

# Inferential Approaches for Network Analysis: AMEN for Latent Factor Models

Shahryar Minhas<sup>a,c,1</sup>, Peter D. Hoff<sup>b,1,2</sup>, and Michael D. Ward<sup>a</sup>

<sup>a</sup>Department of Political Science, Michigan State University, East Lansing, MI 48824, USA; <sup>b</sup>Departments of Statistics, Duke University, Durham, NC 27701, USA; <sup>c</sup>Department of Political Science, Duke University, Durham, NC 27701, USA

This manuscript was compiled on June 22, 2017

There is growing interest in the study of social networks. Network analysis allows scholars to move away from focusing on individual observations to the interrelationships among observations. Many network approaches have been developed in descriptive fashion, but attention to inferential approaches to network analysis has been growing. We introduce a new approach that models interdependencies among observations using additive and multiplicative effects (AME). This approach can be applied to binary, ordinal, and continuous network data, and provides a set of tools for inference from longitudinal networks as well. We review this approach and compare it to an alternative latent variable approach and exponential random graph models. The AME approach is shown a) to be easy to implement; b) interpretable in a general linear model framework; c) computationally straightforward; d) not prone to degeneracy; and e) notably outperforms alternatives in terms of predicting links within a network in an out of sample context and capturing network dependencies.

Social network analysis | bayesian estimation | latent variable models | inference

Network analysis provides a way to represent and study “relational data”, that is data which defines characteristics between pairs of actors. Data structures that extend beyond the actor level are common across many fields in the social sciences. In the study of international relations, for instance, the focus often rests on how countries conflict or cooperate with one another. Yet, the dominant paradigm in international relations and the social sciences more broadly for dealing with such data structures is not a network approach but rather a dyadic design, in which an interaction between a pair of countries is considered independent of interactions between any other pair in the system.

The implication of this assumption is that when, for example, Vietnam and the United States decide to form a trade agreement, they make this decision independently of what they have done with other countries and what other countries in the international system have done among themselves. A common defense of the dyad-only approach is that many events are only bilateral (1), thus alleviating the need for an approach that incorporates interdependencies between observations. The network perspective asserts that even bilateral events and processes take place within a broader system. At a minimum, we do not know whether independence of events and processes characterizes what we observe, thus we should at least examine this assertion. The potential for interdependence

among observations poses a challenge to theoretical as well as statistical modeling since the assumption made by standard approaches used across the social sciences is that observations are, at least, conditionally independent (2). The consequence of ignoring this assumption has been frequently noted already in fields such as political science (3–5), computer science (6, 7), economics (8, 9), and psychology (10, 11).

The goal of our paper is to introduce the additive and multiplicative effects model (AME) and contrast it with popular, alternative approaches for conducting statistical inference on networks, namely, the latent distance model (LDM) developed by (12) and exponential random graph models (ERGMs). To illustrate the contrasts between AME and these alternatives, we apply each to studying a cross-sectional network measuring collaborations during the policy design of the Swiss CO<sub>2</sub> act. By doing so we are able to show that AME provides a superior goodness of fit to the data than alternative approaches, both in terms of ability to predict linkages and capture network dependencies.

In general, the AME approach is a flexible framework that can be used to estimate many different types of cross-sectional and longitudinal network (e.g., binomial, gaussian, and ordinal edges) in a straightforward way. The AME modeling framework can provide a flexible and easy to use scheme through which scholars can study relational data. It addresses the issue of interdependence while still allowing scholars to examine theories that may only be relevant in the monadic or dyadic level. Further, at the network level it accounts for nodal and dyadic dependence patterns, and provides a descriptive visualization of higher-order dependencies such as homophily and stochastic equivalence.

## Addressing Dependencies in Dyadic Data

Relational, or dyadic, data provide measurements of how pairs of actors relate to one another. The easiest way to organize such data is the directed dyadic design in which the unit of analysis is some set of  $n$  actors that have been paired together to form a dataset of  $z$  directed dyads. A tabular design such as this for a set of  $n$  actors,  $\{i, j, k, l\}$  results in  $n \times (n - 1)$  observations, as shown in Table 1.

When modeling relational data, scholars typically employ a generalized linear model (GLM). An assumption we make

### Significance Statement

We introduce an additive and multiplicative (AME) effects model that can be used for conducting inference on cross-sectional and longitudinal networks. The approach we develop is distinctive in that it applies to data that include binary links as well as relations between actors that may be ordinal or Gaussian. Unlike popular, alternative approaches, AME can be easily a) interpreted as a regression, b) captures 1st, 2nd, and 3rd order linkages, c) is easy to compute, and d) out-performs alternative inferential approaches in terms of both predicting linkages and capturing network dependencies.

S.M., P.H., and M.W. designed research, performed research, contributed analytic tools, wrote the paper, and S.M. analyzed data.

The authors declare no conflict of interests.

<sup>2</sup>To whom correspondence should be addressed. E-mail: s7.minhas@gmail.com

**Table 1. Structure of datasets used in canonical design.**

Sender	Receiver	Event
$i$	$j$	$y_{ij}$
$\vdots$	$k$	$y_{ik}$
$\vdots$	$l$	$y_{il}$
$j$	$i$	$y_{ji}$
$\vdots$	$k$	$y_{jk}$
$\vdots$	$l$	$y_{jl}$
$k$	$i$	$y_{ki}$
$\vdots$	$j$	$y_{kj}$
$\vdots$	$l$	$y_{kl}$
$l$	$i$	$y_{li}$
$\vdots$	$j$	$y_{lj}$
$\vdots$	$k$	$y_{lk}$

**Table 2. Adjacency matrix representation of data in Table 1. Senders are represented by the rows and receivers by the columns.**

	$i$	$j$	$k$	$l$
$i$	NA	$y_{ij}$	$y_{ik}$	$y_{il}$
$j$	$y_{ji}$	NA	$y_{jk}$	$y_{jl}$
$k$	$y_{ki}$	$y_{kj}$	NA	$y_{kl}$
$l$	$y_{li}$	$y_{lj}$	$y_{lk}$	NA

when applying this modeling technique is that each of the dyadic observations is conditionally independent. However, this is a strong assumption to make given that events between actors in a network are often interdependent. The dependencies that tend to develop in relational data can be more easily understood when we move away from stacking dyads on top of one another and turn instead to a matrix design as shown in Table 2. Operationally, this type of data structure is represented as a  $n \times n$  matrix,  $\mathbf{Y}$ , where the diagonals are typically undefined. The  $y_{ij}^{th}$  entry defines the relationship sent from  $i$  to  $j$  and can be continuous or discrete.

