link to individual  $j$  generates for individuals that had offered links to  $i$ .<sup>16</sup> Utility is nontransferable and information is complete.

There is then a meeting sequence  $m = \{m^t\}_{t=1}^\infty$ , where  $m^t = (i, j)$  means that  $i$  can offer or dissolve a link to  $j$  in iteration  $t$ . Define  $\mathbb{P}(m^t = ij | g^{t-1}, X) = \rho(g^{t-1}, X_i, X_j)$ , where  $g^{t-1}$  is the network in iteration  $t - 1$ . In addition, the meeting probability between  $i$  and  $j$  does not depend on the existence of a link between them, and each meeting has a positive probability of occurring [i.e.,  $\rho(g^{t-1}, X_i, X_j) = \rho(g_{-ij}^{t-1}, X_i, X_j) > 0, \forall ij$ ]. This ensures that the likelihood function from this model does not depend on the meeting protocol.

Finally, whenever individual  $i$  is offered a meeting with another node, it is supposed that they receive idiosyncratic shocks ( $\epsilon_1, \epsilon_0$ ) to the utility of forming a link ( $W_{ij} = 1$ ) or not ( $W_{ij} = 0$ ). An edge to  $j$  is established if and only if

$$U_i(W_{ij}^t = 1, g_{-ij}^t, X; \theta) + \epsilon_{1t} \geq U_i(W_{ij}^t = 0, g_{-ij}^t, X; \theta) + \epsilon_{0t}.$$

While there were no unobservable errors up to this point, once given the choice the decision by  $i$  acts much like a standard random utility model over the remaining  $|\mathcal{N}_g| - 1$  nodes. (This implies that the average degree is roughly proportional to  $|\mathcal{N}_g|$  and the model produces a dense network.) This revision process then leads to subsequent additions and deletions of edges, and the resulting process forms a Markov chain on networks  $\{g^t\}$ . In the absence of unobservable shocks (i.e., if  $\epsilon_0$  and  $\epsilon_1$  are zero with probability one), the chain converges to one of the Nash equilibria for the game without  $\epsilon$ s. Under the assumption that the  $\epsilon$ s are distributed independently (across time and links) and follow an extreme value type I distribution, the chain converges instead to a unique stationary distribution,

$$\pi(g, X; \theta) = \exp[Q(g, X; \theta) - \mathcal{C}(\theta)],$$

where  $\mathcal{C}(\theta) = \ln\{\sum_{g'} \exp[Q(g', X; \theta)]\}$  and

$$Q(g, X; \theta) = \sum_{(i,j)} W_{ij} u_{ij}^\theta + \sum_{(i,j)} W_{ij} W_{ji} m_{ij}^\theta + \sum_{(i,j,k)} W_{ij} W_{jk} v_{ik}^\theta.$$

As noted by Mele (2017), the Nash equilibria for the game with payoffs thus defined correspond to the maxima of the function  $Q(g, X; \theta)$ .

This distribution can in principle be taken to data on either one or more networks. In fact, the stationary distribution above describes an exponential random graph model, discussed previously.

<sup>16</sup>The author refers to this as popularity: "If individual  $i$  forms a link to  $j$ , he automatically creates an indirect link for all the agents that already have a link to  $i$ . Thus  $i$  generates an externality (positive or negative) for each  $k$  that formed a link to him in previous periods. This externality makes  $i$  more or less popular" (Mele 2017, p. 829).

The model is estimated by Bayesian methods, producing a posterior distribution over parameters of interest. One important difficulty with such models is nonetheless the computation of the normalizing constant  $\mathcal{C}(\theta)$  in  $\pi(g, X; \theta)$ , as the summation is over all possible directed networks between the individuals in the group. With 10 individuals, for example, there are  $2^{90} \approx 10^{27}$  such network configurations.<sup>17</sup> This is relevant because it is necessary to compute such denominators to solve for the posterior distribution either analytically or numerically (e.g., via simulations of the posterior distribution).

Alternative strategies to circumvent this issue or to approximate the denominator in other settings include pseudo-likelihood methods (Besag 1975, Strauss & Ikeda 1990) and variational principles (see Jordan & Wainwright 2008 and Mele & Zhu 2019 for an application in econometrics). These and other protocols are briefly discussed by de Paula (2017). Mele (2017) instead handles the situation by using simulation methods (Markov chain Monte Carlo, or MCMC) (see also Kolaczyk 2009 and references therein). The first challenge is the computation of  $\mathcal{C}(\theta)$  for a given parameter value  $\theta$ . Here, it is possible to design a Metropolis-Hastings algorithm to simulate the distribution of networks without the normalizing constant.

Unfortunately, this simulation protocol is itself not immune to problems. It is well known, for example, that parameter changes in ERGMs may lead to abrupt changes in probable graphs (see Snijders 2002). In addition, in parameter regions where the distribution over networks is multimodal, the convergence of the algorithm is impractically slow for MCMC protocols where networks change only locally from iteration to iteration (see Bhamidi et al. 2011). Here, the modes of the distribution will correspond to the maxima of the function  $Q$ , which are in turn related to the Nash equilibria for the game discussed previously. Hence, while the stationary distribution for the Markov process defined by the meeting protocol is unique, equilibrium multiplicity is also a complicating computational feature here. To accelerate convergence, the article suggests a simulation algorithm that updates networks at larger steps (see Mele 2017, supplemental appendix B.1).

For parameter regions where the distribution is unimodal, on the other hand, Bhamidi et al. (2011) and Chatterjee & Diaconis (2013) show that graph draws are indistinguishable from random network models with independent link formation (i.e., an Erdős-Rényi model) or a mixture of such models. Mele (2017) also shows that a similar phenomenon occurs in a simplified version of his model without covariates and positive utility from mutual links ( $m > 0$ ). This points to an identification issue when estimation is based on a single network in this particular example.<sup>18</sup> The result does not allow for covariates, but the author conjectures that the sign of such externalities (on reciprocated links) will remain relevant.

In the Bayesian estimation, one also needs to approximate the posterior distribution for the parameters of interest. The “inner” simulation of the networks for a given parameter vector discussed above is a component of the “outer” simulation protocol for the posterior distribution presented by Mele (2017). Here, too, the article follows an MCMC procedure based on the exchange algorithm proposed by Murray et al. (2006) to avoid the computation of the normalizing constant. The model is taken to data by Mele (2020), who studies ethnic segregation in high schools using data from Add Health.

