

# Uncovering structure in biological networks

J-J. Daudin<sup>†</sup>, V. Lacroix<sup>‡</sup>, F. Picard<sup>\*</sup>, S. Robin<sup>†</sup>, M-F. Sagot<sup>‡</sup>.

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

<sup>‡</sup> UMR 5558 Biométrie et Biologie Évolutive, Lyon,

<sup>\*</sup>UMR CNRS-8071/INRA-1152, Statistique et Génome, Évry.

Statistics for Systems Biology (SSB) group

INRA-MIG

E. Roquain, S. Schbath,

Stat. et Génome-Évry

E. Birmelé, C. Matias, V. Miele.

# The network revolution

- ▶ **Many scientific fields:**

- biology,

- sociology, physics, "internet".

- ▶ **Nature of the data under study:**

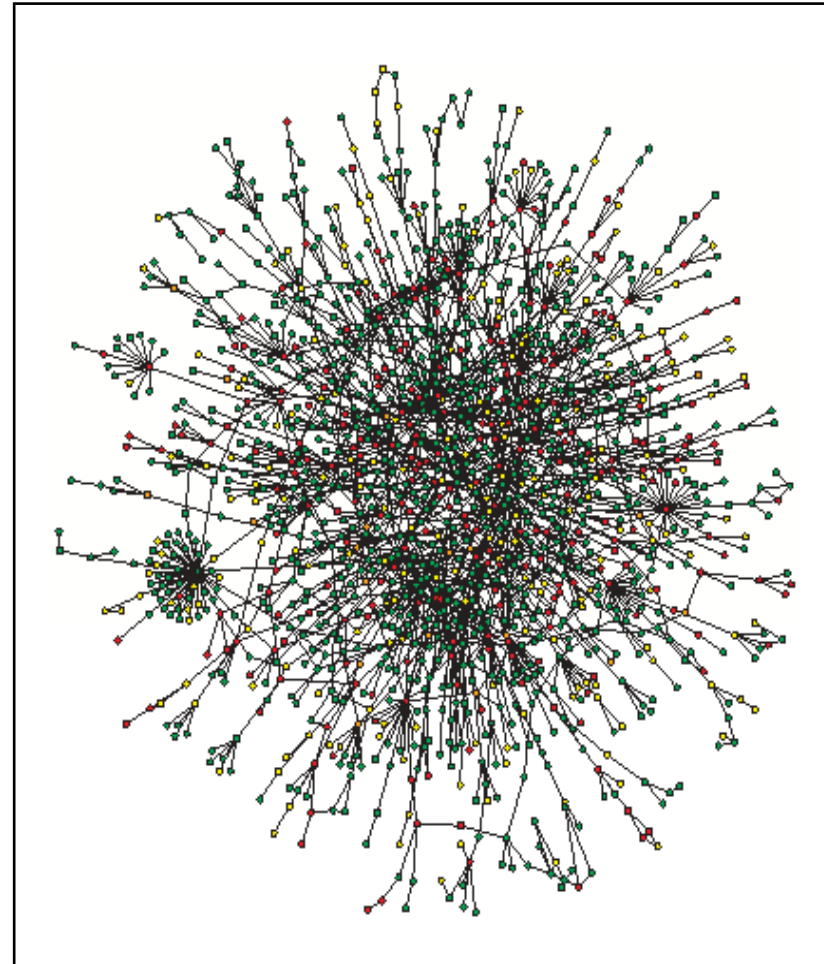
- interactions between  $n$  elements,

- $\mathcal{O}(n^2)$  possible interactions.

- ▶ **Topology of the network:**

- describes the way genes/proteins interact,

- structure/function relationship.



From Barabasi et al. (2004)

## Mathematical tool: random graphs

---

### ► Notations :

→  $V$  a set of vertices in  $\{1, \dots, n\}$ ,

→  $E$  a set of edges in  $\{1, \dots, n\}^2$ ,

→  $\mathbf{X} = (X_{ij})$  the adjacency matrix such that  $\{X_{ij} = 1\} = \mathbb{I}\{i \leftrightarrow j\}$ .

### ► Possible graphs:

→ directed:  $X_{ij} \neq X_{ji}$ ,

→ valuated:  $X_{ij} \in \mathbb{R}$ .