The most common type of dependency that arises in networks are first-order, or nodal dependencies, and these point to the fact that we typically find significant heterogeneity in activity levels across nodes. The implication of this across-node heterogeneity is within-node homogeneity of ties, meaning that values across a row, say  $\{y_{ij}, y_{ik}, y_{il}\}$ , will be more similar to each other than other values in the adjacency matrix because each of these values has a common sender  $i$ . This type of dependency manifests in cases where sender  $i$  tends to be more active or less active in the network than other senders. Similarly, while some actors may be more active in sending ties to others in the network, we might also observe that others are more popular targets, this would manifest in observations down a column,  $\{y_{ji}, y_{ki}, y_{li}\}$ , being more similar. Last, we might also find that actors who are more likely to send ties in a network are also more likely to receive them, meaning that the row and column means of an adjacency matrix may be correlated. Another ubiquitous type of structural interdependency is reciprocity. This is a second-order, or dyadic, dependency relevant only to directed datasets, and asserts that values of  $y_{ij}$  and  $y_{ji}$  may be statistically dependent. The prevalence of these types of potential interactions within directed dyadic data also complicates the basic assumption of observational independence. In the SI Appendix, we include a lengthier

discussion on the limitations of studying relational data via the typical GLM framework.

The presence of these types of interdependencies in relational data complicates the *a priori* assumption of observational independence. Accordingly, inferences drawn from misspecified models that ignore potential interdependencies between dyadic observations are likely to have a number of issues including biased estimates of the effect of independent variables, uncalibrated confidence intervals, and poor predictive performance. By ignoring these interdependencies, we ignore a potentially important part of the data generating process behind relational data.

## Social Relations Model: Additive Part of AME

The relevance of modeling first- and second-order dependencies has long been recognized within fields such as psychology. Warner et al. developed the social relations model (SRM), a type of ANOVA decomposition technique, that facilitates this undertaking (13). The SRM is of particular note as it provides the error structure for the additive effects component of the AME framework that we introduce here. The goal of the SRM is to decompose the variance of observations in an adjacency matrix in terms of heterogeneity across row means (out-degree), heterogeneity along column means (in-degree), correlation between row and column means, and correlations within dyads. Wong and Li & Loken and provide a random effects representation of the SRM (14, 15):

$$\begin{aligned}
 y_{ij} &= \mu + e_{ij} \\
 e_{ij} &= a_i + b_j + \epsilon_{ij} \\
 \{(a_1, b_1), \dots, (a_n, b_n)\} &\stackrel{\text{iid}}{\sim} N(0, \Sigma_{ab}) \\
 \{(\epsilon_{ij}, \epsilon_{ji}) : i \neq j\} &\stackrel{\text{iid}}{\sim} N(0, \Sigma_{\epsilon}), \text{ where} \\
 \Sigma_{ab} &= \begin{pmatrix} \sigma_a^2 & \sigma_{ab} \\ \sigma_{ab} & \sigma_b^2 \end{pmatrix} \quad \Sigma_{\epsilon} = \sigma_{\epsilon}^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}
 \end{aligned} \tag{1}$$

In the above,  $\mu$  provides a baseline measure of the density or sparsity of a network, and  $e_{ij}$  represents residual variation. The residual variation decomposes into parts: a row/sender effect ( $a_i$ ), a column/receiver effect ( $b_j$ ), and a within-dyad effect ( $\epsilon_{ij}$ ). The row and column effects are modeled jointly to account for correlation in how active an actor is in sending and receiving ties. Heterogeneity in the row and column means is captured by  $\sigma_a^2$  and  $\sigma_b^2$ , respectively, and  $\sigma_{ab}$  describes the linear relationship between these two effects (i.e., whether actors who send [receive] a lot of ties also receive [send] a lot of ties). Beyond these first-order dependencies, second-order dependencies are described by  $\sigma_{\epsilon}^2$  and a within dyad correlation, or reciprocity, parameter  $\rho$ .

The SRM covariance structure described in Equation 1 can be incorporated into the systematic component of a GLM framework to produce the social relations regression model (SRRM):  $\beta^T \mathbf{X}_{ij} + a_i + b_j + \epsilon_{ij}$ , where  $\beta^T \mathbf{X}_{ij}$  accommodates the inclusion of dyadic, sender, and receiver covariates (16). The SRRM approach incorporates row, column, and within-dyad dependence in way that is widely used and understood by applied researchers: a regression framework and additive random effects to accommodate variances and covariances often seen in relational data. Furthermore, this handles a diversity of outcome distributions. In the case of binary data

this can be done by utilizing a latent variable representation of a probit regression model.

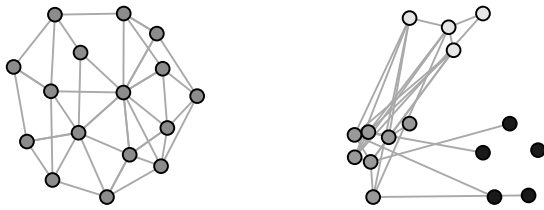
### Latent Factor Model: Multiplicative Part of AME

Missing from the framework provided by the SRM is an accounting of third-order dependence patterns that can arise in relational data. The ubiquity of third-order effects in relational datasets arises from the presence of some set of shared attributes between nodes that affects their probability of interacting with one another. Another reason why we may see the emergence of third-order effects is the “sociology” explanation: that individuals want to close triads because this is putatively a more stable or preferable social situation (17).

For example, one finding from the gravity model of trade is that neighboring countries are more likely to trade with one another; in this case, the shared attribute is simply geographic proximity. A finding common in the political economy literature is that democracies are more likely to form trade agreements with one another, and the shared attribute here is a country’s political system. Both geographic proximity and a country’s political system are examples of homophily, which captures the idea that the relationships between actors with similar characteristics in a network are likely to be stronger than nodes with different characteristics. Homophily can be used to explain the emergence of patterns such as transitivity (“a friend of a friend is a friend”) and balance (“an enemy of a friend is an enemy”).

