

The exchangeable graph model with applications to dynamic network analysis

Edo Airolidi

*Computer Science Department &
Lewis-Sigler Institute for Integrative Genomics
Princeton University*

Joint work with: David Blei, Kathleen Carley, Stephen Fienberg & Eric Xing

IPAM, November 6th, 2007, Los Angeles CA

Key notions

- Complexity of observed connectivity is resolved in a structure of simple motifs and their evolution
- Mixed membership
- Dynamics
 - State-space models
 - Birth-death processes

IPAM, November 6th, 2007, Los Angeles CA

Edo Airolidi

Overview

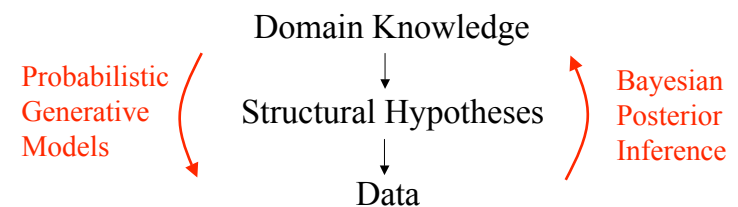
- Problem: how can we think quantitatively about social structure and social dynamics?
- Data:
 - Sampson's monastery data
 - National survey of adolescent health
 - Linked-In
- Disclaimer: do not think probability, statistical methodology or learning, rather think substantive

IPAM, November 6th, 2007, Los Angeles CA

Edo Airolidi

The role of structure

- Structural hypotheses drive inference



IPAM, November 6th, 2007, Los Angeles CA

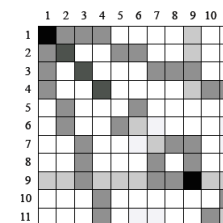
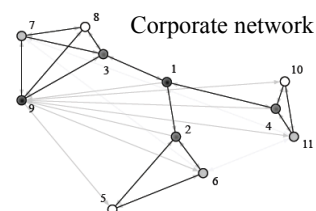
Edo Airolidi

Agenda

- Static network analysis
- Methodological themes
- Dynamics of social failure
- The exchangeable edge model
- Concluding remarks

IPAM, November 6th, 2007, Los Angeles CA

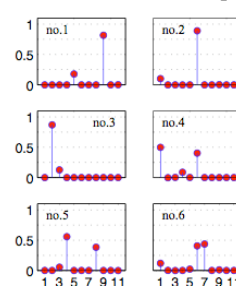
Edo Airolidi



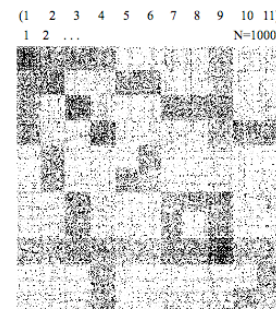
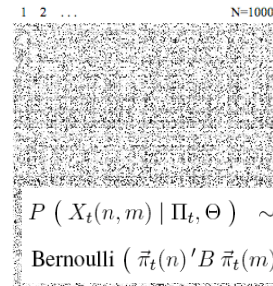
Stochastic blockmodel

Resolved data set

Mixed membership

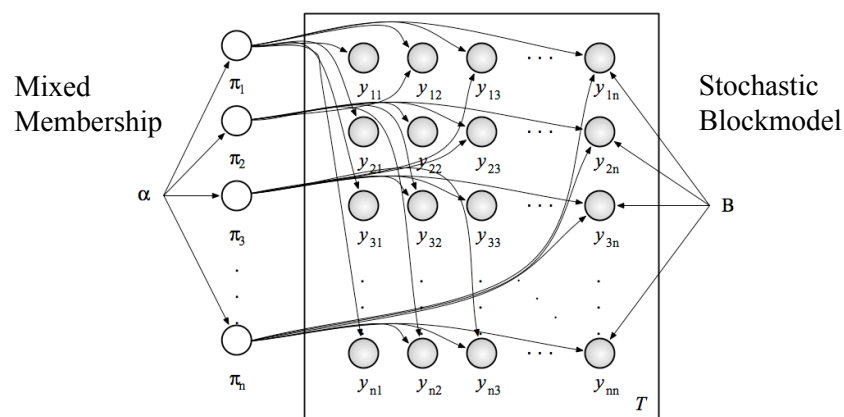


Observed data set



$$P(X_t(n, m) | \Pi_t, \Theta) \sim \text{Bernoulli}(\bar{\pi}_t(n)' B \bar{\pi}_t(m))$$