**4.1.2. Noniterative network formation.** An alternative strategy takes a noniterative perspective whereby network formation is obtained as a simultaneous move game with players as

<sup>17</sup> Similar issues also appear in models like those in Equation 3, which also belong to the exponential family.

<sup>18</sup> When multiple networks are used, identification is attained with variation in sufficient statistics across networks: “My model is identified in the many networks framework under usual regularity conditions, because the likelihood belongs to the exponential family” (Mele 2017, p. 828).

nodes/vertices and the action space as potential links. For example, Leung (2015b) studies a directed network formation model under incomplete information and Bayes-Nash as the solution concept, adapting two-step estimation strategies commonly used in the econometrics of such games (see de Paula 2013). He employs the model to analyze trust networks in India.<sup>19</sup> Also modeling a directed network, Gualdani (2019) examines a complete information model employing Nash equilibrium as the solution concept and studies board interlocks among firms. One important challenge in this context relates to the dimensionality of potential networks, and the article offers suggestions to reduce the computational burden. Here I focus on works using pairwise stability (without transferability).

One avenue is to adapt some of the strategies used in the empirical games literature. To illustrate this, take the simplified directed network formation setting discussed previously where there are two nodes, 1 and 2, and the payoff obtained by  $i = 1, 2$  from a link to  $j \neq i$  is given by  $\gamma + \epsilon_{ij}$  if  $j$  also offers a link to  $i$  and  $\epsilon_{i2}$  otherwise. If  $i$  does not send a link to  $j$ ,  $i$ 's payoff is zero. The payoffs for  $j$  are analogously defined. A Nash equilibrium for this game leads to an econometric model where  $W_{12} = \mathbf{1} (\gamma W_{21} + \epsilon_{12} \geq 0)$  and  $W_{21} = \mathbf{1} (\gamma W_{12} + \epsilon_{21} \geq 0)$ . As discussed earlier, if there is a tendency to reciprocate links (i.e.,  $\gamma$  is positive), there are realizations of  $\epsilon_{12}$  and  $\epsilon_{21}$  for which the model predicts both  $(W_{12}, W_{21}) = (0, 0)$  and  $(W_{12}, W_{21}) = (1, 1)$  as solutions to the system of equations. While this does not tie down a single probability distribution for  $(W_{12}, W_{21})$ , it nonetheless offers bounds on the probabilities for the possible realizations of this vector. The probability for the event  $\{(W_{12}, W_{21}) = (0, 0)\}$  is minimized if  $(0, 0)$  is never selected when  $\epsilon_{12}$  and  $\epsilon_{21}$  are such that other possible equilibria exist. It is on the other hand maximized when  $(0, 0)$  is always chosen when there are other possible solutions for those realizations of  $\epsilon_{12}$  and  $\epsilon_{21}$ . These quantities thus bound the probability for the event  $\{(W_{12}, W_{21}) = (0, 0)\}$ , and those inequalities in turn provide information on the parameters guiding the data-generating process. Parameters that lead to lower probability bounds above the observed frequency in the data, or to upper probability bounds below it, are not consistent with the generating process. At the same time, more than one parameter will typically be consistent with the data, rendering the model partially identified (see Tamer 2010).<sup>20</sup>

To illustrate how this translates into an alternative solution concept, consider, for instance, an undirected network formation game among three individuals ( $|\mathcal{N}_g| = 3$ ) with nontransferable utility and complete information. Let payoffs be given by

$$U_i(g) \equiv \sum_{j \in 1, \dots, n, j \neq i} \delta^{d(i,j;g)-1} (1 + \epsilon_{ij}) - |N_i(g)|,$$

where  $d(i,j;g)$  is the minimum distance between  $i$  and  $j$  in the graph  $g$ ,  $0 < \delta < 1$ ,  $\epsilon_{ij}$  are unobservable (to the researcher) preference shocks, and  $|N_i(g)|$  is the number of direct connections of individual  $i$  in graph  $g$ . Assume that observed networks are pairwise stable. For  $\epsilon_{ij} = \epsilon_{ji}$ ,  $0 < \epsilon_{23} < \delta/(1 - \delta)$ , we can represent the possible pairwise stable equilibria in the space of unobservables as shown in **Figure 2**.

The approach above would produce bounds on  $\delta$  corresponding to the probability that a particular network is pairwise stable (though possibly not unique) (upper bound) and the probability that it is the unique pairwise stable network (lower bound). In **Figure 1** (which does not comprise

<sup>19</sup>In an earlier working paper taking a similar modeling framework, Gilleskie & Zhang (2009) also examine smoking behavior among network members. Ridder & Sheng (2015) and Candelaria & Ura (2018) also study incomplete-information network formation games.

<sup>20</sup>For additional strategies to handle the multiplicity problem, the reader is referred to de Paula (2013).



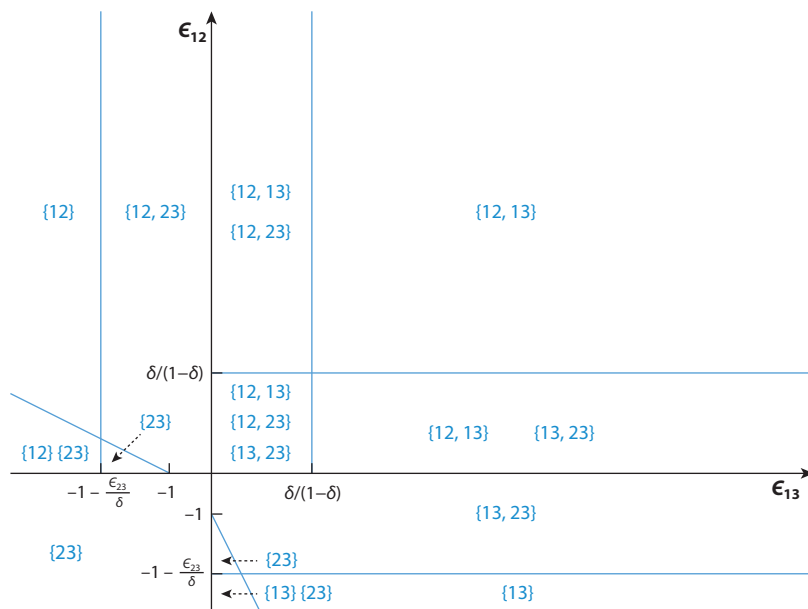


Figure 2

Pairwise stable networks.

the whole space for  $\epsilon$ s), those bounds for the network  $\{12, 13\}$  would be

$$\mathbb{P}(\epsilon_{12}, \epsilon_{13} \geq 0) \geq \mathbb{P}(\{12, 13\}) \geq \mathbb{P}(\epsilon_{12}, \epsilon_{13} \geq \delta/(1 - \delta)),$$

and one could form similar bounds for all (eight) possible networks (exploring the whole space of unobservables). Unfortunately, while this might be conceivable when only three individuals are involved, for even a moderate number of nodes, the number of networks one would need to consider as potential equilibria is computationally intractable. With more than 24 individuals, for instance, there are more potential networks than atoms in the observable universe.

