## Uncovering Structure in Biological Networks

J-J. Daudin<sup>1</sup>, V. Lacroix<sup>2</sup>, <u>M. Mariadassou<sup>1</sup></u>, V. Miele<sup>3</sup>, F. Picard<sup>3</sup>, S. Robin<sup>1</sup>, M-F. Sagot<sup>2</sup>.

<sup>1</sup>UMR INAPG/ENGREF/INRA MIA 518, Paris

<sup>2</sup>Projet HELIX, INRIA Rhône-Alpes

<sup>3</sup>UMR CNRS-8071/INRA-1152, Statistique et Génome, Évry

RIAMS 2006, December 2006

# Outline

## Motivations

- 2 An Explicit Random Graph Model
  - Some Notations
  - ER and ERMG Graph Model

#### Parametric Estimation

- Log-likelihoods and Variational Inference
- Iterative algorithm
- Model Selection Criterion

# 4 Application

# Outline

## Motivations

- 2 An Explicit Random Graph Model
  - Some Notations
  - ER and ERMG Graph Model

#### Parametric Estimation

- Log-likelihoods and Variational Inference
- Iterative algorithm
- Model Selection Criterion

## Application

・ 同 ト ・ ヨ ト ・ ヨ ト

#### Networks...

- Arise in many fields:
  - $\rightarrow$  Biology, Chemistry
  - → Physics, Internet.
- Represent an interaction pattern:
  - $\rightarrow O(n^2)$  interactions
  - $\rightarrow$  between *n* elements.
- Have a topology which:
  - → reflects the structure/function relationship



From Barabási website

# Outline

## Motivations

2 An Explicit Random Graph Model
 • Some Notations
 • ER and ERMG Graph Model

#### Parametric Estimation

- Log-likelihoods and Variational Inference
- Iterative algorithm
- Model Selection Criterion

## Application

< 回 ト < 三 ト < 三

Sac

#### Notations:

- $\rightarrow$  V a set of vertices in  $\{1, \ldots, n\}$
- $\rightarrow$  *E* a set of edges in  $\{1, \ldots, n\}^2$
- $\rightarrow$  **X** = (*X*<sub>*ij*</sub>) the adjacency matrix, with *X*<sub>*ij*</sub> =  $\mathbb{1}(i \sim j)$ .

#### Possible extensions

- $\rightarrow$  Directed graphs:  $X_{ij} \neq X_{ji}$
- $\rightarrow$  Valued graphs:  $X_{ij} \in \mathbb{N}, \mathbb{R}$

## Random graph definition:

→ The joint distribution of the  $X_{ij}$  describes the topology of the network.

ヘロト ヘ戸ト ヘヨト ヘヨト

## The Model

- $\rightarrow$  The oldest and best-known graph model,
- $\rightarrow$  The (*X<sub>ij</sub>*) are independent, with distribution  $\mathcal{B}(p)$ .

## Some Properties and Problems

 $\rightarrow$  Degree  $K_i$  of a node *i* has a Poisson distribution,

$$K_i = \sum_{j \neq i} X_{ij} \sim \mathcal{P}(\lambda)$$

 $\rightarrow$  Clustering coefficient *c* is low: *c* = *p* 

$$c = \Pr\{X_{jk} = 1 | X_{ij} = X_{ik} = 1\} = \Pr\{\nabla | \mathbf{V}\}$$

→ Highly inaccurate to describe real networks.

< ロト < 同ト < ヨト < ヨト

## Vertices heterogeneity

- $\rightarrow$  Hypothesis: the vertices are distributed among Q classes with different connectivity,
- $\rightarrow$  **Z** = (**Z**<sub>*i*</sub>)<sub>*i*</sub>,  $Z_{iq} = \mathbb{1}{i \in q}$  are indep. hidden variables,
- $\rightarrow \alpha = \{\alpha_q\}$ , the *prior* proportions of groups,
- $\rightarrow (\mathbf{Z}_i) \sim \mathcal{M}(1, \alpha).$

## X distribution

- $\rightarrow$  conditional distribution :  $X_{ij}|\{i \in q, j \in l\} \sim \mathcal{B}(\pi_{ql}),$
- $\rightarrow \pi = (\pi_{ql})$  is the connectivity matrix,
- → ERMG : "Erdös-Rényi Mixture for Graphs".