A projection onto $\Pi \times B$



IPAM, November 6th, 2007, Los Angeles CA

Edo Airolidi

Sampson's monastery data

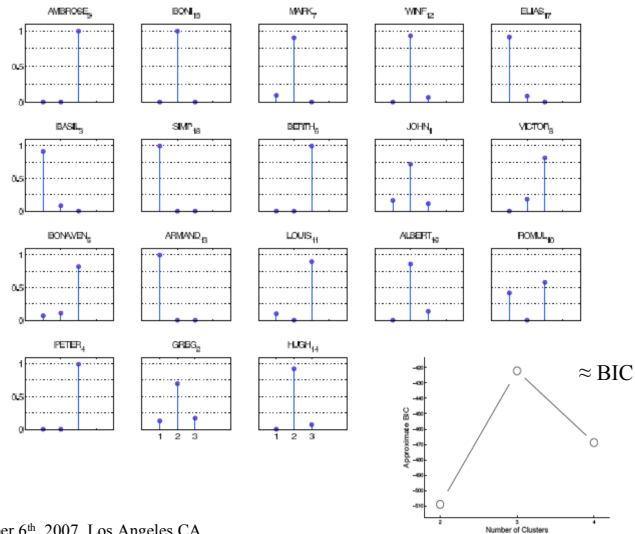
- How many factions are there?
- How do factions relate to one another?
- Who belongs to which faction?

18 Novices
Edge = Like



IPAM, November 6th, 2007, Los Angeles CA

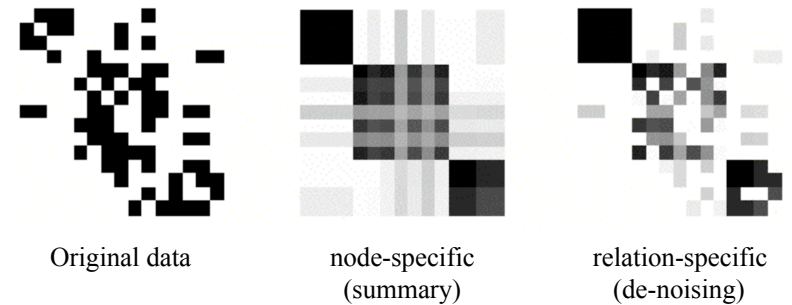
Edo Airolidi



IPAM, November 6th, 2007, Los Angeles CA

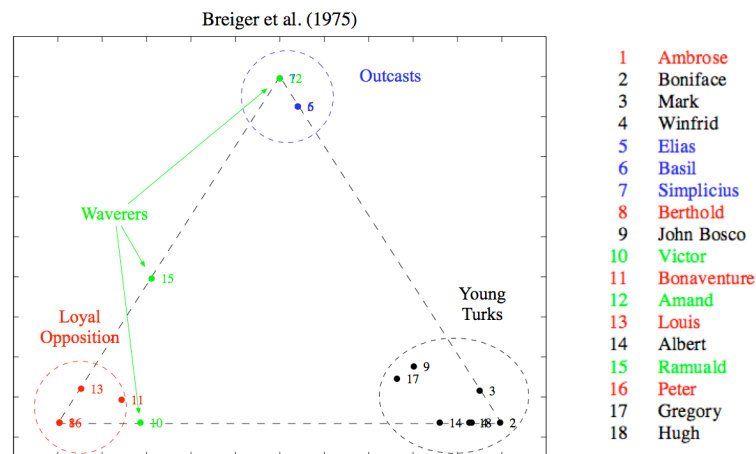
Recovering observed connectivity

- Two model variants (node-specific, relation-specific) provide increasing levels of definition



IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi



IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

National study on adolescents

- A friendship network among 69 students in grades 7-12



Fig. 8. Original matrix of friendship relations (left), and estimated relations obtained by thresholding the posterior expectations $\pi_p' B \pi_q | R$ (center), and $\phi_p' B \phi_q | R$ (right).

IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

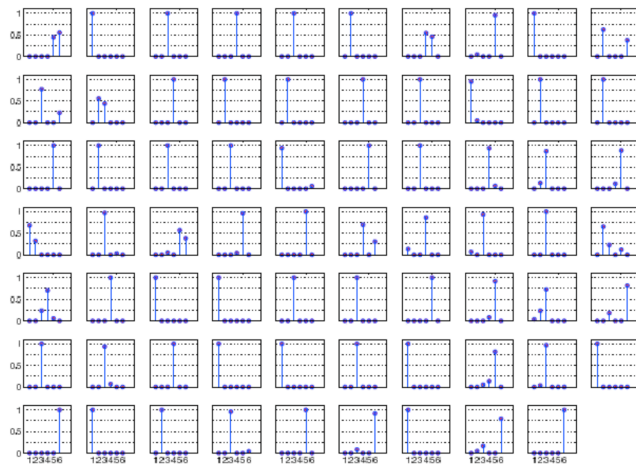


Fig. 7. The posterior mixed membership scores, π , for the 69 students in a school. Each panel correspond to a student; on the Y axis we measure the grade of membership, corresponding to the six grade levels from 7 to 12, on the X axis.

Problem revisited

- Given: A collection of relational measurement on the same sets of objects (units of analysis)
(square matrices, or unipartite graphs, with integer, real or multivariate edge weights)
- Find: (i) A pool of recurrent connectivity patterns among blocks of nodes —how many and what they look like, and (ii) A mapping of nodes to connectivity patterns — at the block level
(PCA for relational data, with symmetry constraints)

Grade	MMSB Clusters						MSB Clusters						LSCM Clusters					
	1	2	3	4	5	6	1	2	3	4	5	6	1	2	3	4	5	6
7	13	1	0	0	0	0	13	1	0	0	0	0	13	1	0	0	0	0
8	0	9	2	0	0	1	0	10	2	0	0	0	0	11	1	0	0	0
9	0	0	16	0	0	0	0	0	10	0	0	6	0	0	7	6	3	0
10	0	0	0	10	0	0	0	0	0	10	0	0	0	0	0	0	3	7
11	0	0	1	0	11	1	0	0	1	0	11	1	0	0	0	0	3	10
12	0	0	0	0	0	4	0	0	0	0	0	4	0	0	0	0	0	4

Table 1: Grade levels versus (highest) expected posterior membership for the 69 students, according to three alternative models. MMSB is the proposed mixed membership stochastic blockmodel, MSB is a simpler stochastic block mixture model (Doreian et al., 2007), and LSCM is the latent space cluster model (Handcock et al., 2007).

$$\hat{B} = \begin{bmatrix} 0.3235 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.3614 & 0.0002 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.2607 & 0.0 & 0.0 & 0.0002 \\ 0.0 & 0.0 & 0.0 & 0.3751 & 0.0009 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0002 & 0.3795 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.3719 \end{bmatrix}$$

Summary

- Observed connectivity structure is described in terms of two main sources of variability:
 - Stochastic blockmodel
 - Blocks and block-to-block connectivity patterns
(the community structure, global, asymmetric)
 - Membership map
 - Nodes-to-blocks map
(mixed membership, object-specific, symmetric)

Agenda

- Static network analysis
- Methodological themes
- Dynamics of social failure
- The exchangeable edge model
- Concluding remarks

- For each node $p \in \mathcal{N}$:
 - Draw a K dimensional mixed membership vector $\vec{\pi}_p \sim \text{Dirichlet}(\vec{\alpha})$.
- For each pair of nodes $(p, q) \in \mathcal{N} \times \mathcal{N}$:
 - Draw membership indicator for the initiator, $\vec{z}_{p \rightarrow q} \sim \text{Multinomial}(\vec{\pi}_p)$.
 - Draw membership indicator for the receiver, $\vec{z}_{q \rightarrow p} \sim \text{Multinomial}(\vec{\pi}_q)$.
 - Sample the value of their interaction, $R(p, q) \sim \text{Bernoulli}(\vec{z}_{p \rightarrow q}^\top B \vec{z}_{p \leftarrow q})$.

$$p(R, \vec{\pi}_{1:N}, Z_{\rightarrow}, Z_{\leftarrow} | \vec{\alpha}, B)$$

$$= \prod_{p,q} P(R(p, q) | \vec{z}_{p \rightarrow q}, \vec{z}_{p \leftarrow q}, B) P(\vec{z}_{p \rightarrow q} | \vec{\pi}_p) P(\vec{z}_{p \leftarrow q} | \vec{\pi}_q) \prod_p P(\vec{\pi}_p | \vec{\alpha}).$$