Given this dimensionality issue, one possible avenue to reduce the computational burden is to focus on smaller subnetworks involving subgroups of individuals. Sheng (2018) develops this strategy in a model for undirected networks. She focuses on pairwise stable networks (with either transferable or nontransferable utility) under complete information and a utility structure given by

$$U_i(g) \equiv \sum_{j \neq i} W_{ij}(n_{ij}^\theta + \epsilon_{ij}) + \frac{1}{|\mathcal{N}_g| - 2} \sum_{j \neq i} W_{ij} \sum_{k \neq i, j} W_{jk} v + \frac{1}{|\mathcal{N}_g| - 2} \sum_{j, k \neq i} W_{ij} W_{ki} W_{jk} \omega.$$

As before,  $u_{ij}^{\theta} \equiv u(X_i, X_j; \theta)$  records the direct utility obtained from linking to other individuals in the group,  $v$  provides benefits to indirect friendships, and  $\omega$  indicates the additional payoff of having common connections that are also linked to each other.<sup>21</sup> The unobservable random variables  $\epsilon$  are independent and follow a known distribution. (Since there are preference shocks

<sup>21</sup>In contrast to the utility specification presented at the beginning of this review, the specification above normalizes the number of connections by  $|\mathcal{N}_g| - 2$  and allows for some double counting: If  $i, j$ , and  $k$  are all connected, each accrues benefits from direct connection with two individuals, indirect connections from each of the other two, and the fact that the other two are also linked to each other.

for each potential link, isolated individuals are unlikely when  $|\mathcal{N}_g|$  is large.) The data are presumed to come from a sample of independent and identically distributed networks.

Pairwise stability (with either transferable utility or not) is determined by the marginal utility of a link:  $\Delta U_{ij}(g) = U_i(g) - U_i(g_{-ij})$  when  $ij \in g$ , and  $\Delta U_{ij}(g) = U_i(g_{+ij}) - U_i(g)$  when  $ij \notin g$ . Let then  $\text{PS}(\Delta U(X, \epsilon))$  denote the set of pairwise stable networks for realizations of  $X$  and  $\epsilon$ . The probability that one observes  $g$  is then given by

$$\begin{aligned} \mathbb{P}(g|X) = & \int_{g \in \text{PS}(\Delta U(X, \epsilon)) \wedge |\text{PS}(\Delta U(X, \epsilon))|=1} dF(\epsilon) \\ & + \int_{g \in \text{PS}(\Delta U(X, \epsilon)) \wedge |\text{PS}(\Delta U(X, \epsilon))|>1} \lambda(g|\text{PS}(\Delta U(X, \epsilon))) dF(\epsilon), \end{aligned}$$

where  $\lambda$  is the probability that  $g$  is selected for realizations of  $\epsilon$  that allow for multiple pairwise stable networks. One can then emulate our previous discussion on empirical games to generate bounds on the probability and guide the computation of identified parameter sets. As previously discussed, these bounds are unfortunately computationally intractable for even moderately sized networks. Instead, Sheng (2018) suggests looking at subgraphs on  $A \subset \mathcal{N}_g$  nodes: the network  $g_A$  comprising  $A$  and all edges in  $g$  linking those nodes. The probability of observing such a subnetwork is then given by

$$\begin{aligned} \mathbb{P}(g_A|X) = & \int_{g_A \in \text{PS}_A(\Delta U(X, \epsilon)) \wedge |\text{PS}_A(\Delta U(X, \epsilon))|=1} dF(\epsilon) \\ & + \int_{g_A \in \text{PS}_A(\Delta U(X, \epsilon)) \wedge |\text{PS}_A(\Delta U(X, \epsilon))|>1} \sum_{g_{-A}} \lambda(g_A|\text{PS}(\Delta U(X, \epsilon))) dF(\epsilon), \end{aligned}$$

where the addition over  $g_{-A}$  sums over the collection of complementary nodes ( $\mathcal{N}_g \setminus A$ ) and edges in  $g$  connecting them to each other and to nodes in  $A$ .  $\text{PS}_A(\Delta U(X, \epsilon))$  is the subset of networks in  $A$  that are part of a network in  $\text{PS}(\Delta U(X, \epsilon))$ . These deliver

$$\int_{g_A \in \text{PS}_A(\Delta U(X, \epsilon)) \wedge |\text{PS}_A(\Delta U(X, \epsilon))|=1} dF(\epsilon) \leq \mathbb{P}(g_A|X_A) \leq \int_{g_A \in \text{PS}_A(\Delta U(X, \epsilon))} dF(\epsilon),$$

where the upper bound is the probability that the subnetwork  $g_A$  on nodes  $A$  pertains to a pairwise stable network, and the lower bound is the probability that the only subnetwork on nodes  $A$  pertaining to a pairwise stable network is  $g_A$ .

Consider for example the network formation game on three individuals presented above. For the subnetwork  $\{12\}$  on nodes 1 and 2, the upper bound is given by the probability that  $\epsilon_{12}$  is greater than zero minus the probability for the triangular region where only  $\{23\}$  is pairwise stable. In this region,  $\{12\}$  always pertains to a pairwise stable network. The lower bound is given by the probability that  $\epsilon_{12}$  is greater than  $\delta/(1 + \delta)$  and  $\epsilon_{13}$  is greater than zero, plus the probability that  $\epsilon_{13}$  is less than zero, minus the subregions where  $\{23\}$  is also a pairwise stable network. In this region, the only network on nodes 1 and 2 that is part of a pairwise stable network is  $\{12\}$ —even though there are multiple pairwise stable networks in the region where  $\epsilon_{12}$  is greater than  $\delta/(1 - \delta)$  and  $\epsilon_{13}$  is between zero and  $\delta/(1 - \delta)$ .