A binary network where actors tend to form ties with others based on some set of shared characteristics often leads to a network graph with a high number of “transitive triads” in which sets of actors  $\{i, j, k\}$  are each linked to one another. The left-most plot in Figure 1 provides a representation of a network that exhibits this type of pattern. The relevant implication of this when it comes to conducting statistical inference is that—unless we are able to specify the list of exogenous variable that may explain this prevalence of triads—the probability of  $j$  and  $k$  forming a tie is not independent of the ties that already exist between those actors and  $i$ .

**Fig. 1.** Graph on the left is a representation of an undirected network that exhibits a high degree of homophily, while on the right we show an undirected network that exhibits stochastic equivalence.



Another third-order dependence pattern that cannot be accounted for in the additive effects framework is stochastic equivalence. A pair of actors  $ij$  are stochastically equivalent if the probability of  $i$  relating to, and being related to, by every other actor is the same as the probability for  $j$ . This refers to the idea that there will be groups of nodes in a network with similar relational patterns. The occurrence of a dependence pattern such as this is not uncommon in the social science applications. Recent work estimates a stochastic equivalence structure to explain the formation of preferential

trade agreements (PTAs) between countries (18). Specifically, they suggest that PTA formation is related to differences in per capita income levels between countries. Countries falling into high, middle, and low income per capita levels will have patterns of PTA formation that are determined by the groups into which they fall. Such a structure is represented in the right-most panel of Figure 1, here the lightly shaded group of nodes at the top can represent high-income countries, nodes on the bottom-left middle-income, and the darkest shade of nodes low-income countries. The behavior of actors in a network can at times be governed by group level dynamics, and failing to account for such dynamics leaves potentially important parts of the data generating process ignored.

To account for third-order dependence patterns within the context of the SRRM we turn to latent variable models, which have become a popular approach for modeling relational data in fields as diverse as biology to computer science to the social sciences. These models assumes that relationships between nodes are mediated by a small number ( $K$ ) of node-specific unobserved latent variables. One reason for their increased usage is that they enable researchers to capture and visualize third-order dependencies in a way that other approaches are not able to replicate. Additionally, the conditional independence assumption facilitates the testing of a variety of nodal and dyadic level theories, and provides a range of computational advantages relative to ERGMs.

A number of latent variable approaches have been developed to represent third-order dependencies in relational data, we focus on two here: the latent distance model and the latent factor model (LFM). For the sake of exposition, we consider the case where relations are symmetric to describe the differences between these approaches. Both of these approaches can be incorporated into an undirected version of the framework that we have been constructing through the inclusion of an additional term to the model for  $y_{ij}$ ,  $\alpha(u_i, u_j)$ , that captures latent third-order characteristics of a network, where  $u_i$  and  $u_j$  are node-specific latent variables. General definitions for how  $\alpha(u_i, u_j)$  is defined for these latent variable models are shown in Equations 2. One other point of note about these approaches is that researchers have to specify a value for  $K$ .

Latent distance model

$$\alpha(\mathbf{u}_i, \mathbf{u}_j) = -|\mathbf{u}_i - \mathbf{u}_j|$$

$$\mathbf{u}_i \in \mathbb{R}^K, i \in \{1, \dots, n\}$$

Latent factor model

$$\alpha(\mathbf{u}_i, \mathbf{u}_j) = \mathbf{u}_i^\top \Lambda \mathbf{u}_j$$

$$\mathbf{u}_i \in \mathbb{R}^K, i \in \{1, \dots, n\}$$

$$\Lambda \text{ a } K \times K \text{ diagonal matrix}$$

In the LDM approach, each node  $i$  has some unknown latent position in  $K$  dimensional space,  $\mathbf{u}_i \in \mathbb{R}^K$ , and the probability of a tie between a pair  $ij$  is a function of the negative Euclidean distance between them:  $-|\mathbf{u}_i - \mathbf{u}_j|$ . Hoff et al. show that because latent distances for a triple of actors obey the triangle inequality, this formulation models the tendencies toward homophily commonly found in social networks. This approach has been operationalized in the **latentnet** R package developed by Krivitsky & Handcock (19). However, this approach also comes with an important shortcoming: it confounds stochastic



equivalence and homophily. Consider two nodes  $i$  and  $j$  that are proximate to one another in  $K$  dimensional Euclidean space, this suggests not only that  $\|\mathbf{u}_i - \mathbf{u}_j\|$  is small but also that  $\|\mathbf{u}_i - \mathbf{u}_l\| \approx \|\mathbf{u}_j - \mathbf{u}_l\|$ , the result being that nodes  $i$  and  $j$  will by construction assumed to possess the same relational patterns with other actors such as  $l$  (i.e., that they are stochastically equivalent). Thus LDMs confound strong ties with stochastic equivalence. This approach cannot adequately model data with many ties between nodes that have different network roles.

An early iteration of the latent factor approach was presented in (16). The revised approach is motivated by an eigenvalue decomposition of a network. An important difference in the earlier approaches compared to the model that we present here is that  $\Lambda$  was taken to be the identity matrix thus stochastic equivalence could not be characterized. The motivation for this alternative framework stems from the fact that many real networks exhibit varying degrees of stochastic equivalence and homophily. In these situations, using the LDM would end up representing only a part of the network structure. In the latent factor model, each actor has an unobserved vector of characteristics,  $\mathbf{u}_i = \{u_{i,1}, \dots, u_{i,K}\}$ , which describe their behavior as an actor in the network. The probability of a tie from  $i$  to  $j$  depends on the extent to which  $\mathbf{u}_i$  and  $\mathbf{u}_j$  are “similar” (i.e., point in the same direction) and on whether the entries of  $\Lambda$  are greater or less than zero.

