Statistical Models for Social Networks with Application to HIV Epidemiology



Department of Statistics University of Washington

Joint work with

Pavel Krivitsky Martina Morris

and the

U. Washington Network Modeling Group

Supported by NIH NIDA Grant DA012831 and NICHD Grant HD041877

NIPS 2007, December 4 2007

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

• Networks are widely used to represent data on relations between interacting actors or nodes.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• Networks are widely used to represent data on relations between interacting actors or nodes.

- The study of social networks is multi-disciplinary
 - plethora of terminologies
 - varied objectives, multitude of frameworks

- Networks are widely used to represent data on relations between interacting actors or nodes.
- The study of social networks is multi-disciplinary
 - plethora of terminologies
 - varied objectives, multitude of frameworks
- Understanding the structure of social relations has been the focus of the social sciences
 - *social structure:* a system of social relations tying distinct social entities to one another

• Interest in understanding how social structure form and evolve

- Networks are widely used to represent data on relations between interacting actors or nodes.
- The study of social networks is multi-disciplinary
 - plethora of terminologies
 - varied objectives, multitude of frameworks
- Understanding the structure of social relations has been the focus of the social sciences
 - *social structure:* a system of social relations tying distinct social entities to one another
 - Interest in understanding how social structure form and evolve
- Attempt to represent the structure in social relations via networks
 - the data is conceptualized as a realization of a network model

- Networks are widely used to represent data on relations between interacting actors or nodes.
- The study of social networks is multi-disciplinary
 - plethora of terminologies
 - varied objectives, multitude of frameworks
- Understanding the structure of social relations has been the focus of the social sciences
 - *social structure:* a system of social relations tying distinct social entities to one another
 - Interest in understanding how social structure form and evolve
- Attempt to represent the structure in social relations via networks
 - the data is conceptualized as a realization of a network model

- The data are of at least three forms:
 - individual-level information on the social entities
 - relational data on pairs of entities
 - population-level data

- Networks are widely used to represent data on relations between interacting actors or nodes.
- The study of social networks is multi-disciplinary
 - plethora of terminologies
 - varied objectives, multitude of frameworks
- Understanding the structure of social relations has been the focus of the social sciences
 - *social structure:* a system of social relations tying distinct social entities to one another
 - Interest in understanding how social structure form and evolve
- Attempt to represent the structure in social relations via networks
 - the data is conceptualized as a realization of a network model

- The data are of at least three forms:
 - individual-level information on the social entities
 - relational data on pairs of entities
 - population-level data

• Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)

- Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)
- Statistical Networks Community (Frank and Strauss 1986; Snijders 1997)

- Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)
- Statistical Networks Community (Frank and Strauss 1986; Snijders 1997)
- Spatial Statistics Community (Besag 1974)

- Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)
- Statistical Networks Community (Frank and Strauss 1986; Snijders 1997)
- Spatial Statistics Community (Besag 1974)
- Statistical Exponential Family Theory (Barndorff-Nielsen 1978)

- Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)
- Statistical Networks Community (Frank and Strauss 1986; Snijders 1997)
- Spatial Statistics Community (Besag 1974)
- Statistical Exponential Family Theory (Barndorff-Nielsen 1978)
- Graphical Modeling Community (Lauritzen and Spiegelhalter 1988, ...)

- Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)
- Statistical Networks Community (Frank and Strauss 1986; Snijders 1997)
- Spatial Statistics Community (Besag 1974)
- Statistical Exponential Family Theory (Barndorff-Nielsen 1978)
- Graphical Modeling Community (Lauritzen and Spiegelhalter 1988, ...)
- Machine Learning Community (Jordan, Jensen, Xing,)

- Social networks community (Heider 1946; Frank 1972; Holland and Leinhardt 1981)
- Statistical Networks Community (Frank and Strauss 1986; Snijders 1997)
- Spatial Statistics Community (Besag 1974)
- Statistical Exponential Family Theory (Barndorff-Nielsen 1978)
- Graphical Modeling Community (Lauritzen and Spiegelhalter 1988, ...)
- Machine Learning Community (Jordan, Jensen, Xing,)
- Physics and Applied Math (Newman, Watts, ...)

Example of Social Relationships between Monks

- Expressed "liking" between 18 monks within an isolated monastery
 - \Rightarrow Sampson (1969)
 - A directed relationship aggregated over a 12 month period before the breakup of the cloister.