In the article, Sheng (2018) imposes exchangeability restrictions on equilibrium selection and payoff primitives that guarantee that these bounds are nontrivial even as the number of individuals in the groups gets larger. Among other things, such exchangeability restrictions also imply a dense network (as in Mele 2017): The total number of links is  $O_p(N^2)$  (see, e.g., Orbanz & Roy 2015). While the sets above are not sharp—i.e., more informative bounds on the parameters

can intuitively be obtained by considering subnetworks on a larger number of nodes—they are potentially computable.<sup>22</sup> Sheng offers an algorithm to perform such computation and indicates a few additional potential simplifications. For example, when  $v$  and  $\omega$  are nonnegative (which guarantees the existence of pairwise stable networks when utility is nontransferable), the game is supermodular and the solution set possesses a maximal and a minimal element. This can be leveraged to reduce the computational complexities here, as done by Miyauchi (2016) (also in the context of pairwise stable network formation) and other authors in the empirical games literature. A Monte Carlo study in Sheng's (2018) paper demonstrates the performance of the algorithm in a TU context with 50 and 100 networks varying in size from 25 individuals to 100.

Another use of subnetworks to circumvent the challenges presented in this setting is proposed by de Paula et al. (2018). The article works on a complete-information NTU model using pairwise stability as the solution concept. The treatment is tailored to handle large networks, and  $\mathcal{N}_g$  is an uncountable set with continuum cardinality, taken to be an approximation for a large group of individuals. The approach described here starts with a large network. Related approaches also focused on large networks, but taken as limits for finite sequences of networks, are those by Leung (2015a), Menzel (2016), and Boucher & Mourifié (2017). To capture sparsity, they restrict payoffs, allowing only for a finite number of links  $L$  such that in equilibrium one obtains a bounded degree graph on the continuum (sometimes referred to as a graphing). Utilities are also assumed to depend only on individual characteristics (and not identities) and on indirect connections only up to a finite distance (depth  $D$ ). This allows one to focus on network types defined by one's local neighborhood, whose cardinality is potentially more manageable than that of networks on individual nodes.

To illustrate this strategy, consider a very simple network formation game where individuals can only form one link and their utility depends only on this link. Here, both  $L$  and  $D$  are one. Nodes are characterized by  $X$ , which takes two values,  $B$  or  $W$ , and there is a continuum of individuals of each type,  $\mu_B$  and  $\mu_W$ . Outcomes are thus given by ordered pairs  $(x, y)$ , where  $x$  is the individual's characteristic and  $y$  the characteristic of their connection. Utilities are given by

$$U_i(g) \equiv u_{xy} + \epsilon_i(y),$$

where again  $x$  marks the individual's characteristic and  $y$  that of their counterpart. If no link is formed, the payoff is normalized to zero. Hence, there are four parameters ( $f_{x,y}$ , with  $x, y \in \{B, W\}$ ) and two preference shocks for each individual,  $\epsilon_i(B)$  and  $\epsilon_i(W)$ .<sup>23</sup>

De Paula et al. (2018) define network types as the local networks surrounding an individual in an observed network. Given the payoff structure, the network type should record payoff-relevant connections. In this simple case with just one individual, the relevant network type is given by the pair  $(x, y)$ . In a more general setting, a network type is characterized by the individual, their direct connections, their connections' direct connections, and so on, together with each

<sup>22</sup> Previous versions of the article, as reviewed by de Paula (2017), also consider simpler bounds:

$$\begin{aligned} \mathbb{P}(W_A = w_A | X_A) &\leq \int_{\exists W_{-A}: w_A \in \text{PS}(\Delta U_A(W_{-A}, X_A, \epsilon_A))} dF(\epsilon_A) \text{ and} \\ \mathbb{P}(g_A | X_A) &\geq \int_{\forall g_{-A}: g_A \in \text{PS}(\Delta U_A(g_{-A}, X_A, \epsilon_A)) \wedge |\text{PS}_A(\Delta U_A(g_{-A}, X_A, \epsilon_A))| = 1} dF(\epsilon_A). \end{aligned} \quad 7.$$

In words, the upper bound for  $\mathbb{P}(g_A | X_A)$  is the probability that subnetwork  $g_A$  is pairwise stable for some  $g_{-A}$ , and the lower bound is the probability that, for any  $g_{-A}$ , only subnetwork  $g_A$  is pairwise stable. These bounds do not require pairwise stability on the rest of the network and are thus easier to compute, but they are less informative than the ones above and might yield trivial bounds for larger groups.

<sup>23</sup> Additional restrictions are imposed on the preference shocks so that in equilibrium, even in large networks, there are isolated individuals.

of these nodes' characteristics up to the payoff depth  $D$ . This thus corresponds to a network on up to  $1 + L + L(L-1) + \dots + L(L-1)^{D-1} = 1 + L \sum_{d=1}^D (L-1)^{d-1}$  nodes and is equal to 2 when  $L = 1$ . The proportion of individuals of each network type is an equilibrium outcome, and one would like to verify which parameter values rationalize the type shares in the data as outcomes of a pairwise stable network.

To do this, the article classifies individuals based on which network types they would not reject. Depending on the preference shocks, a  $B$  individual may be content to have a  $W$  connection ( $f_{BW} + \epsilon_i(W) > 0$ ) but not a  $B$  connection ( $f_{BB} + \epsilon_i(B) < 0$ ). In this case, this individual would not have network type  $BB$  in equilibrium, as this would contradict pairwise stability, but would be content to have  $BW$ , or, should there be no  $W$  individuals to link to, to remain isolated as  $B0$ . The authors thus form preference classes collecting all types that would be acceptable to an individual with given realizations of the preference shocks. In this simple example, the preference class for the individual above would be  $\{BB, B0\}$ . (Since there are no connections to be dropped from an isolated network type, the isolated type is an element for every preference class.)

As discussed above, this corresponds to a partition of the space of unobservable shocks, but only for the individual. Given a distribution for the preference shocks  $\epsilon$ , one can compute (either analytically or numerically) the probability of each preference class for a given individual at a particular parameter value. One can then stipulate how individuals in each preference class are allocated to network types. In the paper by de Paula et al. (2018), this is done using allocation parameters designating the proportion of individuals in a particular preference class that are allocated to a network type. In the example above, there are four preference classes for a  $B$  individual:  $H_1 = \{B0\}$ ,  $H_2 = \{B0, BB\}$ ,  $H_3 = \{B0, BW\}$ , and  $H_4 = \{B0, BB, BW\}$ . Letting  $\alpha_H(t)$  denote the allocation proportion of individuals in preference class  $H$  to type  $t$ , the predicted share of  $BW$  individuals is given by  $\mathbb{P}(H_1|B)\alpha_{H_1}(BW) + \mathbb{P}(H_2|B)\alpha_{H_2}(BW) + \mathbb{P}(H_3|B)\alpha_{H_3}(BW) + \mathbb{P}(H_4|B)\alpha_{H_4}(BW)$  multiplied by the proportion of  $B$  individuals in the group.