More specifically, the similarity in the latent factors,  $\mathbf{u}_i \approx \mathbf{u}_j$ , corresponds to how stochastically equivalent a pair of actors are and the eigenvalue determines whether the network exhibits positive or negative homophily. For example, say that that we estimate a rank-one latent factor model (i.e.,  $K = 1$ ), in this case  $\mathbf{u}_i$  is represented by a scalar  $u_{i,1}$ , similarly,  $\mathbf{u}_j = u_{j,1}$ , and  $\Lambda$  will have just one diagonal element  $\lambda$ . The average effect this will have on  $y_{ij}$  is simply  $\lambda \times u_i \times u_j$ , where a positive value of  $\lambda > 0$  indicates homophily and  $\lambda < 0$  anti-homophily. This approach can represent both homophily and stochastic equivalence, and that the alternative latent variable approaches can be represented as a latent factor model but not vice versa (20). In the directed version of this approach, we use the singular value decomposition, here actors in the network have a vector of latent characteristics to describe their behavior as a sender, denoted by  $\mathbf{u}$ , and as a receiver,  $\mathbf{v}$ :  $\mathbf{u}_i, \mathbf{v}_j \in \mathbb{R}^K$ . This can alter the probability of an interaction between  $ij$  additively:  $\mathbf{u}_i^\top \mathbf{D} \mathbf{v}_j$ , where  $\mathbf{D}$  is a  $K \times K$  diagonal matrix.

Both the latent distance and factor models are “conditional independence models” in that they assume that ties are conditionally independent given all of the observed predictors and unknown node-specific parameters:  $p(Y|X, U) = \prod_{i < j} p(y_{i,j} | x_{i,j}, u_i, u_j)$ . Typical parametric models of this form relate  $y_{i,j}$  to  $(x_{i,j}, u_i, u_j)$  via some sort of link function:

$$p(y_{i,j} | x_{i,j}, u_i, u_j) = f(y_{i,j} : \eta_{i,j})$$

$$\eta_{i,j} = \beta^\top x_{i,j} + \alpha(\mathbf{u}_i, \mathbf{u}_j).$$

The structure of  $\alpha(\mathbf{u}_i, \mathbf{u}_j)$  can result in very different interpretations for any estimates of the regression coefficients  $\beta$ . For example, suppose the latent effects  $\{u_1, \dots, u_n\}$  are near zero on average (if not, their mean can be absorbed into an intercept parameter and row and column additive effects).

Under the LFM, the average value of  $\alpha(\mathbf{u}_i, \mathbf{u}_j) = \mathbf{u}_i^\top \Lambda \mathbf{u}_j$  will be near zero and so we have

$$\eta_{i,j} = \beta^\top x_{i,j} + \mathbf{u}_i^\top \Lambda \mathbf{u}_j$$

$$\bar{\eta} \approx \beta^\top \bar{x},$$

and so the values of  $\beta$  can be interpreted as yielding the “average” value of  $\eta_{i,j}$ . On the other hand, under the LDM

$$\eta_{i,j} = \beta^\top x_{i,j} - \|\mathbf{u}_i - \mathbf{u}_j\|$$

$$\bar{\eta} \approx \beta^\top \bar{x} - \overline{\|\mathbf{u}_i - \mathbf{u}_j\|} < \beta^\top \bar{x}.$$

In this case,  $\beta^\top \bar{x}$  does not represent an “average” value of the predictor  $\eta_{i,j}$ , it represents a maximal value as if all actors were zero distance from each other in the latent social space. For example, consider the simplest case of a normally distributed network outcome with an identity link. In this case,

$$y_{i,j} = \beta^\top x_{i,j} + \alpha(\mathbf{u}_i, \mathbf{u}_j) + \epsilon_{i,j}$$

$$\bar{y} \approx \beta^\top \bar{x} + \overline{\alpha(\mathbf{u}_i, \mathbf{u}_j)}$$

$$\approx \beta^\top \bar{x}.$$

Under the LDM,  $\bar{y} \approx \beta^\top \bar{x} + \overline{\|\mathbf{u}_i - \mathbf{u}_j\|} < \beta^\top \bar{x}$ , and so we no longer can interpret  $\beta$  as representing the linear relationship between  $y$  and  $x$ . Instead, it may be thought of as describing some sort of average hypothetical “maximal” relationship between  $y_{i,j}$  and  $x_{i,j}$ .

Thus the LFM provides two important benefits. First, we are able to capture a wider assortment of dependence patterns that arise in relational data, and, second, parameter interpretation is more straightforward. The AME approach considers the regression model shown in Equation 3:

$$y_{ij} = g(\theta_{ij})$$

$$\theta_{ij} = \beta^\top \mathbf{X}_{ij} + e_{ij}$$

$$e_{ij} = a_i + b_j + \epsilon_{ij} + \alpha(\mathbf{u}_i, \mathbf{v}_j), \text{ where}$$

$$\alpha(\mathbf{u}_i, \mathbf{v}_j) = \mathbf{u}_i^\top \mathbf{D} \mathbf{v}_j = \sum_{k \in K} d_k u_{ik} v_{jk} \quad [3]$$

Using this framework, we are able to model the dyadic observations as conditionally independent given  $\theta$ , where  $\theta$  depends on the unobserved random effects,  $\mathbf{e}$ .  $\mathbf{e}$  is then modeled to account for the potential first, second, and third-order dependencies that we have discussed. As described in Equation 1,  $a_i + b_j + \epsilon_{ij}$ , are the additive random effects in this framework and account for sender, receiver, and within-dyad dependence. The multiplicative effects,  $\mathbf{u}_i^\top \mathbf{D} \mathbf{v}_j$ , are used to capture higher-order dependence patterns that are left over in  $\theta$  after accounting for any known covariate information. A Bayesian procedure in which parameters are iteratively updated using a Gibbs sampler is available in the **amen** R package to estimate this type of generalized linear mixed effects model from Gaussian, binary, ordinal, and other relational data types.

## 1. ERGMs

A popular alternative approach to accounting for network dependence patterns is ERGM. ERGM approaches are useful when researchers are interested in the role that a specific list of network statistics have in giving rise to a certain network. These network statistics could include the number of transitive triads in a network, balanced triads, reciprocal pairs and so on. In the ERGM framework, a set of statistics,  $S(\mathbf{Y})$ , define a model. Given the chosen set of statistics, the probability of observing a particular network dataset  $\mathbf{Y}$  can be expressed as:

$$\Pr(Y = y) = \frac{\exp(\beta^\top S(y))}{\sum_{z \in \mathcal{Y}} \exp(\beta^\top S(z))}, y \in \mathcal{Y} \quad [4]$$

$\beta$  represents a vector of model coefficients for the specified network statistics,  $\mathcal{Y}$  denotes the set of all obtainable networks, and the denominator is used as a normalizing factor (21). This approach provides a way to state that the probability of observing a given network depends on the patterns that it exhibits, which are operationalized in the list of network statistics specified by the researcher. Within this approach one can test the role that a variety of network statistics play in giving rise to a particular network.