Example of Social Relationships between Monks

- Expressed "liking" between 18 monks within an isolated monastery
 - \Rightarrow Sampson (1969)
 - A directed relationship aggregated over a 12 month period before the breakup of the cloister.
- Sampson identified three groups plus:

(T)urks, (L)oyal Opposition, (O)utcasts and (W)averers



Examples of Friendship Relationships

▲□▶ ▲圖▶ ★園▶ ★園▶ - 園 - のへで

Examples of Friendship Relationships

- The National Longitudinal Study of Adolescent Health
 www.cpc.unc.edu/projects/addhealth
 - "Add Health" is a school-based study of the health-related behaviors of adolescents in grades 7 to 12.

- Each nominated up to 5 boys and 5 girls as their friends
- 160 schools: Smallest has 69 adolescents in grades 7-12



School Community Stratum 44 mutual friendships by Grade



School Community Stratum 44 mutual friendships by Race





• *Mutuality* of ties



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- Mutuality of ties
- Individual heterogeneity in the propensity to form ties

- Mutuality of ties
- Individual heterogeneity in the propensity to form ties
- Homophily by actor attributes
 - \Rightarrow Lazarsfeld and Merton, 1954; Freeman, 1996; McPherson et al., 2001
 - higher propensity to form ties between actors with similar attributes e.g., age, gender, geography, major, social-economic status

attributes may be observed or unobserved

- Mutuality of ties
- Individual heterogeneity in the propensity to form ties
- Homophily by actor attributes
 - \Rightarrow Lazarsfeld and Merton, 1954; Freeman, 1996; McPherson et al., 2001
 - higher propensity to form ties between actors with similar attributes e.g., age, gender, geography, major, social-economic status

- attributes may be observed or unobserved
- Transitivity of relationships
 - friends of friends have a higher propensity to be friends

- Mutuality of ties
- Individual heterogeneity in the propensity to form ties
- Homophily by actor attributes
 - \Rightarrow Lazarsfeld and Merton, 1954; Freeman, 1996; McPherson et al., 2001
 - higher propensity to form ties between actors with similar attributes e.g., age, gender, geography, major, social-economic status
 - attributes may be observed or unobserved
- Transitivity of relationships
 - friends of friends have a higher propensity to be friends
- *Balance* of relationships \Rightarrow Heider (1946)
 - people feel comfortable if they agree with others whom they like

- Mutuality of ties
- Individual heterogeneity in the propensity to form ties
- Homophily by actor attributes
 - \Rightarrow Lazarsfeld and Merton, 1954; Freeman, 1996; McPherson et al., 2001
 - higher propensity to form ties between actors with similar attributes e.g., age, gender, geography, major, social-economic status
 - attributes may be observed or unobserved
- Transitivity of relationships
 - friends of friends have a higher propensity to be friends
- *Balance* of relationships \Rightarrow Heider (1946)
 - people feel comfortable if they agree with others whom they like
- *Context* is important \Rightarrow Simmel (1908)
 - triad, not the dyad, is the fundamental social unit

The Choice of Models depends on the objectives

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Primary interest in the nature of relationships:
 - How the behavior of individuals depends on their location in the social network
 - How the qualities of the individuals influence the social structure

The Choice of Models depends on the objectives

- Primary interest in the nature of relationships:
 - How the behavior of individuals depends on their location in the social network
 - How the qualities of the individuals influence the social structure
- Secondary interest is in how network structure influences processes that develop over a network

- spread of HIV and other STDs
- diffusion of technical innovations
- spread of computer viruses

The Choice of Models depends on the objectives

- Primary interest in the nature of relationships:
 - How the behavior of individuals depends on their location in the social network
 - How the qualities of the individuals influence the social structure
- Secondary interest is in how network structure influences processes that develop over a network
 - spread of HIV and other STDs
 - diffusion of technical innovations
 - spread of computer viruses
- Tertiary interest in the effect of *interventions* on network structure and processes that develop over a network

Perspectives to keep in mind

- Network-specific versus Population-process
 - Network-specific: interest focuses only on the actual network under study
 - Population-process: the network is part of a population of networks and the latter is the focus of interest
 - the network is conceptualized as a realization of a social process

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Notation

A *social network* is defined as a set of *n* social "actors" and a social relationship between each pair of actors.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Notation

A *social network* is defined as a set of *n* social "actors" and a social relationship between each pair of actors.