The key here is to impose restrictions on the allocation parameters  $\alpha_H(t)$  to be satisfied for a given profile of network type shares to be consistent with pairwise stability. These restrictions are necessary conditions for pairwise stability in the article. For example, nodes can only be allocated to network types that pertain to their preference class: If  $t \notin H$ ,  $\alpha_H(t) = 0$ . This corresponds to the condition that links should be beneficial to both individuals and have  $\alpha_{H_1}(BW) = \alpha_{H_2}(BW) = 0$ . Second, given any pair of network types that could feasibly add a link to each other (i.e., an isolated individual of either characteristic,  $B$  or  $W$ , in the example), the measure of individuals who would prefer to do so must be zero for at least one of the types. This corresponds to the condition that nonexistent links should be detrimental to at least one of the individuals. Another way to express this condition is to require that the product of the measure of individuals of one type that would benefit from adding links to individuals of the other type should be zero. This translates into a quadratic objective function, which in equilibrium has to be zero once allocation parameters are adequately chosen. Finally, the predicted proportions of network types ought to match the observed proportions of types in the network, which in turn defines a set of linear constraints.

Once these observations are put together, verifying whether a parameter vector delivers the observed network type shares consistent with a pairwise network corresponds to solving a quadratic program on the allocation parameters with constraints requiring the allocation parameters to be positive, to add up to one, and to generate predicted shares matching the data. If the data are rationalized as a pairwise stable network for a given parameter vector, the optimized objective function is zero. The cardinality of the problem, while still nontrivial, is related to the cardinality of the preference classes and network types rather than the potential networks on  $|\mathcal{N}_g|$  individuals. If this is done at each putative parameter, one can then collect those parameters that rationalize the data

and form the identified set. Whereas the quadratic program above is regrettably not a convex one (since the matrix in the quadratic form is not necessarily positive definite), the article offers simulation evidence for networks as large as  $|\mathcal{N}_g| = 500$ . Anderson & Richards-Shubik (2019) apply this framework to the analysis of coauthorships in economics.<sup>24</sup>

One interesting point to note refers to the source of statistical uncertainty here. Since there are infinitely many nodes, were the network to be completely observed, network type shares would be perfectly measured. On the other hand, if there is sampling uncertainty in the measurement of network type shares because only a sample of individuals (and their network types) is collected, sampling uncertainty in type shares would transfer to the estimation of structural parameters. This relates to what is sometimes referred to as a design-based paradigm in statistics (see Kolaczyk 2009). The randomness here obtains from the probability ascribed by the survey scheme to the sampling of the various individuals in the network. The partial observability of the network is nonetheless not uncommon and is rather the norm in many settings.<sup>25</sup>

## 5. DISCUSSION

The study of network formation models is an active area of research. While this review presents a selective collection of recent developments in the field, several research questions of interest remain and are briefly discussed below.

Since networks mediate and have their formation informed by several outcomes of interest, econometric methods articulating both the formation of networks and the outcomes influenced by them are potentially important in several areas. Studies offering such articulation include those by Goldsmith-Pinkham & Imbens (2013) and Hsieh & Lee (2016), who use a dyadic network formation model and focus on educational achievement; Gilleskie & Zhang (2009), who use Bayes-Nash network formation and focus on smoking initiation; and Badev (2018) and Hsieh et al. (2019), who use exponential random graph network formation models along the lines of Mele (2017) to focus on smoking and, in the latter paper, also on academic achievement.<sup>26</sup> Multiplicity is not a salient issue in the econometric network formation protocols in these works (or in the econometric system of equations determining outcomes there). This nonetheless remains a possibility in other (strategic interaction or conditionally specified) econometric models of network formation. Multiplicity is also possible in nonlinear systems determining the outcomes given the networks. In those cases, partial identification in either stage (i.e., network formation or interactions) may be transmitted to other parameters in the model. This point is illustrated, for example, in the context of an empirical entry-exit game in industrial organization with potentially multiple equilibria studied by Ciliberto et al. (2018).<sup>27</sup>

Relatedly, there is also scope to expand on the catalog of equilibrium notions depending on the protocol for outcome determination and empirical setting of interest. For instance, Ho & Lee

<sup>24</sup>The article also offers a method to account for sampling uncertainty in the type shares.

<sup>25</sup>Because it is also focused on subnetworks, the approach developed by Sheng (2018) may also be adapted to such contexts.

<sup>26</sup>For an exposition on econometric methods for outcomes modulated by networks, the reader is referred to Bramoullé et al. (2020) in this volume. In those cases, the literature tends to implicitly assume econometric exogeneity between the unobservable variables determining outcomes and the econometric errors determining the networks that mediate those. This obtains, for instance, when the unobservables determining networks and outcomes are unrelated. If such a scenario is not empirically adequate, an instrumental variable, a control function, or a model for network formation may provide potential solutions (see, again, Bramoullé et al. 2020 for more details).

<sup>27</sup>Previous versions of this paper were circulated with the title “Inference on market power in markets with multiple equilibria.”

(2019) (see also Ghili 2018 and Liebman 2018) articulate network formation and a bargaining framework developed by Horn & Wolinsky (1988) to examine interactions between hospitals and insurers in the United States. Econometric analysis (and, to a certain extent, theoretical developments) around these and related models remains an area of interest.

Finally, the network itself may be an outcome of interest in a treatment effects context in which programs may lead to changes in the network, which in turn may modulate other outcomes of concern. Comola & Prina (2019), for example, study the effect of a savings product in Nepal on consumption through a randomized field experiment.<sup>28</sup> They find that insurance-motivated connections are likely to be rewired after the intervention, and taking this into account may improve one's understanding of the program's effect on consumption. They offer a treatment response framework to analyze this phenomenon, and further econometrics work on networks in a potential outcomes setting (possibly also accounting for issues such as multiplicity) would also be a useful area of study.