One issue that arises when conducting statistical inference with this model is in the calculation of the normalizing factor, which is what ensures that the expression above corresponds to a legitimate probability distribution. For even a trivially sized directed network that has only 20 actors, calculating the denominator means summing over  $2^{20 \times (20-1)} = 2^{380}$  possible networks, or, to put it another way, more than the total number of atoms in the universe. One of the first approaches to deal with this issue was a computationally fast pseudo-likelihood approach developed by Strauss & Ikeda (22). However, this approach ignores the interdependent nature of observations in relational data, as a result, many have argued that the standard errors remain unreliable (23). Additionally, there is no asymptotic theory underlying this approach on which to base the construction of confidence intervals and hypothesis tests (24). The pseudo-likelihood approach has become increasingly unpopular in recent years among those in the network analysis community, particularly, as simulation based techniques have developed—though it has not disappeared. One favored approach in the literature is to approximate the MLE using Markov Chain Monte Carlo techniques, also referred to as MCMC-MLE (25–27).

The MCMC-MLE approach is an advancement but notable problems remain. Unlike latent variable models, Rastelli et al. (28) argue that ERGMs are not even able to represent transitivity asymptotically as their clustering coefficient converges to zero. Chatterjee & Diaconis (29) have shown that MCMC procedures can take an exponential time to converge for broad classes of ERGMs unless the dyadic observations are independent. This is a result of the fact that MCMC procedures visit an infinitesimally small portion of the set of possible graphs. A related issue when estimating ERGMs is that the estimated model can become degenerate even if the observed graph is not degenerate. This means that the model is placing a large amount of probability on a small subset of networks that fall in the set of obtainable networks,  $\mathcal{Y}$ , but share little resemblance with the observed network (30). For example, most of the probability may be placed on empty graphs, no edges between

nodes, or nearly complete graphs, almost every node is connected by an edge. Some have argued that model degeneracy is simply a result of model misspecification (31, 32). This points to an important caveat in interpreting the implications of an often cited basis for ERGM, the Hammersley-Clifford theorem. Though this theorem ensures that any network can be represented through an ERGM, it says nothing about the complexity of the sufficient statistics ( $S(y)$ ) required to do so. Failure to properly account for higher-order dependence structures through an appropriate specification can at best lead to model degeneracy, which provides an obvious indication that the specification needs to be altered, and at worst deliver a result that converges but does not appropriately capture the interdependencies in the network. The consequence of the latter case is a set of inferences that will continue to be biased as a result of unmeasured heterogeneity, thus defeating the major motivation for pursuing an inferential network model in the first place.

In the following section we undertake a comparison of the latent distance model, ERGM, and the AME model. In doing so, we are able to compare and contrast these various approaches.

## Empirical Comparison

To contrast AME with these alternative approaches, we utilize a cross-sectional network measuring whether an actor indicated that they collaborated with each other during the policy design of the Swiss CO<sub>2</sub> act (33). This is a directed relational matrix as an actor  $i$  can indicate that they collaborated with  $j$  but  $j$  may not have stated that they collaborated with  $i$ . The Swiss government proposed this act in 1995 with the goal of undertaking a 10% reduction in CO<sub>2</sub> emissions by 2012. The act was accepted in the Swiss Parliament in 2000 and implemented in 2008. Ingold (33), and subsequent work by Ingold & Fischer (34), sought to determine what drives collaboration among actors trying to affect climate change policy. The set of actors included in this network are those that were identified by experts as holding an important position in Swiss climate policy. In total, Ingold identifies 34 relevant actors: five state actors, eleven industry and business representatives, seven environmental NGOs and civil society organizations, five political parties, and six scientific institutions and consultants. We follow Ingold & Fischer and Cranmer et al. (35) in developing a model specification to understand and predict link formation in this network. We do not review the specification in detail here, instead we just provide a summary of the variables to be included and the theoretical expectations of their effects in the SI Appendix.

The LDM we fit on this network includes a two-dimensional Euclidean distance metric. The ERGM specification for this network includes the same exogenous variables as LDM, but also includes a number of endogenous variables. The AME model we fit here includes the same exogenous covariates and accounts for nodal and dyadic heterogeneity using the SRM. Third-order effects are represented by the latent factor model with  $K = 2$ . Parameter estimates for these three approaches are shown in Table 3.

The first point to note is that, in general, the parameter estimates returned by the AME are similar to those of ERGM but quite different from the LDM. For example, while the LDM returns a result for the `Opposition/alliance` variable that di-

verges from ERGM, the AME returns a result that is not only similar to those approaches but in line with the theoretical expectations of Ingold & Fischer. Similar discrepancies between LDM and other approaches appear for parameters such as **Influence attribution** and **Alter's influence degree**. Each of these discrepancies are resolved when using AME. In part, this is a result of how the LDM approach complicates the interpretation of the effect of exogenous variables. In the SI Appendix, we show that these differences persist even when incorporating sender and receiver random effects.