### ► Random graph definition :

→ the distribution of  $\mathbf{X}$  describes the topology of the network.

### ► Erdős Rényi (ER) model (1959) :

→  $(X_{ij})$  independent, with Bernoulli distribution  $\mathcal{B}(p)$ .

## The ER model and the degree distribution of real networks

- ▶ Degree distribution  $p_k$ :

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

- ▶ Real networks:

→ heterogeneous connectivity,

→ Scale-free networks:  $K_i \sim k^{-\gamma}$ .

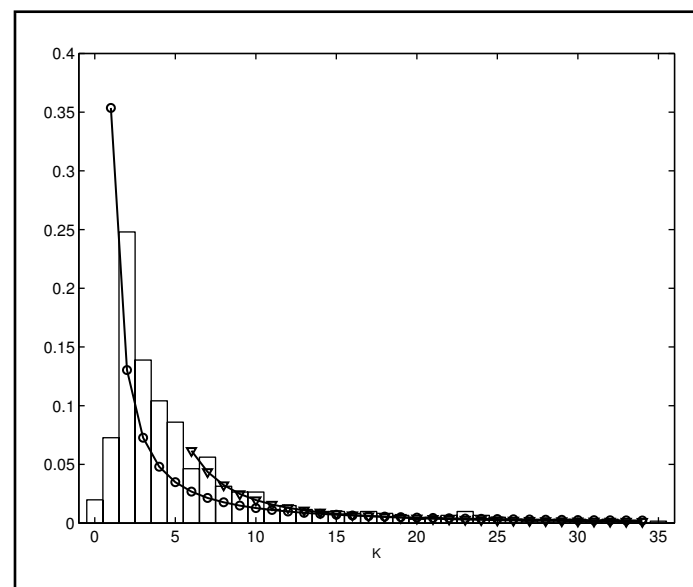
→ Mechanistic interpretation.

→ No consensus for the form of  $p_k$ .

- ▶ Current strategies :

→ description of networks using  $p_k$ ,

→ theoretical results when  $p_k$  is fixed



Network of metabolic reactions (E. Coli).

$p_k$  does not give the distribution of  $X_{ij}$ .

## The ER model and clustering

► The clustering coefficient  $c$ :

$$\rightarrow \Pr\{X_{jk} = 1 | X_{ij} = X_{ik} = 1\},$$

$$\rightarrow \Pr\{\nabla | \mathbf{V}\},$$

$$\rightarrow c = p \text{ in ER.}$$

→ real networks : high clustering coef.

→ Interpretation ?

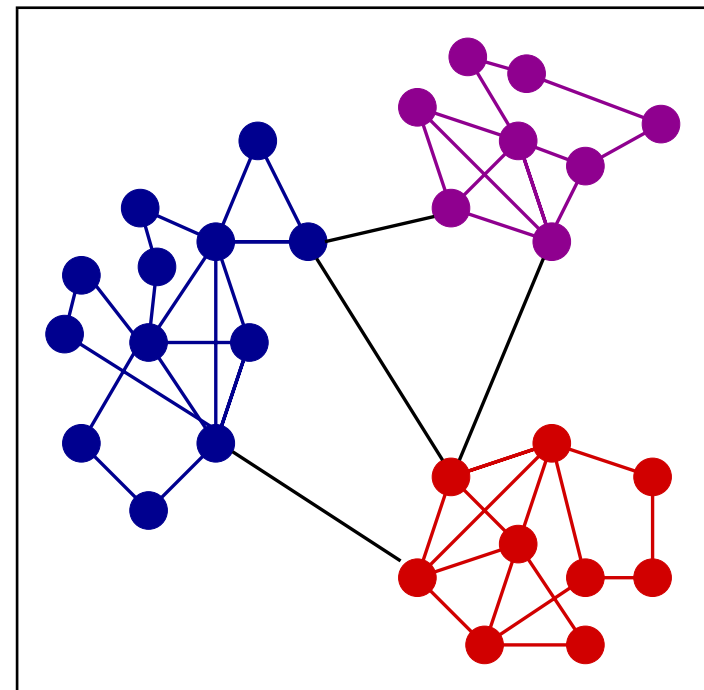
► Community structure/modularity:

→ heterogeneity intra/inter-clusters,

→ modularity of biological networks,

→ current strategies are algorithmic,

→ choosing the number of modules ?