## DISCLOSURE STATEMENT

The author is not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

## ACKNOWLEDGMENTS

I would like to thank Marcel Fafchamps, Seth Richards-Shubik, Jörg Stoye, Martin Weidner, and a reviewer for comments on earlier drafts. Financial support from the Economic and Social Research Council ESRC grant RES-589-28-0001 to the Centre for Microdata Methods and Practice and from ESRC Large Research Grant ES/P008909/1 is gratefully acknowledged.

## LITERATURE CITED

- Anderson KA, Richards-Shubik S. 2019. *Collaborative production in science: an empirical analysis of coauthorships in economics*. Work. Pap., Carnegie Mellon Univ., Pittsburgh, PA
- Arellano M, Honoré BE. 2001. Panel data models: some recent developments. In *Handbook of Econometrics*, Vol. 5, ed. J Heckman, E Leamer, pp. 3229–96. Amsterdam: Elsevier
- Arnold BC, Press SJ. 1989. Compatible conditional distributions. *J. Am. Stat. Assoc.* 84:152–56
- Aronow P, Samii C, Assenova V. 2015. Cluster-robust variance estimation for dyadic data. *Political Anal.* 23:564–77
- Atalay E, Hortacsu A, Roberts J, Syverson C. 2011. Network structure of production. *PNAS* 108:5199–202
- Badev A. 2018. *Nash equilibria on (un)stable networks*. Work. Pap., Fed. Reserve Board, Washington, DC
- Ballester C, Calvó-Armengol A, Zenou Y. 2006. Who's who in networks. Wanted: the key player. *Econometrica* 74:1403–17
- Banerjee A, Chandrasekhar A, Duflo E, Jackson M. 2014. *Gossip: identifying central individuals in a social network*. Work. Pap., Mass. Inst. Technol., Cambridge
- Barabási A, Albert R. 1999. Emergence of scaling in random networks. *Science* 286:509–12
- Besag J. 1975. Statistical analysis of non-lattice data. *Statistician* 24:179–95
- Bhamidi S, Bresler G, Sly A. 2011. Mixing time of exponential random graphs. *Ann. Appl. Probab.* 21:2146–70
- Bhattacharya S, Bickel PJ. 2015. Subsampling bootstrap of count features of networks. *Ann. Stat.* 43:2384–411
- Bickel PJ, Chen A, Levina E. 2011. The method of moments and degree distributions for network models. *Ann. Stat.* 39:2280–301

<sup>28</sup>This paper was previously circulated with the title “Do interventions change the network? A dynamic peer effect model accounting for network changes.”

- Bloch F, Jackson MO. 2006. Definitions of equilibrium in network formation games. *Int. J. Game Theory* 34:305–18
- Bollobás B. 2001. *Random Graphs*. Cambridge, UK: Cambridge Univ. Press
- Bonacich P. 1972. Factoring and weighting approaches to status scores and clique identification. *J. Math. Sociol.* 2:113–20
- Bonacich P. 1987. Power and centrality: a family of measures. *Am. J. Sociol.* 92:1170–82
- Boucher V, Mourifié I. 2017. My friend far, far away: asymptotic properties of pairwise stable networks. *Econom. J.* 20:S14–46
- Bramoullé Y, Djebbari H, Fortin B. 2020. Peer effects in networks: a survey. *Annu. Rev. Econ.* 12:603–29
- Brin S, Page L. 1998. The anatomy of a large scale hypertextual Web search engine. *Comput. Netw.* 30:107–17
- Cameron AC, Miller DL. 2014. *Robust inference for dyadic data*. Work. Pap., Univ. Calif., Davis
- Candelaria LE, Ura T. 2018. *Identification and inference of network formation games with misclassified links*. Work. Pap., Univ. Warwick, Coventry, UK
- Chandrasekhar A. 2015. Econometrics of network formation. In *Oxford Handbook on the Economics of Networks*, ed. Y Bramoullé, A Galeotti, B Rogers, pp. 303–57. Oxford, UK: Oxford Univ. Press
- Chandrasekhar AG, Jackson M. 2016. *A network formation model based on subgraphs*. Work. Pap., Stanford Univ., Stanford, CA
- Charbonneau KB. 2017. Multiple fixed effects in binary response panel data models. *Econometrics J.* 20:S1–13
- Chatterjee S, Diaconis P. 2013. Estimating and understanding exponential random graph models. *Ann. Stat.* 41:2428–61
- Chatterjee S, Diaconis P, Sly A. 2011. Random graphs with a given degree sequence. *Ann. Appl. Probab.* 21:1400–35
- Chiappori PA. 2019. *Matching with Transfers: The Economics of Love and Marriage*. Princeton, NJ: Princeton Univ. Press
- Christakis NA, Fowler JH, Imbens GW, Kalyanaraman K. 2020. An empirical model for strategic network formation. In *The Econometric Analysis of Network Data*, ed. BS Graham, Á de Paula, pp. 123–48. Amsterdam: Elsevier
- Chwe M. 2000. Communication and coordination in social networks. *Rev. Econ. Stud.* 67:1–16
- Ciliberto F, Murry C, Tamer E. 2018. *Market structure and competition in airline markets*. Work. Pap., Univ. Va., Charlottesville
- Comola M, Fafchamps M. 2014. Testing unilateral and bilateral link formation. *Econ. J.* 124:954–76
- Comola M, Fafchamps M. 2017. The missing transfers: estimating mis-reporting in dyadic data. *Econ. Dev. Cult. Change* 65:549–82
- Comola M, Prina S. 2019. *Treatment effect accounting for network changes*. Work. Pap., Paris Sch. Econ., Paris
- Cragg JG. 1971. Some statistical models for limited dependent variables with application to the demand for durable goods. *Econometrica* 39:829–44
- de Paula Á. 2013. Econometric analysis of games with multiple equilibria. *Annu. Rev. Econ.* 5:107–31
- de Paula Á. 2017. Econometrics of network models. In *Advances in Economics and Econometrics: Theory and Applications*, ed. B Honore, A Pakes, M Piazzesi, L Samuelson, pp. 268–323. Cambridge, UK: Cambridge Univ. Press
- de Paula Á, Richards-Shubik S, Tamer E. 2018. Identifying preferences in networks with bounded degrees. *Econometrica* 86:263–88
- Dzanski A. 2019. An empirical model of dyadic link formation in a network with unobserved heterogeneity. *Rev. Econ. Stat.* 101:763–76
- Erdős P, Rényi A. 1959. On random graphs. *Publ. Math. Debr.* 6:290–97
- Erdős P, Rényi A. 1960. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.* 5:17–61
- Fafchamps M, Gubert F. 2007. The formation of risk sharing networks. *J. Dev. Econ.* 83:326–50
- Fafchamps M, Lund S. 2003. Risk-sharing networks in rural Philippines. *J. Dev. Econ.* 71:261–87
- Fernández-Val I, Weidner M. 2016. Individual and time effects in nonlinear panel models with large  $n$ . *J. Econom.* 192:291–312
- Fernández-Val I, Weidner M. 2018. Fixed effects estimation of large- $n$  panel data models. *Annu. Rev. Econ.* 10:109–38