$$Y_{ij} = \begin{cases} 1 & \text{relationship from actor } i \text{ to actor } j \\ 0 & \text{otherwise} \end{cases}$$

▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - のへで

Notation

A *social network* is defined as a set of *n* social "actors" and a social relationship between each pair of actors.

$$Y_{ij} = \begin{cases} 1 & \text{relationship from actor } i \text{ to actor } j \\ 0 & \text{otherwise} \end{cases}$$

call Y ≡ [Y_{ij}]_{n×n} a sociomatrix
 a N = n(n − 1) binary array

Notation

A *social network* is defined as a set of *n* social "actors" and a social relationship between each pair of actors.

$$Y_{ij} = \begin{cases} 1 & \text{relationship from actor } i \text{ to actor } j \\ 0 & \text{otherwise} \end{cases}$$

• call $Y \equiv [Y_{ij}]_{n \times n}$ a sociomatrix

• a N = n(n-1) binary array

• The basic problem of stochastic modeling is to specify a distribution for Y i.e., P(Y = y)

A Framework for Network Modeling

Let \mathcal{Y} be the sample space of Y e.g. $\{0,1\}^N$ Any model-class for the multivariate distribution of Y can be *parametrized* in the form:

$$egin{aligned} & \mathcal{P}_\eta(Y=y) = rac{\exp\{\eta \cdot g(y)\}}{\kappa(\eta,\mathcal{Y})} & y \in \mathcal{Y} \end{aligned}$$

Besag (1974), Frank and Strauss (1986)

- $\eta \in \Lambda \subset R^q$ q-vector of parameters
- g(y) q-vector of network statistics.
 ⇒ g(Y) are jointly sufficient for the model
- ullet For a "saturated" model-class $q=2^{|\mathcal{Y}|}-1$
- $\kappa(\eta, \mathcal{Y})$ distribution normalizing constant

$$\kappa(\eta, \mathcal{Y}) = \sum_{y \in \mathcal{Y}} \exp\{\eta \cdot g(y)\}$$
Simple model-classes for social networks

Homogeneous Bernoulli graph (Rényi-Erdős model)

• Y_{ij} are independent and equally likely with log-odds $\eta = \text{logit}[P_{\eta}(Y_{ij} = 1)]$

$$egin{aligned} & P_\eta(Y=y) = rac{\mathrm{e}^{\eta\sum_{i,j}y_{ij}}}{\kappa(\eta,\mathcal{Y})} & \qquad y\in\mathcal{Y} \end{aligned}$$

where q = 1, $g(y) = \sum_{i,j} y_{ij}$, $\kappa(\eta, \mathcal{Y}) = [1 + \exp(\eta)]^N$

 homogeneity means it is unlikely to be proposed as a model for real phenomena

Dyad-independence models with attributes

• Y_{ij} are independent but depend on dyadic covariates $x_{k,ij}$

$$P_\eta(Y=y) = rac{e^{\sum_{k=1}^q \eta_k g_k(y)}}{\kappa(\eta,\mathcal{Y})} \qquad y\in\mathcal{Y}$$

Dyad-independence models with attributes

• Y_{ij} are independent but depend on dyadic covariates $x_{k,ij}$

$$P_\eta(Y=y) = rac{e^{\sum_{k=1}^q \eta_k g_k(y)}}{\kappa(\eta,\mathcal{Y})} \qquad y\in\mathcal{Y}$$

$$g_k(y) = \sum_{i,j} x_{k,ij} y_{ij}, \quad k = 1, \ldots, q$$

Dyad-independence models with attributes

• Y_{ij} are independent but depend on dyadic covariates $x_{k,ij}$

$$egin{aligned} & \mathcal{P}_\eta(\mathbf{Y}=\mathbf{y}) = rac{e^{\sum_{k=1}^q \eta_k g_k(\mathbf{y})}}{\kappa(\eta,\mathcal{Y})} & \qquad \mathbf{y}\in\mathcal{Y} \end{aligned}$$

$$g_k(y) = \sum_{i,j} x_{k,ij} y_{ij}, \quad k = 1, \dots, q$$

$$\kappa(\eta,\mathcal{Y}) = \prod_{i,j} [1 + \exp(\sum_{k=1}^{q} \eta_k x_{k,ij})]$$

Of course,

$$logit[P_{\eta}(Y_{ij}=1)] = \sum_{k} \eta_k x_{k,ij}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Some history of exponential family models for social networks

Holland and Leinhardt (1981) proposed a general dyad independence model

– Also an homogeneous version they refer to as the " p^1 " model

$$P_{\eta}(Y = y) = \frac{\exp\{\rho \sum_{i < j} y_{ij}y_{ji} + \phi y_{++} + \sum_{i} \alpha_{i}y_{i+} + \sum_{j} \beta_{j}y_{+j}\}}{\kappa(\rho, \alpha, \beta, \phi)}$$

where $\eta = (\rho, \alpha, \beta, \phi)$.