Community structure

## ERMG: a new probabilistic model

---

### ► Modelling connection heterogeneity

→ hypothesis: there exists a hidden structure into  $Q$  classes of connectivity,

→  $\mathbf{Z} = (\mathbf{Z}_i)_i$ ,  $Z_{iq} = \mathbb{I}\{i \in q\}$  are indep. hidden variables,

→  $\boldsymbol{\alpha} = \{\alpha_q\}$ , the *prior* proportions of groups,

→  $(\mathbf{Z}_i) \sim \mathcal{M}(1, \boldsymbol{\alpha})$ .

### ► X distribution

→ conditional distribution :  $X_{ij} | \{Z_{iq}Z_{j\ell} = 1\} \sim \mathcal{B}(\pi_{q\ell})$ ,

→  $\boldsymbol{\pi} = (\pi_{q\ell})$  is the connectivity matrix.

→ Marginal distribution :  $X_{ij} \sim \sum_{q\ell} \alpha_q \alpha_\ell \mathcal{B}(\pi_{q\ell})$ ,

→ ERMG : "Erdős-Rényi Mixture for Graphs".

## Some properties of ERMG

---

► **Degree distribution**

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

→  $K_i \sim \sum_q \alpha_q \mathcal{P}(\lambda_q)$ .

→ The mixture distribution of  $K_i$  is a sub-product of ERMG.

→ It models the observed heterogeneity among degrees with an intuitive interpretation.

► **Clustering coefficient** : ERMG allows us to derive a probabilistic definition:

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

## Parameters estimation

---

► **Log-likelihood(s) of the model:**

→ Observed data :  $\mathcal{L}(\mathbf{X}) = \log (\sum_{\mathbf{Z}} \exp \mathcal{L}(\mathbf{X}, \mathbf{Z}))$ .

→ Complete data :  $\mathcal{Q}(\mathbf{X}) = \mathbb{E} [\mathcal{L}(\mathbf{X}, \mathbf{Z}) | \mathbf{X}]$ .

→ EM-like strategies require the knowledge of  $\text{Pr}(\mathbf{Z} | \mathbf{X})$ .

→ In our case, this distribution is not tractable (no conditional independence).

► **Variational methods:**

→  $\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]$  chosen such that  $KL(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}], \text{Pr}(\mathbf{Z} | \mathbf{X}))$  is minimal.

→ Optimizing  $\mathcal{J}(\mathcal{R}_{\mathbf{X}})$  w.r.t.  $\mathcal{R}_{\mathbf{X}}$  gives an approximation of  $\mathcal{L}(\mathbf{X})$  such that:

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

→ If  $\mathcal{R}_{\mathbf{X}}[\mathbf{Z}] = \text{Pr}(\mathbf{Z} | \mathbf{X})$  then  $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{L}(\mathbf{X})$ .

→ Moreover  $\mathcal{J}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) = \mathcal{H}(\mathcal{R}_{\mathbf{X}}[\mathbf{Z}]) - \sum_{\mathbf{Z}} \mathcal{R}_{\mathbf{X}}[\mathbf{Z}] \mathcal{L}(\mathbf{X}, \mathbf{Z})$  (tractable)



## An iterative algorithm

---

(h) **Optimizing**  $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$  **w.r.t.**  $\mathcal{R}_X[\mathbf{Z}]$ :

→ Restriction to a "comfortable" class of distributions,

→  $\mathcal{R}_X[\mathbf{Z}] = \prod_i h(\mathbf{Z}_i; \boldsymbol{\tau}_i)$ , with  $h(\bullet; \boldsymbol{\tau}_i)$  the multinomial distribution.

→  $\tau_{iq}$  is the variational parameter to optimize using a fixed-point algorithm:

$$\tilde{\tau}_{iq} = \Pr\{Z_{iq} = 1 | \mathbf{X}, \tilde{\mathbf{Z}}^i\}.$$

→  $\tilde{\tau}_i$  is an approximation of the conditional expectation:  $\tilde{\tau}_i = \mathbb{E}_{\mathcal{R}_X}[\mathbf{Z}_i]$ .