## ERMG is a model to easily generate graphs

#### Vertices heterogeneity

- $\rightarrow$  Hypothesis: the vertices are distributed among Q classes with different connectivity,
- $\rightarrow$  **Z** = (**Z**<sub>*i*</sub>)<sub>*i*</sub>,  $Z_{iq} = \mathbb{1}{i \in q}$  are indep. hidden variables,
- $\rightarrow \alpha = \{\alpha_q\}$ , the *prior* proportions of groups,
- $\rightarrow (\mathbf{Z}_i) \sim \mathcal{M}(1, \alpha).$

## X distribution

- $\rightarrow$  conditional distribution :  $X_{ij}|\{i \in q, j \in l\} \sim \mathcal{B}(\pi_{ql}),$
- $\rightarrow \pi = (\pi_{ql})$  is the connectivity matrix,
- → ERMG : "Erdös-Rényi Mixture for Graphs".

## ERMG is a model to easily generate graphs

#### • Degree distribution

$$\rightarrow K_i|\{Z_{iq}=1\} \sim \mathcal{P}(\lambda_q), \, \lambda_q=(n-1)\bar{\pi}_q, \, \bar{\pi}_q=\sum_l \alpha_l \pi_{ql},$$

$$\rightarrow K_i \sim \sum_q \alpha_q \mathcal{P}(\lambda_q),$$

- $\rightarrow$  Mixture distribution of  $K_i$  is a sub-product of ERMG,
- Clustering coefficient: ERMG and the probabilistic definition give:

$$c = \sum_{q,l,m} \alpha_q \alpha_l \alpha_m \pi_{ql} \pi_{qm} \pi_{lm} \left| \sum_{q,l,m} \alpha_q \alpha_l \alpha_m \pi_{ql} \pi_{qm} \right|.$$

# ERMG couterpart to some topologies

Description	Network	Q	π	Clustering coef.
Random		1	р	р
Stars		4	$\left(\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	0
Clusters (affiliation networks)		2	$\left(\begin{array}{cc}1 & \varepsilon\\ \varepsilon & 1\end{array}\right)$	$\frac{1+3\varepsilon^2}{(1+\varepsilon)^2}$

<ロト < 回ト < 回ト < 回ト

3

DQC

# Outline

## Motivations

- 2 An Explicit Random Graph Model
  - Some Notations
  - ER and ERMG Graph Model

#### Parametric Estimation

- Log-likelihoods and Variational Inference
- Iterative algorithm
- Model Selection Criterion

## Application

→ Ξ > < Ξ >

# Log-Likelihood of the model

First Idea: Use maximum likelihood estimators

Complete data likelihood

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}) = \sum_{i} \sum_{q} Z_{iq} \ln \alpha_{q} + \sum_{i < j} \sum_{q, l} Z_{iq} Z_{jl} \ln b(\pi_{ql}, X_{ij})$$

where 
$$b(\pi_{ql}, X_{ij}) = \pi_{ql}^{X_{ij}} (1 - \pi_{ql})^{(1 - X_{ij})}$$

Observed data likelihood

$$\mathcal{L}(\mathbf{X}) = \ln \sum_{\mathbf{Z}} \exp \mathcal{L}(\mathbf{X}, \mathbf{Z})$$

- The observed data likelihood requires a sum over *Q*<sup>*n*</sup> terms, and is thus untractable
- EM-like strategies require the knowledge of Pr(Z|X), also untractable (no conditional independence) and thus also fail.