- ϕ controls the expected number of edges
- ρ represent the expected tendency toward reciprocation
- $-\alpha_i$ productivity of node *i*; β_j attractiveness of node *j*

Some history of exponential family models for social networks

Holland and Leinhardt (1981) proposed a general dyad independence model

– Also an homogeneous version they refer to as the " p^{1} " model

$$P_{\eta}(Y = y) = \frac{\exp\{\rho \sum_{i < j} y_{ij}y_{ji} + \phi y_{++} + \sum_{i} \alpha_{i}y_{i+} + \sum_{j} \beta_{j}y_{+j}\}}{\kappa(\rho, \alpha, \beta, \phi)}$$

where $\eta = (\rho, \alpha, \beta, \phi)$.

- ϕ controls the expected number of edges
- ρ represent the expected tendency toward reciprocation
- $-\alpha_i$ productivity of node *i*; β_j attractiveness of node *j*
 - Much related work and generalizations

Actor Markov statistics

 \Rightarrow Frank and Strauss (1986)

- motivated by notions of "symmetry" and "homogeneity"

Actor Markov statistics

- \Rightarrow Frank and Strauss (1986)
 - motivated by notions of "symmetry" and "homogeneity"
 - Y_{ij} in Y that do not share an actor are conditionally independent given the rest of the network

Actor Markov statistics

- \Rightarrow Frank and Strauss (1986)
 - motivated by notions of "symmetry" and "homogeneity"
 - Y_{ij} in Y that do not share an actor are conditionally independent given the rest of the network
- \Rightarrow analogous to nearest neighbor ideas in spatial statistics

Actor Markov statistics

- \Rightarrow Frank and Strauss (1986)
 - motivated by notions of "symmetry" and "homogeneity"
 - Y_{ij} in Y that do not share an actor are conditionally independent given the rest of the network
- \Rightarrow analogous to nearest neighbor ideas in spatial statistics
- Degree distribution: $d_k(y) = proportion of actors of degree k in y.$

Actor Markov statistics

- \Rightarrow Frank and Strauss (1986)
 - motivated by notions of "symmetry" and "homogeneity"
 - Y_{ij} in Y that do not share an actor are conditionally independent given the rest of the network
- \Rightarrow analogous to nearest neighbor ideas in spatial statistics
- Degree distribution: $d_k(y) = proportion of actors of degree k in y$.
- k-star distribution: $s_k(y) = proportion of k-stars in the graph y.$ (In particular,

(日) (同) (三) (三) (三) (○) (○)

 $\mathrm{s}_1=$ proportion of edges that exist between pairs of actors.)

Actor Markov statistics

- \Rightarrow Frank and Strauss (1986)
 - motivated by notions of "symmetry" and "homogeneity"
 - Y_{ij} in Y that do not share an actor are conditionally independent given the rest of the network
- \Rightarrow analogous to nearest neighbor ideas in spatial statistics
- Degree distribution: $d_k(y) = proportion of actors of degree k in y$.
- k-star distribution: $s_k(y) = proportion of k-stars in the graph y.$ (In particular,
 - $\mathrm{s}_1 = \mathsf{proportion}$ of edges that exist between pairs of actors.)

triangles:

 $t_1(y) =$ proportion of triads that from a complete sub-graph in y.



Other statistics motivated by conditional independence

- \Rightarrow Pattison and Robins (2002), Butts (2005)
- \Rightarrow Snijders, Pattison, Robins and Handcock (2004)
 - Y_{uj} and Y_{iv} in Y are conditionally independent given the rest of the network if they could not produce a cycle in the network

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Other statistics motivated by conditional independence

- \Rightarrow Pattison and Robins (2002), Butts (2005)
- \Rightarrow Snijders, Pattison, Robins and Handcock (2004)
 - Y_{uj} and Y_{iv} in Y are conditionally independent given the rest of the network if they could not produce a cycle in the network



Partial conditional dependence when four-cycle is created

This produces statistics of the form:

• edgewise shared partner distribution: $\exp_k(y) =$ proportion of edges between actors with exactly k shared partners k = 0, 1, ...