(h+1) **Optimizing**  $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$  **w.r.t.**  $(\boldsymbol{\alpha}, \boldsymbol{\pi})$ :

→ Constraint:  $\sum_q \alpha_q = 1$ ,

→  $\tilde{\alpha}_q = \sum_i \tilde{\tau}_{iq} / n$ ,

→  $\tilde{\pi}_{q\ell} = \sum_{ij} \tilde{\tau}_{iq} \tilde{\tau}_{j\ell} X_{ij} / \sum_{ij} \tilde{\tau}_{iq} \tilde{\tau}_{j\ell}$ .

## Model selection criterion

---

► We derive a statistical criterion to select the number of classes, using the integrated likelihood of the complete data:

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}|m_Q) = \int_{\Theta} \mathcal{L}(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}, m_Q)g(\boldsymbol{\theta}|m_Q)d\boldsymbol{\theta}.$$

► This likelihood can be split:  $\mathcal{L}(\mathbf{X}, \mathbf{Z}|m_Q) = \mathcal{L}(\mathbf{X}|\mathbf{Z}, m_Q) + \mathcal{L}(\mathbf{Z}|m_Q)$ .

► These terms can be penalized separately :

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

$$\mathcal{L}(\mathbf{Z}|m_Q) \rightarrow \text{pen}_{\mathbf{Z}} = (Q-1) \log(n).$$

$$ICL(m_Q) = \max_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{Z}}|\boldsymbol{\theta}, m_Q) - \frac{Q(Q+1)}{4} \log \frac{n(n-1)}{2} - \frac{Q-1}{2} \log(n).$$

## Application

---

### ► Reaction Network of E.Coli :

→ data from <http://www.biocyc.org/>,

→  $n = 605$  vertices (reactions) and 1 782 edges.

→ 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).

### ► ERMG results:

→ ICL gives  $\hat{Q} = 21$  classes.

→ Most classes correspond to pseudo-cliques.

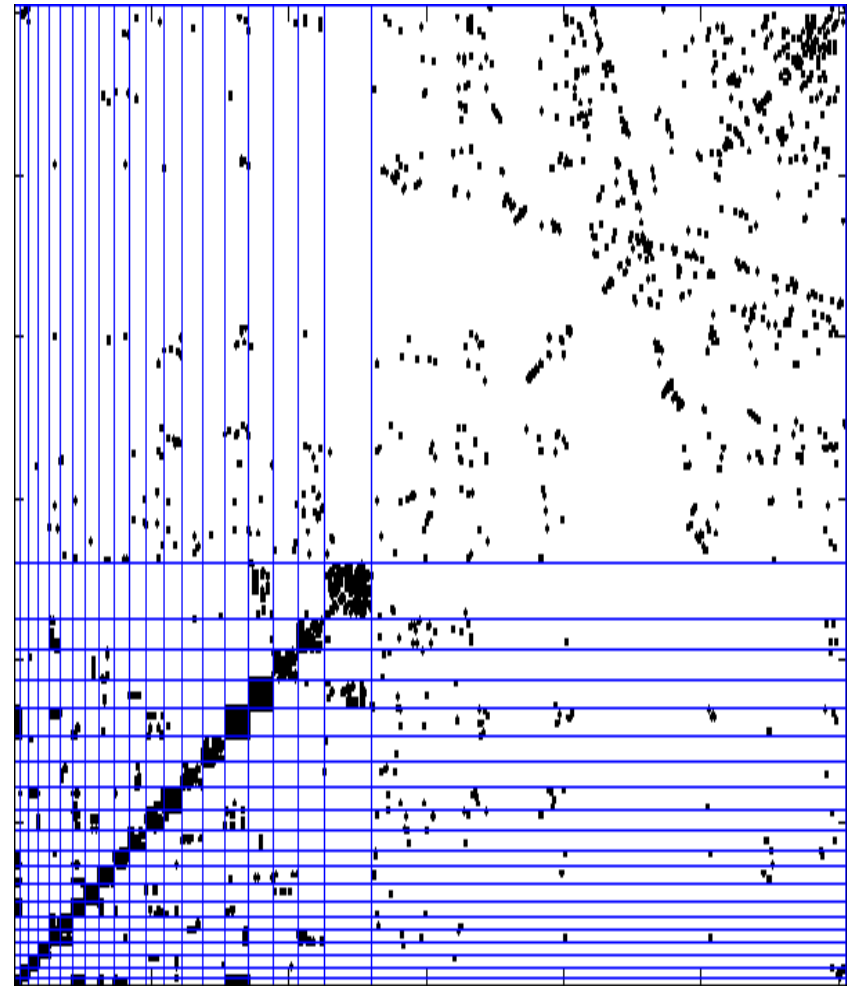
→ Interpretation of the connectivity structure of classes ?