Mariadassou (INA-PG)

# Log-Likelihood of the model

First Idea: Use maximum likelihood estimators

Complete data likelihood

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}) = \sum_{i} \sum_{q} Z_{iq} \ln \alpha_{q} + \sum_{i < j} \sum_{q, l} Z_{iq} Z_{jl} \ln b(\pi_{ql}, X_{ij})$$

where 
$$b(\pi_{ql}, X_{ij}) = \pi_{ql}^{X_{ij}} (1 - \pi_{ql})^{(1 - X_{ij})}$$

Observed data likelihood

$$\pounds(X) = ln \sum_{Z} exp \pounds(X, Z)$$

- The observed data likelihood requires a sum over *Q*<sup>*n*</sup> terms, and is thus untractable
- EM-like strategies require the knowledge of Pr(Z|X), also untractable (no conditional independence) and thus also fail.

Mariadassou (INA-PG)

# Log-Likelihood of the model

First Idea: Use maximum likelihood estimators

Complete data likelihood

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}) = \sum_{i} \sum_{q} Z_{iq} \ln \alpha_{q} + \sum_{i < j} \sum_{q, l} Z_{iq} Z_{jl} \ln b(\pi_{ql}, X_{ij})$$

where 
$$b(\pi_{ql}, X_{ij}) = \pi_{ql}^{X_{ij}} (1 - \pi_{ql})^{(1 - X_{ij})}$$

Observed data likelihood

$$\pounds(X) = ln \sum_{Z} exp \pounds(X, Z)$$

- The observed data likelihood requires a sum over *Q<sup>n</sup>* terms, and is thus untractable
- EM-like strategies require the knowledge of Pr(Z|X), also untractable (no conditional independence) and thus also fail.

Mariadassou (INA-PG) Uncovering Structure in Biological Networks RIAMS 2006, 29/12/2006 12 / 21

# **Main Idea:** Replace complicated Pr(Z|X) by a simple $\mathcal{R}_X[Z]$ such that $KL(\mathcal{R}_X[Z], Pr(Z|X))$ is minimal.

• Optimize in  $\mathcal{R}_X$  the function  $\mathcal{J}(\mathcal{R}_X)$  given by :

$$\begin{aligned} \mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) &= \mathcal{L}(\mathbf{X}) - KL(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}], \Pr(\mathbf{Z}|\mathbf{X})) \\ &= \mathcal{H}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) - \sum_{\mathbf{Z}} \mathcal{R}_{\mathbf{X}}[\mathbf{Z}]\mathcal{L}(\mathbf{X}, \mathbf{Z}) \end{aligned}$$

• For simple  $\mathcal{R}_X$ ,  $\mathcal{J}(\mathcal{R}_X[\mathbb{Z}])$  is tractable,

• At best,  $\mathcal{R}_{\mathbf{X}} = \Pr(\mathbf{Z}|\mathbf{X})$  and  $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{L}(\mathbf{X})$ .

< ロト < 同ト < ヨト < ヨト

**Main Idea:** Replace complicated Pr(Z|X) by a simple  $\mathcal{R}_X[Z]$  such that  $KL(\mathcal{R}_X[Z], Pr(Z|X))$  is minimal.

• Optimize in  $\mathcal{R}_X$  the function  $\mathcal{J}(\mathcal{R}_X)$  given by :

$$\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{L}(\mathbf{X}) - KL(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}], \Pr(\mathbf{Z}|\mathbf{X}))$$
$$= \mathcal{H}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) - \sum_{\mathbf{Z}} \mathcal{R}_{\mathbf{X}}[\mathbf{Z}]\mathcal{L}(\mathbf{X}, \mathbf{Z})$$

• For simple  $\mathcal{R}_X$ ,  $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$  is tractable,

• At best,  $\mathcal{R}_{\mathbf{X}} = \Pr(\mathbf{Z}|\mathbf{X})$  and  $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{L}(\mathbf{X})$ .

Mariadassou (INA-PG)

Uncovering Structure in Biological Networks RIAMS 2006, 29/12/2006 13 / 21

**Main Idea:** Replace complicated Pr(Z|X) by a simple  $\mathcal{R}_X[Z]$  such that  $KL(\mathcal{R}_X[Z], Pr(Z|X))$  is minimal.