Figure: The actors in the non-directed (i, j) edge have 5 shared partners

• dyadwise shared partner distribution: $dsp_k(y) = proportion of dyads with exactly k shared partners$ k = 0, 1, ...

- identify social constructs or features
- based on intuitive notions or partial appeal to substantive theory

(ロ)、(型)、(E)、(E)、 E) の(の)

- identify social constructs or features
- based on intuitive notions or partial appeal to substantive theory
- Clusters of edges are often *transitive*: Recall t₁(y) is the proportion of triangles amongst triads

$$t_1(y) = \frac{1}{\binom{g}{3}} \sum_{\{i,j,k\} \in \binom{g}{3}} y_{ij} y_{ik} y_{jk}$$

- identify social constructs or features
- based on intuitive notions or partial appeal to substantive theory
- Clusters of edges are often *transitive*: Recall t₁(y) is the proportion of triangles amongst triads

$$t_1(y) = \frac{1}{\binom{g}{3}} \sum_{\{i,j,k\} \in \binom{g}{3}} y_{ij} y_{ik} y_{jk}$$

A closely related quantity is the proportion of triangles amongst 2-stars

$$C(y) = \frac{3 \times t_1(y)}{s_2(y)}$$

- identify social constructs or features
- based on intuitive notions or partial appeal to substantive theory
- Clusters of edges are often *transitive*: Recall t₁(y) is the proportion of triangles amongst triads

$$t_1(y) = \frac{1}{\binom{g}{3}} \sum_{\{i,j,k\} \in \binom{g}{3}} y_{ij} y_{ik} y_{jk}$$

A closely related quantity is the proportion of triangles amongst 2-stars

$$C(y) = \frac{3 \times t_1(y)}{s_2(y)}$$

Also called mean clustering coefficient

Example: A simple model-class with transitivity

$$n = 50$$
 actors $N = 1225$ pairs 10^{369} graphs
 $P(Y = y) = \frac{\exp\{\eta_1 E(y) + \eta_2 C(y)\}}{\kappa(\eta_1, \eta_2)}$ $y \in \mathcal{Y}$

where

$$E(x)$$
 is the density of edges $(0-1)$
 $C(x)$ is the triangle percent $(0-100)$

- If we set the density of the graph to have about 50 edges then the expected triangle percent is 3.8%
- Suppose we set the triangle percent large to reflect transitivity in the graph: 38%

How can we tell if the model is useful?

• Does this model capture transitivity and density in a flexible way?

◆□▶ ◆□▶ ◆∃▶ ◆∃▶ = のへで

How can we tell if the model is useful?

• Does this model capture transitivity and density in a flexible way?

- By construction, on average, graphs from this model have average density 4% and average triangle percent 38%
- If the model is a good representation of transitivity and density we expect the graphs drawn from the model to be close to these values.

• What do graphs produced by this model look like?



Distribution of Graphs from this model

Curved Exponential Family Models

Suppose that η is modeled as a function of a lower dimensional parameter: $\theta \in R^p$

$$P(Y = y) = rac{\exp\{\eta(\theta) \cdot g(y)\}}{\kappa(\theta, \mathcal{Y})}$$
 $y \in \mathcal{Y}$

Hunter and Handcock (2004)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Curved Exponential Family Models

Suppose that η is modeled as a function of a lower dimensional parameter: $\theta \in R^p$

$$P(Y = y) = rac{\exp\{\eta(\theta) \cdot g(y)\}}{\kappa(\theta, \mathcal{Y})}$$
 $y \in \mathcal{Y}$

 $Hunter \ and \ Handcock \ (2004)$ Suppose we focus on a model for network degree distribution and clustering

$$\log \left[P_{\theta}(Y=y)\right] = \eta(\phi) \cdot d(y) + \nu C(y) - \log c(\phi, \nu, \mathcal{Y}), \qquad (1)$$

where $d(x) = \{d_1(x), \ldots, d_{n-1}(x)\}$ are the network degree distribution counts.