Inference on mixed membership

- Define: observations $Y = R$, latent variables $X = (\Pi, Z)$, and underlying constants $\Theta = (\alpha, B)$

$$\begin{aligned} p(Y|\Theta) &= \log \int_{\mathcal{X}} p(Y, X|\Theta) dX \\ &= \log \int_{\mathcal{X}} q(X) \frac{p(Y, X|\Theta)}{q(X)} dX \quad (\text{for any } q) \\ &\geq \int_{\mathcal{X}} q(X) \log \frac{p(Y, X|\Theta)}{q(X)} dX \quad (\text{Jensen's}) \\ &= \mathbb{E}_q [\log p(Y, X|\Theta) - \log q(X)] =: \mathcal{L}(q, \Theta) \end{aligned}$$

Variational approximation

- The idea is to maximize lower bound over (X, Θ)
- Alas, not possible to compute

$$q^{(t)} = p(X|Y, \Theta^{(t-1)}),$$

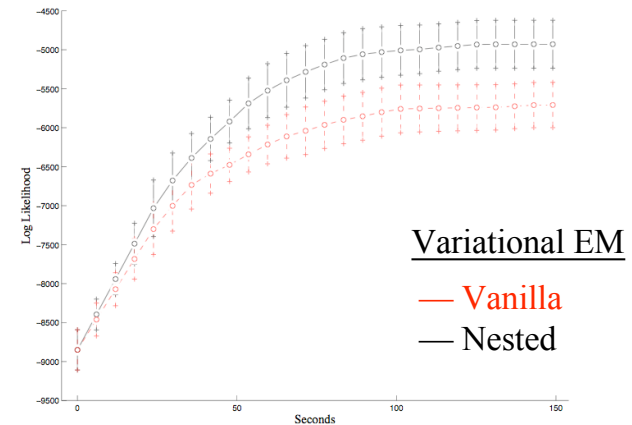
- Posit parametric approximation for q using free parameters Δ

$$q^{(t)} \approx q_{\Delta^*(Y)}^{(t)}(X) = p(X|Y).$$

Large scale computation

- Masses of data
 - 750K observations in a small problem ($N=871$)
 - 2.5M observations in a medium problem ($N=1567$)
 - Introduce parameter ρ to deal with sparsity
- Variational inference [Jordan et al., 2001]
 - Naïve implementation does not work
 - Develop a novel “nested” variational EM algorithm

Large scale computation



1. Stochastic blockmodel, B

- Captures salient structure, at the block level
(collapse nodes into groups, or blocks)
- Node-specific connectivity patterns are instances of (multiple) block-to-block connectivity patterns
- Connectivity among nodes within the same block is only specified on average

2. Mixed membership, Π

- Extends the idea of a mixture
 - Mixture: variability of data *top-down*; global weights
 - MM: variability of data *bottom-up*, unit-specific weights
- Unit-specific descriptions useful for prediction
- Sparsity: to induce parsimony in the mixed membership map between nodes and patterns
 - Enforced via prior distribution, or other means

3. Allocation paradigms, $Pr(\Pi)$

- Alternative specifications of mixed membership lead to different interpretations
- The simplex
 - Intuition: finite resources, more constrained
- The unit hyper-cube
 - Intuition: relevance, less constrained

IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

Agenda

- Static network analysis
- Methodological themes
- Dynamics of social failure
- The exchangeable edge model
- Concluding remarks

IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

Modeling social dynamics

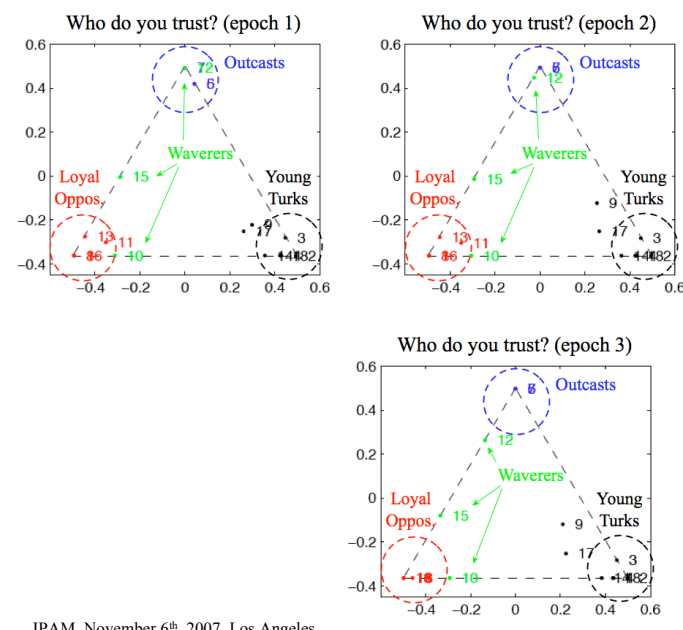
- Mixed membership analysis reduces pair-wise measurements to node-specific attributes
- Introduce smooth temporal evolution