• Optimize in  $\mathcal{R}_X$  the function  $\mathcal{J}(\mathcal{R}_X)$  given by :

$$\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{L}(\mathbf{X}) - KL(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}], \Pr(\mathbf{Z}|\mathbf{X}))$$
$$= \mathcal{H}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) - \sum_{\mathbf{Z}} \mathcal{R}_{\mathbf{X}}[\mathbf{Z}]\mathcal{L}(\mathbf{X}, \mathbf{Z})$$

- For simple  $\mathcal{R}_X$ ,  $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$  is tractable,
- At best,  $\mathcal{R}_{\mathbf{X}} = \Pr(\mathbf{Z}|\mathbf{X})$  and  $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{L}(\mathbf{X})$ .

# 2 Steps Iterative Algorithm

## • Step 1 Optimize $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$ w.r.t. $\mathcal{R}_X[\mathbf{Z}]$ :

- → Restriction to a "comfortable" class of functions,
- $\rightarrow \mathcal{R}_{\mathbf{X}}[\mathbf{Z}] = \prod_{i} h(\mathbf{Z}_{i}; \tau_{i})$ , with  $h(.; \tau_{i})$  the multinomial distribution,
- $\rightarrow \tau_{iq}$  is a variational parameter to be optimized using a fixed point algorithm:

$$\widehat{ au_{iq}} \propto lpha_q \prod_{j 
eq i} \prod_{l=1}^{Q} b(\pi_{ql}, X_{ij})^{\widetilde{ au}_{jl}}$$

• Step 2 Optimize  $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}])$  w.r.t.  $(\alpha, \pi)$ :

 $\rightarrow$  Constraint:  $\sum_{q} \alpha_{q} = 1$ 

$$\begin{array}{lll} \tilde{\alpha}_{q} & = & \sum_{i} \tilde{\tau}_{iq}/n \\ \tilde{\pi}_{ql} & = & \sum_{ij}^{i} \tilde{\tau}_{iq} \tilde{\tau}_{jl} X_{ij} / \sum_{ij} \tilde{\tau}_{iq} \tilde{\tau}_{jl} \end{array}$$

Mariadassou (INA-PG)

# 2 Steps Iterative Algorithm

## • Step 1 Optimize $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$ w.r.t. $\mathcal{R}_X[\mathbf{Z}]$ :

- → Restriction to a "comfortable" class of functions,
- $\rightarrow \mathcal{R}_{\mathbf{X}}[\mathbf{Z}] = \prod_{i} h(\mathbf{Z}_{i}; \tau_{i})$ , with  $h(.; \tau_{i})$  the multinomial distribution,
- $\rightarrow \tau_{iq}$  is a variational parameter to be optimized using a fixed point algorithm:

$$\tilde{\tau}_{iq} \propto \alpha_q \prod_{j \neq i} \prod_{l=1}^{Q} b(\pi_{ql}, X_{ij})^{\tilde{\tau}_{jl}}$$

## • Step 2 Optimize $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}])$ w.r.t. $(\alpha, \pi)$ :

 $\rightarrow$  Constraint:  $\sum_{q} \alpha_q = 1$ 

Mariadassou (INA-PG)

- We derive a statistical BIC-like criterion to select the number of classes:
- The likelihood can be split:  $\mathcal{L}(\mathbf{X}, \mathbf{Z}|Q) = \mathcal{L}(\mathbf{X}|\mathbf{Z}, Q) + \mathcal{L}(\mathbf{Z}|Q)$ .
- These terms can be penalized separately:

$$\mathcal{L}(\mathbf{X}|\mathbf{Z}, Q) \rightarrow \text{pen}_{\mathbf{X}|\mathbf{Z}} = \frac{Q(Q+1)}{2}\log\frac{n(n-1)}{2}$$
$$\mathcal{L}(\mathbf{Z}|Q) \rightarrow \text{pen}_{\mathbf{Z}} = (Q-1)\log(n)$$

$$ICL(Q) = \max_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{Z}} | \boldsymbol{\theta}, m_Q) - \frac{1}{2} \left( \frac{Q(Q+1)}{2} \log \frac{n(n-1)}{2} - (Q-1) \log(n) \right)$$