→ Degree distribution ?

→ Clustering coefficient?

## Biological interpretation of the groups

- ▶ **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**



## Detailed example with pyruvate

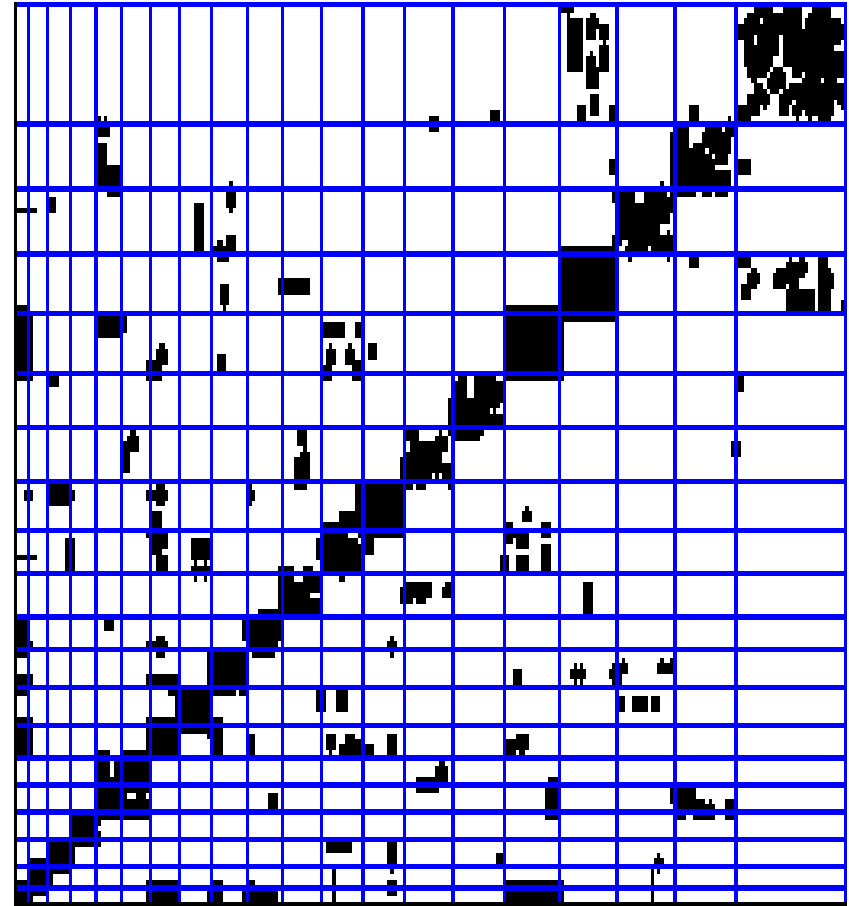
► **Biological interpretation:**

→ classes 1 and 16 constitute a clique which corresponds to a single compound (pyruvate).

