Uncovering Latent Structure in Valued Graphs

M. Mariadassou¹, S. Robin¹, C. Vacher²

¹UMR AgroParisTech/INRA MIA 518, Paris, France

²INRA, UMR1202 BioGeCo, Bordeaux, France

NIPS Graphs Workshop, Whistler, December 2008



Motivations for the Study of Networks

Networks...

Do exist:

- → Electricity, transport,
- → www networks,
- Represent or model interactions:
 - → Social relations,
 - → Metabolic pathways,
 - → Chemical reactions,
- Are hard to interpret as such:
 - → Many nodes/edges,
 - → Structure is hidden.



From Barabási website



NIPS08 2/17

Goal: Simple Representation of the Graph



Nodes

- *n* nodes labeled in {1,...,*n*};
- Q classes with different connectivity;
- $Z = Z_{iq}$ discrete latent variable, $Z_{iq} = 1$ if node *i* belongs to class *q*;

•
$$(Z_{i1},\ldots,Z_{iQ}) \sim \mathcal{M}(1,\alpha_1,\ldots,\alpha_Q).$$

Edges

X_{ij} = value of the edge from node *i* to node *j*: *i* → *j* Conditionally to *Z*, the *X_{ii}*s are independent with:

 $X_{ij}|i \in q, j \in l \sim f_{\theta_{ql}}$

Mariadassou (AgroParisTech)

NIPS08 4/17

Image: A matrix and a matrix

AgroParisTech

Nodes

- *n* nodes labeled in {1,...,n};
- Q classes with different connectivity;
- $Z = Z_{iq}$ discrete latent variable, $Z_{iq} = 1$ if node *i* belongs to class *q*;

•
$$(Z_{i1},\ldots,Z_{iQ}) \sim \mathcal{M}(1,\alpha_1,\ldots,\alpha_Q)$$

Edges

- X_{ij} = value of the edge from node *i* to node *j*: $i \xrightarrow{X_{ij}} j$
- Conditionally to Z, the X_{ij} s are independent with:

$$X_{ij}|i \in q, j \in l \sim f_{\theta_{ql}}$$

I > <
 I >
 I

Mariadassou (AgroParisTech)

AgroParisTech

Example

Nodes:



Edge (Poisson)



590

Mariadassou (AgroParisTech)

Uncovering Structure in Valued Graphs

э NIPS08 5/17

イロト イヨト イヨト イヨト

Example



Flexibility of MixNet

Classical Distributions:

- $\rightarrow f_{\theta}$ can be any probability distribution;
- \rightarrow Bernoulli (interaction graph): presence/absence of an edge;
- → Multinomial (labelled graph): nature of the connection (friend, lover, colleague);
- → Poisson (count): in coauthorship networks, number of copublished papers;
- \rightarrow Gaussian (intensity): intensity of the connection (airport network);

MixNet easily generates graphs.

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Pseudo-Likelihood

Maximum Likelihood Approach

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 *Qⁿ* terms, and is thus untractable;
- EM-like strategies require the knowledge of Pr(Z|X), also untractable (no conditional independence) and also fail.

Variational Inference: Pseudo Likelihood

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

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

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

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:

$$ilde{ au_{iq}} \propto lpha_q \prod_{j
eq i} \prod_{l=1}^Q f_{ heta_{ql}}(X_{ij})^{ ilde{ au_{jl}}}$$

• Step 2 Optimize $\mathcal{J}(\mathcal{R}_X[\mathbf{Z}])$ w.r.t. (α, θ) :

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

$$\begin{aligned} \widetilde{\alpha}_{q} &= \sum_{i} \widetilde{\tau}_{iq}/n \\ \widetilde{\theta}_{ql} &= \arg \max_{\theta} \sum_{ij} \widetilde{\tau}_{iq} \widetilde{\tau}_{jl} \log f_{\theta}(X_{ij}) \end{aligned}$$

 \rightarrow Simple expression of $\tilde{\theta}_{ql}$ for classical distributions



Mariadassou (AgroParisTech)

Uncovering Structure in Valued Graphs

NIPS08 8 / 17

2 Steps Iterative Algorithm

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

- \rightarrow 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:

$$ilde{ au_{iq}} \propto lpha_q \prod_{j \neq i} \prod_{l=1}^{\mathcal{Q}} f_{ heta_{ql}}(X_{ij})^{ ilde{ au}_{jl}}$$

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

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

$$\begin{aligned} \tilde{\alpha}_{q} &= \sum_{i} \tilde{\tau}_{iq} / n \\ \tilde{\theta}_{ql} &= \arg \max_{\theta} \sum_{ij} \tilde{\tau}_{iq} \tilde{\tau}_{jl} \log f_{\theta}(X_{ij}) \end{aligned}$$

 \rightarrow Simple expression of $\tilde{\theta}_{ql}$ for classical distributions.



- 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:

$$\begin{aligned} \mathcal{L}(\mathbf{X}|\mathbf{Z}, Q) &\to \quad \mathsf{pen}_{\mathbf{X}|\mathbf{Z}} P_Q \log n(n-1) \\ \mathcal{L}(\mathbf{Z}|Q) &\to \quad \mathsf{pen}_{\mathbf{Z}} = (Q-1) \log(n) \end{aligned}$$

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



< < >> < <</p>

- 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:

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

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

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



Number of Classes

- → Undirected graph with Q^* = 3 classes and Poisson edges;
- \rightarrow *n* = 50, 100, 500, 1000 vertices;
- $\rightarrow \alpha_q = (57.1\%, 28, 6\%, 14, 3\%);$
- \rightarrow Connectivity matrix of the form 1

		Q	
n	2	3	4
50	82	17	1
100	7	90	3
500	0	100	0
1000	0	100	0

Table: Frequency of selected Q for various n.



- **Dataset** Interactions between 154 fungi and 51 trees European species. Fungus *f* is connected to tree *t* if it has been collected on it.
- **Projected Graphs** For each species, we define the projected graph: for trees $X_{tt'}$ = Number of common fungi, for fungi $X_{ff'}$ = Number of common trees,
- **Poisson Model** For both species, we assume that the intensities have Poisson distributions: $X \simeq \mathcal{P}(\lambda_{ql})$.
- Number of classes. The ICL criterion selects:
 - 7 classes for trees;
 - 9 classes for fungi.

Results of Clustering



Fungi Network

Trees Network





Sac

Mariadassou (AgroParisTech)

Uncovering Structure in Valued Graphs

NIPS08 12/17

Clustered Interaction Network



- Tree classes are associated with phylogenetic history;
- Fungus classes are associated with nutritional strategy;
- Consistent with (known) tree speciation prior to fungus speciation;
- Identify generalist and specialist fungi classes.

- Flexible probabilistic model to detect structure in complex valued graphs;
- Pseudo-likelihood estimators computed obtained through variational EM;
- A statistical criterion to choose the number of classes;
- Package available at http://pbil.univ-lyon1.fr/software/MixNet.



▲ 同 ▶ ▲ 国 ▶

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:

→ Interpretation of the connectivity structure of classes?

MixNet results:

- \rightarrow ICL gives $\hat{Q} = 21$ classes,
- → Most classes correspond to pseudo-cliques,



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	ϵ	1.0



Adjacency matrix (sample)



NIPS08 17 / 17