- Frank O, Strauss D. 1986. Markov graphs. *J. Am. Stat. Assoc.* 81:832–42
- Gao WY. 2020. Nonparametric identification in index models of link formation. *J. Econom.* 215:399–413
- Ghili S. 2018. *Network formation and bargaining in vertical markets: the case of narrow networks in health insurance*. Work. Pap., Yale Univ., New Haven, CT
- Gilbert E. 1959. Random graphs. *Ann. Math. Stat.* 30:1141–44
- Gilleskie D, Zhang YS. 2009. *Friendship formation and smoking initiation among teens*. Work. Pap., Univ. N.C., Chapel Hill
- Gilmore M. 2016. The curse of the Ramones. *Rolling Stone*, May 19. <https://www.rollingstone.com/culture/culture-news/the-curse-of-the-ramones-165741/>
- Goldsmith-Pinkham P, Imbens G. 2013. Social networks and the identification of peer effects. *J. Bus. Econ. Stat.* 31:253–64
- Gould P. 1967. On the geographical interpretation of eigenvalues. *Trans. Inst. Br. Geogr.* 42:53–86
- Graham BS. 2015. Methods of identification in social networks. *Annu. Rev. Econ.* 7:465–85
- Graham BS. 2016. *Homophily and transitivity in dynamic network formation*. Work. Pap., Univ. Calif., Berkeley
- Graham BS. 2017. An econometric model of network formation with degree heterogeneity. *Econometrica* 85:1033–63
- Graham BS. 2020. The econometrics of networks. In *Handbook of Econometrics*, Vol. 7A, ed. S Durlauf, L Hansen, JJ Heckman, R Matzkin. Amsterdam: Elsevier. In press
- Graham BS, de Paula Á, eds. 2020. *The Econometric Analysis of Network Data*. Amsterdam: Elsevier
- Gualdani C. 2019. *An econometric model of network formation with an application to board interlocks between firms*. Work. Pap., Toulouse Sch. Econ., Toulouse, Fr.
- Han A. 1987. Non-parametric analysis of a generalized regression model: the maximum rank correlation estimator. *J. Econom.* 35:303–16
- Han X, Hsieh CS, Ko S. 2019. *Spatial modeling approach for dynamic network formation and interactions*. Work. Pap., Natl. Taiwan Univ., Taipei
- Ho K, Lee R. 2019. Equilibrium provider networks: bargaining and exclusion in health care markets. *Am. Econ. Rev.* 109:473–522
- Hoff P. 2005. Bilinear mixed-effects models for dyadic data. *J. Am. Stat. Assoc.* 100:286–95
- Holland PW, Leinhardt S. 1976. Local structure in social networks. *Sociol. Methodol.* 7:1–45
- Holland PW, Leinhardt S. 1981. An exponential family of probability distributions for directed graphs. *J. Am. Stat. Assoc.* 76:33–65
- Horn H, Wolinsky A. 1988. Bilateral monopolies and incentives for merger. *RAND J. Econ.* 19:408–19
- Hsieh CS, Lee LF. 2016. A social interactions model with endogenous friendship formation and selectivity. *J. Appl. Econom.* 31:301–19
- Hsieh CS, Lee LF, Boucher V. 2019. *Specification and estimation of network formation and network interaction models with the exponential probability distribution*. Work. Pap., Ohio State Univ., Columbus
- Jackson M. 2009. *Social and Economic Networks*. Princeton, NJ: Princeton Univ. Press
- Jackson MO, Watts A. 2002. The evolution of social and economic networks. *J. Econ. Theory* 106:265–95
- Jackson MO, Wolinsky A. 1996. A strategic model of social and economic networks. *J. Econ. Theory* 71:44–74
- Jochmans K. 2018. Semiparametric analysis of network formation. *J. Bus. Econ. Stat.* 36:705–13
- Jordan MI, Wainwright MJ. 2008. *Graphical Models, Exponential Families, and Variational Inference*. Delft, Neth.: Now Publ.
- Katz L. 1953. A new status index derived from sociometric analysis. *Psychometrika* 18:39–43
- Kolaczyk E. 2009. *Statistical Analysis of Network Data*. Berlin: Springer-Verlag
- Leung M. 2015a. *A random-field approach to inference in large models of network formation*. Work. Pap., Stanford Univ., Stanford, CA
- Leung M. 2015b. Two-step estimation of network formation models with incomplete information. *J. Econom.* 188:182–95
- Liebman E. 2018. *Bargaining in markets with exclusion: an analysis of health insurance networks*. Work. Pap., Univ. Ga., Athens
- Manski CF. 1993. Identification of endogenous social effects: the reflection problem. *Rev. Econ. Stud.* 60:531–42