→ They are split into 2 sub-cliques because of their connection with classes 7 (CO<sub>2</sub>) and 10 (AcetylCoA)

► **Connectivity matrix (sample):**

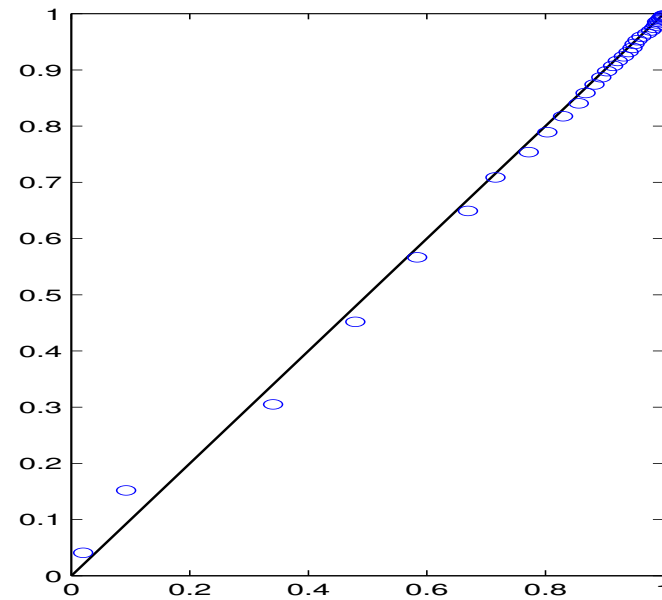
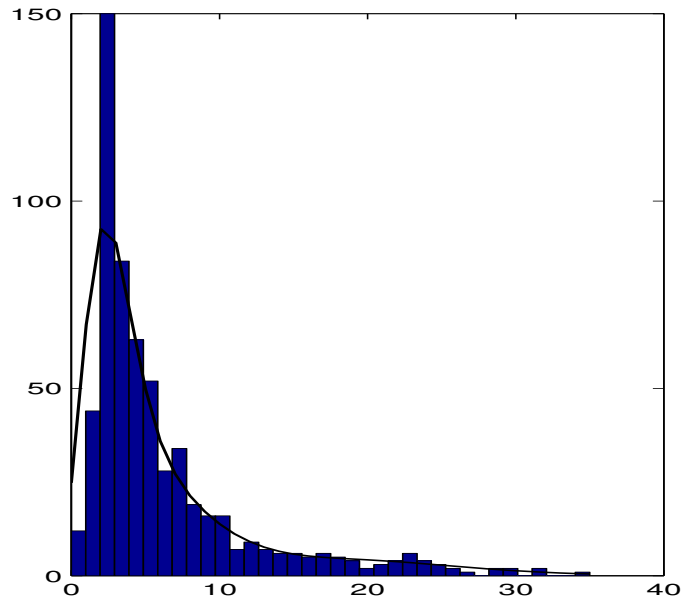
$q, \ell$	1	7	10	16
1	1.0			
7	.11	.65		
10	.43		.67	
16	1.0	.01		1.0



Adjacency matrix (sample)

# Is ERMG more realistic than other models ?

## Degree distribution (histogram and PP-plot)



## Clustering coefficient

Empirical	ERMG ( $Q = 6$ )	ERMG ( $Q = 21$ )	ER ( $Q = 1$ )
0.626	0.436	0.544	0.0098

## Conclusions

---

► **Flexibility of ERMG:**

- ERMG is a probabilistic model which captures features of real-networks,
- it can be used to model various network topologies,
- it constitutes a promising alternative to existing methods.

► **Estimation and Model selection:**

- variational approaches allow us to compute approximate MLE estimators when the dependency structure can not be simplified.
- We developed a statistical criterion to choose the number of classes (ICL).

► **Extensions:**

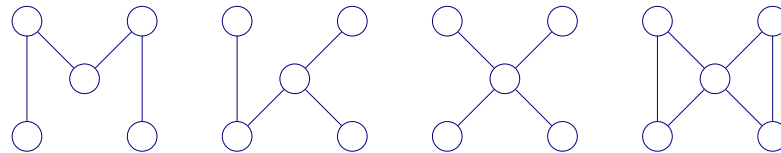
- directed graphs and regulation networks,
- valuated graphs.

## Perspectives: network motifs

---

- ▶ Network motifs provide insights regarding the local organization of a network.

- ▶ Examples of motifs



- ▶ Denoting  $N(\mathbf{m})$ , the count of motif  $\mathbf{m}$ ,  
→ Is  $N_{obs}(\mathbf{m})$  exceptional ?
- ▶ **Need of a probabilistic model under  $\mathcal{H}_0$ :**  
→ ERMG can be used for this purpose.