Curved Exponential Family Models

Suppose that η is modeled as a function of a lower dimensional parameter: $\theta \in R^p$

$$P(Y = y) = \frac{\exp\{\eta(\theta) \cdot g(y)\}}{\kappa(\theta, \mathcal{Y})} \qquad y \in \mathcal{Y}$$

 $Hunter \ and \ Handcock \ (2004)$ Suppose we focus on a model for network degree distribution and clustering

$$\log \left[P_{\theta}(Y=y)\right] = \eta(\phi) \cdot d(y) + \nu C(y) - \log c(\phi, \nu, \mathcal{Y}), \qquad (1)$$

where $d(x) = \{d_1(x), \ldots, d_{n-1}(x)\}$ are the network degree distribution counts.

Any degree distribution can be specified by n - 1 or less independent parameters.

Statistical Inference for η

Base inference on the loglikelihood function,

$$\ell(\eta) = \eta \cdot g(y_{\mathrm{obs}}) - \log \kappa(\eta)$$

$$\kappa(\eta) = \sum_{\text{all possible}} \exp\{\eta \cdot g(z)\}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

graphs z

Mean-value representation of the model

Let $P_{\nu}(K = k)$ be the PMF of K, the number of ties that a randomly chosen node in the network has.

An alternative parameterization: (ϕ, ρ) where the mapping is:

$$\rho = \mathbf{E}_{\phi,\rho} \left[C(X) \right] = \sum_{y \in \mathcal{Y}} C(y) \exp \left[\eta(\phi) \cdot d(y) + \nu C(y) \right] \ge 0$$
 (2)

$$P_{\nu}(K = k) = \mathbf{E}_{\phi,\rho} [d_k(Y)] \qquad k = 0, \dots, n-1$$
(3)

– ρ is the mean clustering coefficient over networks in \mathcal{Y} .

 $-\nu$ controls the parametrization of the degree distribution

Illustrations of good models within this model-class

- village-level structure
 - *n* = 50
 - mean clustering coefficient = 15% degree distribution: Yule with scaling exponent 3.
- larger-level structure
 - -n = 1000
 - mean clustering coefficient = 15% degree distribution: Yule with scaling exponent 3.
- Attribute mixing
 - Two-sex populations
 - mean clustering coefficient = 15% degree distribution: Yule with scaling exponent 3.





Yule with zero clustering coefficient conditional on degree

Yule with clustering coefficient 15%





うくで







tripercent = 60.6

Heterosexual Yule with modest correlation

Heterosexual Yule with negative correlation





◆□ > ◆□ > ◆ Ξ > ◆ Ξ > Ξ のへで

Application to a Protein-Protein Interaction Network

- By interact is meant that two amino acid chains were experimentally identified to bind to each other.
- The network is for *E. Coli* and is drawn from the "Database of Interacting Proteins (DIP)" http://dip.doe-mbi.ucla.edu
- For simplicity we focus on proteins that interact with themselves and have at least one other interaction

- 108 proteins and 94 interactions.



Figure: A protein - protein interaction network for *E. Coli*. The nodes represent proteins and the ties indicate that the two proteins are known to interact with each other.

Statistical Inference and Simulation

- Simulate using a Metropolis-Hastings algorithm (Handcock 2002).
- Here base inference on the likelihood function
- For computational reasons, approximate the likelihood via Markov Chain Monte Carlo (MCMC)
- Use maximum likelihood estimates (Geyer and Thompson 1992)

| Parameter | est. | s.e. | |
|--------------------------------|-------|--------|--|
| Scaling decay rate (ϕ) | 3.034 | 0.3108 | |
| Correlation Coefficient (u) | 1.176 | 0.1457 | |

Table: MCMC maximum likelihood parameter estimates for the protein-protein interaction network.

Approximating the loglikelihood

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Approximating the loglikelihood

- Suppose $Y_1, Y_2, \ldots, Y_m \overset{\text{i.i.d.}}{\sim} P_{\eta_0}(Y = y)$ for some η_0 .
- Using the LOLN, the difference in log-likelihoods is

$$\ell(\eta) - \ell(\eta_0) = \log \frac{\kappa(\eta_0)}{\kappa(\eta)}$$
- Suppose $Y_1, Y_2, \ldots, Y_m \overset{\text{i.i.d.}}{\sim} P_{\eta_0}(Y = y)$ for some η_0 .
- Using the LOLN, the difference in log-likelihoods is

$$\begin{split} \ell(\eta) - \ell(\eta_0) &= \log \frac{\kappa(\eta_0)}{\kappa(\eta)} \\ &= \log \mathsf{E}_{\eta_0} \left(\exp \left\{ (\eta_0 - \eta) \cdot g(Y) \right\} \right) \end{split}$$