**Table 3. \*  $p < 0.05$ . ERGM results are shown with standard errors in parentheses. LDM and AME are shown with 95% posterior credible intervals provided in brackets.**

	LDM	ERGM	AME
Intercept/Edges	0.94* [0.09; 1.82]	-12.17* (1.40)	-3.39* [-4.38; -2.50]
<b>Conflicting policy preferences</b>			
Business vs. NGO	-1.37* [-2.42; -0.41]	-1.11* (0.51)	-1.37* [-2.44; -0.47]
Opposition/alliance	0.00 [-0.40; 0.39]	1.22* (0.20)	1.08* [0.72; 1.47]
Preference dissimilarity	-1.76* [-2.62; -0.90]	-0.44 (0.39)	-0.79* [-1.55; -0.08]
<b>Transaction costs</b>			
Joint forum participation	1.51* [0.86; 2.17]	0.90* (0.28)	0.92* [0.40; 1.47]
<b>Influence</b>			
Influence attribution	0.08 [-0.40; 0.55]	1.00* (0.21)	1.09* [0.69; 1.53]
Alter's influence indegree	0.01 [-0.03; 0.04]	0.21* (0.04)	0.11* [0.07; 0.15]
Influence absolute diff.	0.04 [-0.01; 0.09]	-0.05* (0.01)	-0.07* [-0.11; -0.03]
Alter = Government actor	-0.46 [-1.08; 0.14]	1.04* (0.34)	0.55 [-0.07; 1.15]
<b>Functional requirements</b>			
Ego = Environmental NGO	-0.60 [-1.32; 0.09]	0.79* (0.17)	0.67 [-0.38; 1.71]
Same actor type	1.17* [0.63; 1.71]	0.99* (0.23)	1.04* [0.63; 1.50]
<b>Endogenous dependencies</b>			
Mutuality		0.81* (0.25)	0.39 [-0.12; 0.96]
Outdegree popularity		0.95* (0.09)	
Twopaths		-0.04* (0.02)	
GWdegree (2.0)		3.42* (1.47)	
GWESP (1.0)		0.58* (0.16)	
GWdegree (0.5)		8.42* (2.11)	

There are also a few differences between the parameter estimates that result from the ERGM and AME. Using the AME we find evidence that **Preference dissimilarity** is associated with a reduced probability of collaboration between a pair of actors, which is in line with the theoretical expectations of Ingold & Fischer. Additionally, the AME results differ from ERGM for the nodal effects related to whether a receiver of a collaboration is a government actor, **Alter=Government actor**, and whether the sender is an environmental NGO,

Ego=Environmental NGO.

**Tie Formation Prediction.** Next, we utilize a cross-validation procedure to assess the out-of-sample performance for each of the models presented in Table 3 as follows:

- Randomly divide the  $n \times (n - 1)$  data points into  $S$  sets of roughly equal size, letting  $s_{ij}$  be the set to which pair  $\{ij\}$  is assigned.
- For each  $s \in \{1, \dots, S\}$ :
  - Obtain estimates of the model parameters conditional on  $\{y_{ij} : s_{ij} \neq s\}$ , the data on pairs not in set  $s$ .
  - For pairs  $\{kl\}$  in set  $s$ , let  $\hat{y}_{kl} = E[y_{kl} | \{y_{ij} : s_{ij} \neq s\}]$ , the predicted value of  $y_{kl}$  obtained using data not in set  $s$ .

The procedure summarized in the steps above generates a sociomatrix of out-of-sample predictions of the observed data. Each entry  $\hat{y}_{ij}$  is a predicted value obtained from using a subset of the data that does not include  $y_{ij}$ . In this application we set  $S$  to 45 which corresponds to randomly excluding approximately 2% of the data from the estimation. Such a low number of observations were excluded in every fold because excluding any more observations would cause the ERGM specification to result in a degenerate model that empirically can not be fit. This highlights the computational difficulties associated with ERGMs in the presence of even small levels of missingness.

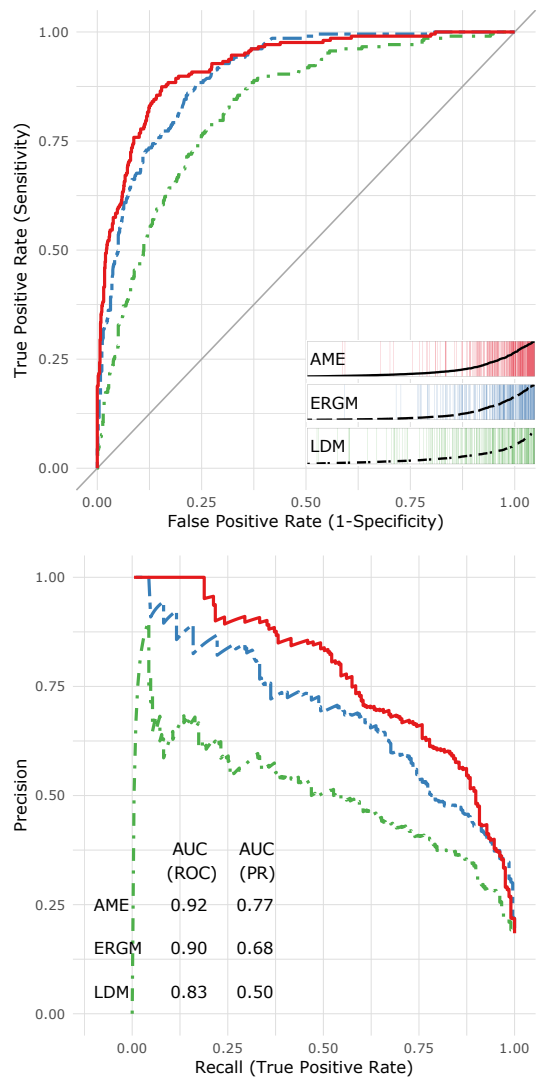
Using the set of out-of-sample predictions we generate from the cross-validation procedure, we provide a series of tests to assess model fit. The left-most plot in Figure 2 compares the five approaches in terms of their ability to predict the out-of-sample occurrence of collaboration based on Receiver Operating Characteristic (ROC) curves. ROC curves provide a comparison of the trade-off between the True Positive Rate (TPR), sensitivity, False Positive Rate (FPR), 1-specificity, for each model. Models that have a better fit according to this test should have curves that follow the left-hand border and then the top border of the ROC space. On this diagnostic, the AME model performs best closely followed by ERGM. The LDM approach lags notably behind the other specifications.

A more intuitive visualization of the differences between these modeling approaches can be gleaned through examining the separation plots included on the right-bottom edge of the ROC plot. This visualization tool plots each of the observations, in this case actor pairs, in the dataset according to their predicted value from left (low values) to right (high values). Models with a good fit should have all network links, here these are colored by the modeling approach, towards the right of the plot. Using this type of visualization we can again see that the AME and ERGM models performs better than the alternatives.