Mariadassou (INA-PG)

Uncovering Structure in Biological Networks RIAMS 2006, 29/12/2006 15 / 21

< ロト < 同ト < ヨト < ヨト

- We derive a statistical BIC-like criterion to select the number of classes:
- The likelihood can be split:  $\mathcal{L}(\mathbf{X}, \mathbf{Z}|Q) = \mathcal{L}(\mathbf{X}|\mathbf{Z}, Q) + \mathcal{L}(\mathbf{Z}|Q).$
- These terms can be penalized separately:

$$\mathcal{L}(\mathbf{X}|\mathbf{Z}, Q) \rightarrow \text{pen}_{\mathbf{X}|\mathbf{Z}} = \frac{Q(Q+1)}{2} \log \frac{n(n-1)}{2} \\ \mathcal{L}(\mathbf{Z}|Q) \rightarrow \text{pen}_{\mathbf{Z}} = (Q-1) \log(n)$$

$$ICL(Q) = \max_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{Z}} | \boldsymbol{\theta}, m_Q) - \frac{1}{2} \left( \frac{Q(Q+1)}{2} \log \frac{n(n-1)}{2} - (Q-1) \log(n) \right)$$

Mariadassou (INA-PG)

< ロト < 同ト < ヨト < ヨト

# Outline

## Motivations

- 2 An Explicit Random Graph Model
  - Some Notations
  - ER and ERMG Graph Model

## Parametric Estimation

- Log-likelihoods and Variational Inference
- Iterative algorithm
- Model Selection Criterion

# Application

< 6 b

A B < A B </p>

## Data and short results

#### Reaction Network of E.Coli:

- → data from http://www.biocyc.org/,
- $\rightarrow$  *n* = 605 vertices (reactions) and 1 782 edges.
- $\rightarrow$  2 reactions *i* and *j* are connected if the product of *i* is the substrate of *j* (cofactors excluded),
- → V. Lacroix and M.-F. Sagot (INRIA Hélix).

## Question:

 $\rightarrow\,$  Interpretation of the connectivity structure of classes?

## ERMG results:

- $\rightarrow$  ICL gives  $\hat{Q} = 21$  classes,
- → Most classes correspond to pseudo-cliques,
- → Clustering coefficient and degree distribution asses a good fit of the model to the data.

< ロト < 同ト < ヨト < ヨト

500

# Biological interpretation of the groups I

- Dot-plot representation
  - → adjacency matrix (sorted)
- Biological interpretation:
  - → Groups 1 to 20 gather reactions involving all the same compound either as a substrate or as a product,
  - → A compound (chorismate, pyruvate, ATP,*etc*) can be associated to each group.
- The structure of the metabolic network is governed by the compounds.



# Biological interpretation of the groups II

- → Classes 1 and 16 constitute s single clique corresponding to a single compound (pyruvate),
- → They are split into two classes because they interact differently with classes 7 (CO2) and 10 (AcetylCoA)
- $\rightarrow$  Connectivity matrix (sample):

q, l	1	7	10	16
1	1.0			
7	.11	.65		
10	.43		.67	
16	1.0	.01	$\epsilon$	1.0



Adjacency matrix (sample)

# Goodness of Fit of the ERMG Model

#### Degree distribution (histogram and PP-plot)



**Clustering coefficient** 

# Summary

## Flexibility of ERMG

- Probabilistic model which captures features of real-networks,
- Models various network topologies,
- A promising alternative to existing methods.

## Estimation and Model selection

- Variational approaches to compute approximate MLE when dependencies are complex,
- A statistical criterion to choose the number of classes (ICL).

#### **Extensions**

- Directed graphs and valued graphs
- Network motifs (cf. Sophie Schbath's talk)

Mariadassou (INA-PG)

Uncovering Structure in Biological Networks

RIAMS 2006, 29/12/2006 21 / 21

-