- Suppose $Y_1, Y_2, \ldots, Y_m \overset{\text{i.i.d.}}{\sim} P_{\eta_0}(Y = y)$ for some η_0 .
- Using the LOLN, the difference in log-likelihoods is

$$\begin{split} \ell(\eta) - \ell(\eta_0) &= \log \frac{\kappa(\eta_0)}{\kappa(\eta)} \\ &= \log \mathbf{E}_{\eta_0} \left(\exp \left\{ (\eta_0 - \eta) \cdot g(Y) \right\} \right) \\ &\approx \log \frac{1}{M} \sum_{i=1}^M \exp \left\{ (\eta_0 - \eta) \cdot (g(Y_i) - g(y_{\text{obs}})) \right\} \\ &\equiv \tilde{\ell}(\eta) - \tilde{\ell}(\eta_0). \end{split}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- Suppose $Y_1, Y_2, \ldots, Y_m \overset{\text{i.i.d.}}{\sim} P_{\eta_0}(Y = y)$ for some η_0 .
- Using the LOLN, the difference in log-likelihoods is

$$\begin{split} \ell(\eta) - \ell(\eta_0) &= \log \frac{\kappa(\eta_0)}{\kappa(\eta)} \\ &= \log \mathsf{E}_{\eta_0} \left(\exp \left\{ (\eta_0 - \eta) \cdot g(Y) \right\} \right) \\ &\approx \log \frac{1}{M} \sum_{i=1}^M \exp \left\{ (\eta_0 - \eta) \cdot (g(Y_i) - g(y_{\text{obs}})) \right\} \\ &\equiv \tilde{\ell}(\eta) - \tilde{\ell}(\eta_0). \end{split}$$

Simulate Y₁, Y₂,..., Y_m using a MCMC (Metropolis-Hastings) algorithm ⇒ Handcock (2002).

- Suppose $Y_1, Y_2, \ldots, Y_m \overset{\text{i.i.d.}}{\sim} P_{\eta_0}(Y = y)$ for some η_0 .
- Using the LOLN, the difference in log-likelihoods is

$$\begin{split} \ell(\eta) - \ell(\eta_0) &= \log \frac{\kappa(\eta_0)}{\kappa(\eta)} \\ &= \log \mathsf{E}_{\eta_0} \left(\exp \left\{ (\eta_0 - \eta) \cdot g(Y) \right\} \right) \\ &\approx \log \frac{1}{M} \sum_{i=1}^M \exp \left\{ (\eta_0 - \eta) \cdot (g(Y_i) - g(y_{\text{obs}})) \right\} \\ &\equiv \tilde{\ell}(\eta) - \tilde{\ell}(\eta_0). \end{split}$$

- Simulate Y₁, Y₂,..., Y_m using a MCMC (Metropolis-Hastings) algorithm ⇒ Handcock (2002).
- Approximate the MLE $\hat{\eta} = \operatorname{argmax}_{\eta} \{ \tilde{\ell}(\eta) \tilde{\ell}(\eta_0) \}$ (MC-MLE) \Rightarrow Geyer and Thompson (1992)

Modeling Network Dynamics

• Suppose we wish to represent the dynamics at t = 0, 1, ..., T time points

$$Y_{ijt} = \begin{cases} 1 & \text{relationship from actor } i \text{ to actor } j \text{ at time } t \\ 0 & \text{otherwise} \end{cases}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Modeling Network Dynamics

• Suppose we wish to represent the dynamics at t = 0, 1, ..., T time points

$$Y_{ijt} = \begin{cases} 1 & \text{relationship from actor } i \text{ to actor } j \text{ at time } t \\ 0 & \text{otherwise} \end{cases}$$

- Need a model that
 - has the correct cross-sectional statistics
 - has the correct durations for relationships
 - realistic dissolution and formation of relationships