The last diagnostic we highlight to assess predictive performance are precision-recall (PR) curves. In both ROC and PR space we utilize the TPR, also referred to as recall—though in the former it is plotted on the y-axis and the latter the x-axis. The difference, however, is that in ROC space we utilize the FPR, while in PR space we use precision. FPR measures the fraction of negative examples that are misclassified as positive, while precision measures the fraction of examples classified

as positive that are truly positive. PR curves are useful in situations where correctly predicting events is more interesting than simply predicting non-events (36). This is especially relevant in the context of studying many relational datasets in political science such as conflict, because events in such data are extremely sparse and it is relatively easy to correctly predict non-events. In the case of our application dataset, the vast majority of dyads, 80%, do not have a network linkage, which points to the relevance of assessing performance using the PR curves as we do in the right-most plot of Figure 2. We can see that the relative-ordering of the models remains similar but the differences in how well they perform become much more stark. Here we find that the AME approach performs notably better in actually predicting network linkages than each of the alternatives. Area under the curve (AUC) statistics are provided in Figure 2 and these also highlight AME's superior out-of-sample performance.

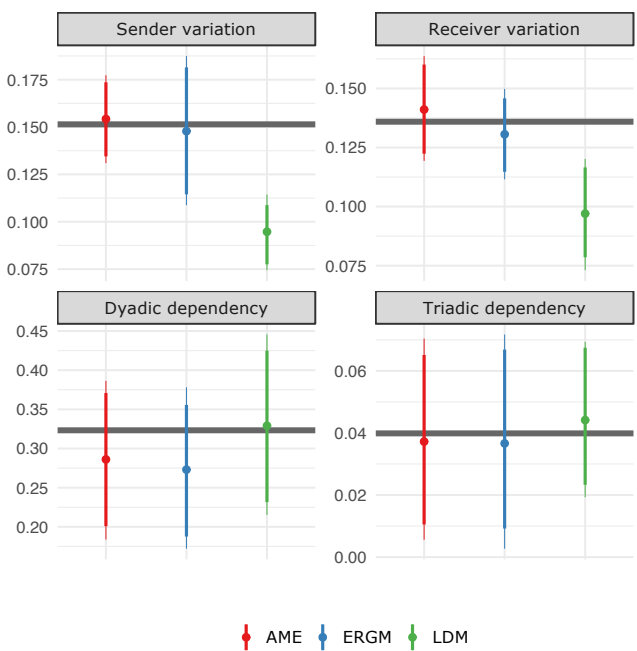
**Fig. 2.** Assessments of out-of-sample predictive performance using ROC curves, separation plots, and PR curves. AUC statistics are provided as well for both curves.



**Capturing Network Attributes.** We also assess which of these models best captures the network features of the dependent

variable. To do this, we compare the observed network with a set of networks simulated from the estimated models. We simulate 1,000 networks from the three models and compare how well they align with the observed network in terms of four network statistics: (1) the empirical standard deviation of the row means (i.e., heterogeneity of nodes in terms of the ties they send); (2) the empirical standard deviation of the column means (i.e., heterogeneity of nodes in terms of the ties they receive); (3) the empirical within-dyad correlation (i.e., measure of reciprocity in the network); and (4) a normalized measure of triadic dependence. A comparison of the LDM, ERGM, and AME models among these four statistics is shown in Figure 3. In the SI Appendix, we compare the ability of these models to capture network attribute across a wider array of statistics, and the results are consistent with what we present below.

**Fig. 3.** Network goodness of fit summary using *amen*.



Here it becomes quickly apparent that the LDM model fails to capture how active and popular actors are in the Swiss climate change mitigation network. Further even after incorporating random sender and receiver effects into the LDM framework this problem is not completely resolved, see the SI Appendix for details. The AME and ERGM specifications again both tend to do equally well. If when running this diagnostic, we found that the AME model did not adequately represent the observed network this would indicate that we might want to increase  $K$  to better account for network interdependencies. No changes to the model specification as described by the exogenous covariates a researcher has chosen would be necessary. If the ERGM results did not align with the diagnostic presented in Figure 3, then this would indicate that an incorrect set of endogenous dependencies have been specified. Failing to identify (or find) the right specification will leave the researcher with the problems we introduced earlier.



## Conclusion

The AME approach to estimation and inference in network data provides a number of benefits over alternative approaches. Specifically, it provides a modeling framework for dyadic data that is based on familiar statistical tools such as linear regression, GLM, random effects, and factor models. We have an understanding of how each of these tools work, they are numerically more stable than ERGM approaches, and more general than alternative latent variable models. Further the estimation procedure utilized in AME avoids complicating interpretation of parameter estimates for exogenous variables. For researchers in the social sciences this is of primary interest, as many studies that employ relational data still have conceptualizations that are monadic or dyadic in nature. Additionally, through the application dataset utilized herein we show that the AME approach outperforms both ERGM and LDM in out-of-sample prediction, and also is better able to capture network dependencies than the LDM.

More broadly, relational data structures are composed of actors that are part of a system. It is unlikely that this system can be viewed simply as a collection of isolated actors or pairs of actors. The assumption that dependencies between observations occur can at the very least be examined. Failure to take into account interdependencies leads to biased parameter estimates and poor fitting models. By using standard diagnostics such as shown in Figure 3, one can easily assess whether

an assumption of independence is reasonable. We stress this point because a common misunderstanding that seems to have emerged within the social science literature relying on dyadic data is that a network based approach is only necessary if one has theoretical explanations that extend beyond the dyadic. This is not at all the case and findings that continue to employ a dyadic design may misrepresent the effects of the very variables that they are interested in. The AME approach that we have detailed here provides a statistically familiar way for scholars to account for unobserved network structures in relational data. Additionally, through this approach we can visualize these dependencies in order to better learn about the network patterns that remain in the event of interest after having accounted for observed covariates.

When compared to other network based approaches, AME is easier to specify and utilize. It is also more straightforward to interpret since it does not require interpretation of unusual features such as *three-stars* which fall outside of the normal language for discussing social science. Further, the **amen** package facilitates the modeling of longitudinal network data. In sum, excuses for continuing to treat relational data as conditionally independent are no longer valid.