- Mele A. 2017. A structural model of dense network formation. *Econometrica* 85:825–50
- Mele A. 2020. Does school desegregation promote diverse interactions? An equilibrium model of segregation within schools. *Am. Econ. J. Econ. Policy* 12:228–57
- Mele A, Zhu L. 2019. *Approximate variational estimation for a model of network formation*. Work. Pap., Johns Hopkins Univ., Baltimore, MD
- Menzel K. 2016. *Strategic network formation with many agents*. Work. Pap., New York Univ., New York
- Miyauchi Y. 2016. Structural estimation of a pairwise stable network with nonnegative externality. *J. Econom.* 195:224–35
- Murray I, Ghahramani Z, MacKay D. 2006. MCMC for doubly-intractable distributions. In *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence*, ed. R Dechter, T Richardson, pp. 359–66. Arlington, VA: AUAI Press
- Orbanz P, Roy DM. 2015. Bayesian models of graphs, arrays and other exchangeable randoms structures. *IEEE Trans. Pattern Anal. Mach. Intell.* 37:437–61
- Pelican A, Graham B. 2019. *Testing for strategic interaction in social and economic network formation*. Work. Pap., Univ. Calif., Berkeley
- Poirier D. 1980. Partial observability in bivariate probit models. *J. Econom.* 12:209–17
- Rasch G. 1960. *Probabilistic Models for Some Intelligence and Attainment Tests*. Chicago: Univ. Chicago Press
- Ridder G, Sheng S. 2015. *Estimation of large network formation games*. Work. Pap., Univ. South. Calif., Los Angeles
- Robins G, Pattison P, Kalish Y, Lusher D. 2007. An introduction to exponential random graph ( $p^*$ ) models for social networks. *Soc. Netw.* 29:173–91
- Sheng S. 2018. *A structural econometric analysis of network formation games through subnetworks*. Work. Pap., Univ. Calif., Los Angeles
- Shi Z, Chen X. 2016. *A structural network pairwise regression model with individual heterogeneity*. Work. Pap., Chin. Univ. Hong Kong, Hong Kong
- Snijders T. 2002. Markov chain Monte Carlo estimation of exponential random graph models. *J. Soc. Struct.* 3:1–40
- Strauss D, Ikeda M. 1990. Pseudolikelihood estimation for social networks. *J. Am. Stat. Assoc.* 85:204–12
- Tabord-Meehan M. 2019. Inference with dyadic data: asymptotic behavior of the dyadic-robust  $t$ -statistic. *J. Bus. Econ. Stat.* 37:671–80
- Tamer E. 2010. Partial identification in econometrics. *Annu. Rev. Econ.* 2:167–95
- Toth P. 2018. *Semiparametric estimation in network formation models with homophily and degree homophily*. Work. Pap., Univ. Tex., Austin
- Yan T, Jiang B, Fienberg SE, Leng C. 2019. Statistical inference in a directed network model with covariates. *J. Am. Stat. Assoc.* 114:857–68
- Yan T, Xu J. 2013. A central limit theorem in the  $\beta$ -model for undirected random graphs with a diverging number of vertices. *Biometrika* 100:519–24
- Zenou Y. 2016. Key players. In *The Oxford Handbook of the Economics of Networks*, ed. Y Bramoullé, A Galeotti, B Rogers, pp. 244–74. Oxford, UK: Oxford Univ. Press
- Zheng T, Salganik MJ, Gelman A. 2006. How many people do you know in prison? Using overdispersion in count data to estimate social structure in networks. *J. Am. Stat. Assoc.* 101:409–23



# Contents

|   |     |
|---|-----|
| Economics with a Moral Compass? Welfare Economics: Past, Present,<br>and Future<br><i>Amartya Sen, Angus Deaton, and Timothy Besley</i> .....     | 1   |
| Trade Policy in American Economic History<br><i>Douglas A. Irwin</i> .....  | 23  |
| An Econometric Perspective on Algorithmic Subsampling<br><i>Sokbae Lee and Serena Ng</i> .....  | 45  |
| Behavioral Implications of Causal Misperceptions<br><i>Ran Spiegler</i> .....   | 81  |
| Poverty and the Labor Market: Today and Yesterday<br><i>Robert C. Allen</i> .....   | 107 |
| The Econometrics of Static Games<br><i>Andrés Aradillas-López</i> .....   | 135 |
| On Measuring Global Poverty<br><i>Martin Ravallion</i> .....  | 167 |
| Taxation and the Superrich<br><i>Florian Scheuer and Joel Slemrod</i> .....   | 189 |
| How Distortions Alter the Impacts of International Trade in Developing<br>Countries<br><i>David Atkin and Amit K. Khandelwal</i> .....            | 213 |
| Robust Decision Theory and Econometrics<br><i>Gary Chamberlain</i> .....  | 239 |
| Cities in the Developing World<br><i>Gharad Bryan, Edward Glaeser, and Nick Tsivanidis</i> .....  | 273 |
| New Developments in Revealed Preference Theory: Decisions Under<br>Risk, Uncertainty, and Intertemporal Choice<br><i>Federico Echenique</i> ..... | 299 |
| Computing Economic Equilibria Using Projection Methods<br><i>Alena Miftakbova, Karl Schmedders, and Malte Schumacher</i> .....                    | 317 |

|   |     |
|---|-----|
| Social Identity and Economic Policy<br><i>Moses Shayo</i> .....   | 355 |
| Empirical Models of Lobbying<br><i>Matilde Bombardini and Francesco Trebbi</i> .....  | 391 |
| Political Effects of the Internet and Social Media<br><i>Ekaterina Zhuravskaya, Maria Petrova, and Ruben Enikolopov</i> ..... | 415 |
| Nash Equilibrium in Discontinuous Games<br><i>Philip J. Reny</i> .....  | 439 |
| Revealed Preference Analysis of School Choice Models<br><i>Nikhil Agarwal and Paulo Somaini</i> .....                         | 471 |
| Social Networks and Migration<br><i>Kaivan Munshi</i> .....   | 503 |
| Informality: Causes and Consequences for Development<br><i>Gabriel Ulyssea</i> .....  | 525 |
| The Theory and Empirics of the Marriage Market<br><i>Pierre-André Chiappori</i> .....   | 547 |
| Modeling Imprecision in Perception, Valuation, and Choice<br><i>Michael Woodford</i> .....                                    | 579 |
| Peer Effects in Networks: A Survey<br><i>Yann Bramoullé, Habiba Djebbari, and Bernard Fortin</i> .....                        | 603 |
| Alternative Work Arrangements<br><i>Alexandre Mas and Amanda Pallais</i> .....  | 631 |
| Shotgun Wedding: Fiscal and Monetary Policy<br><i>Marco Bassetto and Thomas J. Sargent</i> .....                              | 659 |
| Social Identity, Group Behavior, and Teams<br><i>Gary Charness and Yan Chen</i> .....   | 691 |
| Aspirations and Economic Behavior<br><i>Garance Genicot and Debraj Ray</i> .....  | 715 |
| The Search Theory of Over-the-Counter Markets<br><i>Pierre-Olivier Weill</i> .....  | 747 |
| Econometric Models of Network Formation<br><i>Áureo de Paula</i> .....  | 775 |
| Dynamic Taxation<br><i>Stefanie Stantcheva</i> .....  | 801 |
| Capital Flows and Leverage<br><i>Şebnem Kalemli-Özcan and Jun Hee Kwak</i> .....  | 833 |