A Naive Model for Longitudinal Network Data

Consider a dynamic variant of the cross-sectional ERGM:

$$P_{\eta}(Y_{t+1} = y_{t+1}|Y_t = y_t) = \frac{\exp\left(\eta_{t+1} \cdot g\left(y_{t+1}; y_t\right)\right)}{\sum_{s \in \mathcal{Y}} \exp\left(\eta_{t+1} \cdot g\left(x; y_t\right)\right)} \ t = 2, \dots, T$$

where $g_k(y_{t+1}; y_t)$ are statistics formed from y_{t+1} given y_t

- Robins and Pattison (2000) Discrete temporal ERGM
- Morris and Handcock (2001) Discrete temporal ERGM
- Hanneke and Xing (2006) Discrete temporal ERGM
- Guo, Hanneke, Fu and Xing (2007)] Hidden temporal ERGM

Two-Phase Dynamic Model

Consider a Markovian model with transition probabilities from Y_t to Y_{t+1}

governed by simultaneous dissolution and formation phases

Formation Phase



 β completely controls the *incidence*

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Dissolution Phase



 γ complete controls the durations of partnerships

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Simultaneous Formation and Dissolution



Simultaneous Formation and Dissolution

Transition Probability

 $\Pr(Y_1 = y_1 | Y_0 = y_0; \beta, \gamma) = p_-(y_1 \cap y_0 | y_0; \gamma) \times p_+(y_1 \cup y_0 | y_0; \beta)$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Markov Process

$Y_0 \rightarrow Y_1 \rightarrow Y_2 \rightarrow Y_3 \rightarrow \ldots$

Equilibrium Distribution

$Y_t \xrightarrow{\mathsf{D}} Y \sim \mathsf{Pr}(Y = y; \beta, \gamma)$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙

$Prevalence = Incidence \times Duration$

$\begin{array}{rcl} {\sf Prevalence} & = & {\sf Incidence} & \times & {\sf Duration} \\ & & || \\ & & {\sf Formation} \end{array}$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - のへで

- Data: The National Longitudinal Study of Adolescent Health
 - Wave III (with retrospective duration information)
 - Take into account age, sex, race (white/non-white) and age-sex "mixing" patterns

- Data: The National Longitudinal Study of Adolescent Health
 - Wave III (with retrospective duration information)
 - Take into account age, sex, race (white/non-white) and age-sex "mixing" patterns
- Estimate the parameters of the model based on the likelihood.

- Data: The National Longitudinal Study of Adolescent Health
 - Wave III (with retrospective duration information)
 - Take into account age, sex, race (white/non-white) and age-sex "mixing" patterns
- Estimate the parameters of the model based on the likelihood.
- Consider a (quasi)population of 10000 people (about half men and women)

- Data: The National Longitudinal Study of Adolescent Health
 - Wave III (with retrospective duration information)
 - Take into account age, sex, race (white/non-white) and age-sex "mixing" patterns
- Estimate the parameters of the model based on the likelihood.
- Consider a (quasi)population of 10000 people (about half men and women)

- Simulate dynamics of sexual networks over 10 years
 - the time step is daily (3650 steps)

- Data: The National Longitudinal Study of Adolescent Health
 - Wave III (with retrospective duration information)
 - Take into account age, sex, race (white/non-white) and age-sex "mixing" patterns
- Estimate the parameters of the model based on the likelihood.
- Consider a (quasi)population of 10000 people (about half men and women)
- Simulate dynamics of sexual networks over 10 years
 - the time step is daily (3650 steps)
- Simulate disease spread based on 10 "seeds" (2 non-white)
 - as daily have good control over micro-structure of transmission

- Data: The National Longitudinal Study of Adolescent Health
 - Wave III (with retrospective duration information)
 - Take into account age, sex, race (white/non-white) and age-sex "mixing" patterns
- Estimate the parameters of the model based on the likelihood.
- Consider a (quasi)population of 10000 people (about half men and women)
- Simulate dynamics of sexual networks over 10 years
 - the time step is daily (3650 steps)
- Simulate disease spread based on 10 "seeds" (2 non-white)
 - as daily have good control over micro-structure of transmission

• Visualize only those that become infected

Conclusions and Challenges

- Network models are a very constructive way to represent (social) theory
- Some seemingly simple models are not so.
- Large and deep literatures exist are often ignored
- Simple models are being used to capture structural properties
- The inclusion of attributes is very important
 - actor attributes
 - dyad attributes e.g. homophily, race, location
 - structural terms e.g. transitive homophily
- Software: A suite of R packages to implement this is available: statnetproject.org
- See the papers at:statnetproject.org/users_guide.shtml To appear as a special Issue of the *Journal of Statistical Software*, Volume 24.