$$X_t(n, m) \quad \text{s.t.} \quad n, m = 1, \dots, N = 18 \quad \text{and} \quad t = 1, 2, 3.$$

$$P(\vec{\pi}_0(n) | \Theta) \sim \mathbf{f} \circ \text{Gaussian}(\vec{0}, A),$$

$$P(\vec{\pi}_t(n) | \vec{\pi}_{t-1}(n), \Theta) \sim \mathbf{f} \circ [\text{Gaussian}(\vec{0}, A) + \mathbf{f}^{-1} \circ \vec{\pi}_{t-1}(n)],$$

$$P(X_t(n, m) | \Pi_t, \Theta) \sim \text{Bernoulli}(\vec{\pi}_t(n)' B \vec{\pi}_t(m)),$$



IPAM, November 6th, 2007, Los Angeles

Social failure in isolated communities

- Analysis suggests elements of a dynamic theory of social failure in isolated communities:
 1. Fragmented social structure
 2. Progressive polarization
 3. Interstitial members as traitors

An abstraction exercise

- Goal: new model of randomness for graphs
- What are the essential features of our models?
 1. Node attributes
 2. Scarcity (sparsity)
 3. Latent variables

Agenda

- Static network analysis
- Methodological themes
- Dynamics of social failure
- The exchangeable edge model
- Concluding remarks

The exchangeable graph model

- Random graphs (Erdos-Renyi-Gilbert)
 - For all (n,m) do
 $Y(n,m) \sim \text{Bernoulli}(p)$
- Exchangeable graphs
 - For all n do
 $X_k(n) \sim \text{Bernoulli}(p), k = 1 \dots K$
 - For all (n,m) do
 $Y(n,m) = f(X(n), X(m))$

Some results

- Emergence of the giant component
- Emergence of community structure
 - No phase transition
- Lognormal graphs
 - True limit connectivity
 - Scale-free graphs as approximation
- Study connectivity induced by imputing edges

IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

Agenda

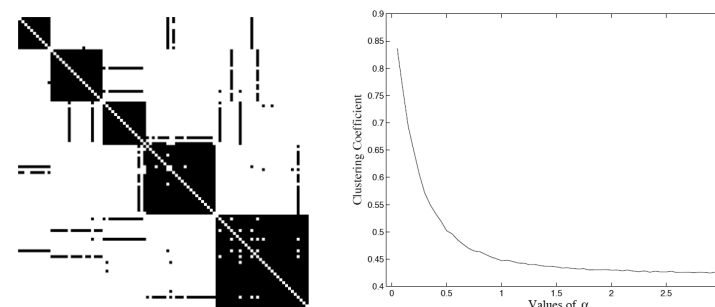
- Static network analysis
- Methodological themes
- Dynamics of social failure
- The exchangeable edge model
- Concluding remarks

IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

Emergence of community structure

- As negative correlation* among node-specific bit strings increases, communities emerge



Related work in biology

1. Inferring protein function from interacts (patches of connectivity correspond to stable complexes)
2. Statistical discovery of signaling pathways from an ensemble of weakly informative data sources

Data: interactions (e.g. Y2H), node attributes (e.g. microarrays, domains), path constraints (e.g. RNAi)

Idea: signaling pathways as latent graphs

IPAM, November 6th, 2007, Los Angeles CA

Edo Airoldi

Take home points

- Mixed membership analysis as a quantitative tool for exploring static/dynamic social networks
- The exchangeable graph model as a new paradigm for theoretical explorations of graph connectivity

Manuscripts on arXiv:

1. Stochastic blockmodel: [stat.ME 0705.4485](#)
2. Exchangeable graph model: [email me](#)