**ACKNOWLEDGMENTS.** S.M. and M.W. acknowledge support from National Science Foundation (NSF) Award 1259266 and P.H. acknowledges support from NSF Award 1505136.

- Diehl PF, Wright TM (2016) A conditional defense of the dyadic approach. *International Studies Quarterly*.
- Snijders TA (2011) Statistical models for social networks. *Annual Review of Sociology* 37:131–53.
- Beck N, Katz JN, Tucker R (1998) Taking time seriously: Time-series-cross-section analysis with a binary dependent variable. *American Journal of Political Science* 42(2):1260–1288.
- Signorino C (1999) Strategic interaction and the statistical analysis of international conflict. *American Political Science Review* 92(2):279–298.
- Aronow PM, Samii C, Assenova VA (2015) Cluster-robust variance estimation for dyadic data. *Political Analysis* 23(4):564–577.
- Bonabeau E (2002) Agent-based modeling: Methods and techniques for simulating human systems. *Proceedings of the National Academy of Sciences* 99(suppl 3):7280–7287.
- Brandes U, Erlebach T (2005) *Network Analysis: Methodological Foundations*. (Springer Science & Business Media) Vol. 3418.
- Goyal S (2012) *Connections: an introduction to the economics of networks*. (Princeton University Press).
- Jackson M (2014) Networks in the understanding of economic behaviors. *The Journal of Economic Perspectives* 28(4):3–22.
- Pattison P, Wasserman S (1999) Logit models and logistic regressions for social networks. ii. multivariate relations. *British Journal of Mathematical and Statistical Psychology* 52:169–194.
- Kenny DA, Kashy DA, Cook WL (2006) *Dyadic Data Analysis*. (Guilford Press, New York).
- Hoff PD, Raftery AE, Handcock MS (2002) Latent space approaches to social network analysis. *Journal of the American Statistical Association* 97(460):1090–1098.
- Warner R, Kenny D, Stoto M (1979) A new round robin analysis of variance for social interaction data. *Journal of Personality and Social Psychology* 37:1742–1757.
- Wong GY (1982) Round robin analysis of variance via maximum likelihood. *Journal of the American Statistical Association* 77(380):714–724.
- Li H, Loken E (2002) A unified theory of statistical analysis and inference for variance component models for dyadic data. *Statistica Sinica* 12(2):519–535.
- Hoff PD (2005) Bilinear mixed-effects models for dyadic data. *Journal of the American Statistical Association* 100(4690):286–295.
- Wasserman S, Faust K (1994) *Social Network Analysis: Methods and Applications*. (Cambridge University Press, Cambridge).
- Manger MS, Pickup MA, Snijders TA (2012) A hierarchy of preferences: A longitudinal network analysis approach to PTA formation. *Journal of Conflict Resolution* 56(5):852–877.
- Krivitsky PN, Handcock MS (2015) *latentnet: Latent Position and Cluster Models for Statistical Networks* (The Statnet Project (<http://www.statnet.org>)). R package version 2.7.1.
- Hoff PD (2008) Modeling homophily and stochastic equivalence in symmetric relational data in *Advances in Neural Information Processing Systems 20*, Processing Systems 21, eds. Platt JC, Koller D, Singer Y, Roweis ST. (MIT Press, Cambridge, MA, USA), pp. 657–664.
- Hunter D, Handcock M, Butts C, Goodreau SM, Morris M (2008) *ergm: A package to fit, simulate and diagnose exponential-family models for networks*. *Journal of Statistical Software* 24(3):1–29.
- Strauss D, Ikeda M (1990) Pseudolikelihood estimation for social networks. *Journal of the American Statistical Association* 85:204–212.
- Van Duijn MA, Gile KJ, Handcock MS (2009) A framework for the comparison of maximum pseudo-likelihood and maximum likelihood estimation of exponential family random graph models. *Social Networks* 31(1):52–62.
- Kolaczyk ED (2009) *Statistical Analysis of Network Data: Methods and Models*. (Springer Verlag, Berlin).
- Geyer CJ, Thompson EA (1992) Constrained Monte Carlo maximum likelihood for dependent data, (with discussion) maximum likelihood for dependent data, (with discussion). *Journal of the Royal Statistical Society, Series B, Methodological* 54:657–699.
- Snijders TA (2002) Markov chain monte carlo estimation of exponential random, graph models. *Journal of Social Structure* 3(2):via web, zeeb.library.cmu.edu:7850/JoSS/snijders/Mcpstar.pdf.
- Handcock MS (2003) Statistical models for social networks: Inference and degeneracy in *Dynamic Social Network Modeling and Analysis*, Committee on Human Factors, Board on Behavioral, Cognitive, and Sensory Sciences, eds. Ronald B, Kathlene C, Pip P. (National Academy Press, Washington D.C.) Vol. 126, pp. 229–252.
- Rastelli R, Friel N, Raftery AE (2016) Properties of latent variable network models. *Network Science* pp. 1–26.
- Chatterjee S, Diaconis P (2013) Estimating and understanding exponential random graph models. *The Annals of Statistics* 41(5):2428–2461.
- Schweinberger M (2011) Instability, sensitivity, and degeneracy of discrete exponential families. *Journal of the American Statistical Association* 106(496):1361–1370.
- Goodreau SM, Handcock MS, Hunter, David R. and Butts CT, Morris M (2008) A statnet tutorial. *Journal of Statistical Software* 24(9):1.
- Handcock MS, Hunter DR, Butts CT, Goodreau SM, Morris M (2008) *statnet: Software tools for the representation, visualization, analysis and simulation of network data*. *Journal of Statistical Software* 24(1):1548.
- Ingold K (2008) *Les mécanismes de décision: Le cas de la politique climatique Suisse. Politikanalysen*. (Rüegger Verlag, Zürich).
- Ingold K, Fischer M (2014) Drivers of collaboration to mitigate climate change: An illustration of swiss climate policy over 15 years. *Global Environmental Change* 24:88–98.
- Cranmer SJ, Leifeld P, McClurg SD, Rolfe M (2016) Navigating the range of statistical tools for inferential network analysis. *American Journal of Political Science*.
- Davis J, Goadrich M (2006) The relationship between Precision-Recall and ROC curves in *Proceedings of the 23rd International Conference on Machine Learning*. (ACM), pp. 233–240.