

Bayesian Methods for Graph Clustering

P. Latouche, E. Birmelé

Laboratoire "Statistique et Génome" (UMR CNRS 8071, INRA 1152)

Journées MAS, August 2008

1 Introduction

- Real networks
- Random graph models
- The MixNet model
- Maximum likelihood estimation

2 Bayesian View of MixNet

- Bayesian probabilistic model
- Variational inference
- Model selection

3 Applications

- Affiliation models
- Metabolic network of E. coli

1 Introduction

- Real networks
- Random graph models
- The MixNet model
- Maximum likelihood estimation

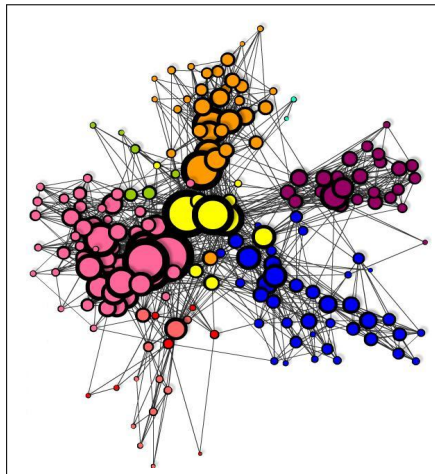
2 Bayesian View of MixNet

- Bayesian probabilistic model
- Variational inference
- Model selection

3 Applications

- Affiliation models
- Metabolic network of E. coli

- **Many scientific fields :**
 - World Wide Web,
 - Biology, sociology, physics.
- **Nature of data under study:**
 - interactions between n objects,
 - $\mathcal{O}(n^2)$ possible interactions.
- **Network topology :**
 - describes the way nodes interact, structure/function relationship.



Sample of 250 blogs (nodes) with their links (edges)

of the French political Blogosphere.



The MixNet probabilistic model

• Origin

- model developed by J. Daudin et al. (2008),
- ER model generalization,
- application fields: biology, internet, social network...

• Modelling connection heterogeneity

- hyp.: there exists a hidden structure with Q classes,
- $\mathbf{Z} = (\mathbf{Z}_i)_i$, $Z_{iq} = \mathbb{I}\{i \in q\}$ are indep. hidden variables,
- $\alpha = \{\alpha_q\}$, the *prior* proportions of groups,
- $(\mathbf{Z}_i) \sim \mathcal{M}(1, \alpha)$.

• Distribution of \mathbf{X}

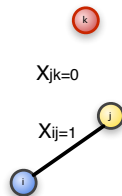
- Conditional distribution:

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

where $\mathcal{B}()$ is the *Bernoulli distribution*,

- $X_{ij} | Z$ are independant.

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



Decomposition and variational EM

$$\ln p(\mathbf{X}|\alpha, \pi) = \mathcal{L}(q(\cdot); \alpha, \pi) + \text{KL}(q(\cdot) \parallel p(\cdot|\mathbf{X}, \alpha, \pi)),$$

where

$$\mathcal{L}(q(\cdot); \alpha, \pi) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\alpha, \pi)}{q(\mathbf{Z})} \right\},$$

and

$$\text{KL}(q(\cdot) \parallel p(\cdot|\mathbf{X}, \alpha, \pi)) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \alpha, \pi)}{q(\mathbf{Z})} \right\}.$$

Approximation

$$q(\mathbf{Z}) = \prod_{i=1}^N q(\mathbf{Z}_i) = \prod_{i=1}^N \mathcal{M}(\mathbf{Z}_i; 1, \tau_i).$$

Criteria

Since $p(\mathbf{X}|\alpha, \pi)$ is not tractable, many criteria cannot be used:

- 1 Akaike Information Criterion: $AIC = \ln p(\mathbf{X}|\alpha_{ML}, \pi_{ML}) - M$.
- 2 Bayesian Information Criterion:
 $BIC = \ln p(\mathbf{X}|\alpha_{MAP}, \pi_{MAP}) - \frac{1}{2}M \ln n$.

Integrated Classification Likelihood (ICL)

- 1 Following the work of Biernacki et al. (2000), Mariadassou and Robin (2007) used a criterion based on an asymptotic approximation of the Integrated Classification Likelihood (ICL).
- 2 $ICL = \max_{\alpha, \pi} \ln p(\mathbf{X}, \tilde{\mathbf{Z}}|\alpha, \pi) - \frac{1}{2} \left(\frac{Q(Q+1)}{2} \ln(n(n+1)) - (Q-1) \ln(n) \right)$

1 Introduction

- Real networks
- Random graph models
- The MixNet model
- Maximum likelihood estimation

2 Bayesian View of MixNet

- Bayesian probabilistic model
- Variational inference
- Model selection

3 Applications

- Affiliation models
- Metabolic network of E. coli

- **Mixing coefficients:** $\alpha \sim \text{Dirichlet}(\alpha; \mathbf{n}^0)$
 - $\mathbf{n}^0 = (n_1^0, \dots, n_Q^0)$.
 - n_q^0 is the prior number of vertices in class q .
- **Connectivity matrix:** $\pi \sim \prod_{q,l} \text{Beta}(\pi_{ql}; \eta_{ql}^0, \zeta_{ql}^0)$
 - η_{ql}^0 is the prior number of edges connecting vertices of class q to vertices of class l .
 - ζ_{ql}^0 is the prior number of *non*-edges connecting vertices of class q to vertices of class l .

Decomposition

$$\ln p(\mathbf{X}) = \mathcal{L}(q(\cdot)) + \text{KL}(q(\cdot) \parallel p(\cdot|\mathbf{X})),$$

where

$$\mathcal{L}(q(\cdot)) = \sum_{\mathbf{Z}} \int \int q(\mathbf{Z}, \alpha, \pi) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}, \alpha, \pi)}{q(\mathbf{Z}, \alpha, \pi)} \right\} d\alpha d\pi,$$

and

$$\text{KL}(q(\cdot) \parallel p(\cdot|\mathbf{X})) = - \sum_{\mathbf{Z}} \int \int q(\mathbf{Z}, \alpha, \pi) \ln \left\{ \frac{p(\mathbf{Z}, \alpha, \pi|\mathbf{X})}{q(\mathbf{Z}, \alpha, \pi)} \right\} d\alpha d\pi.$$

Factorization

$$q(\mathbf{Z}, \alpha, \pi) = q(\alpha)q(\pi)q(\mathbf{Z}) = q(\alpha)q(\pi) \prod_{i=1}^N q(\mathbf{Z}_i).$$

Optimization

- 1 $\ln \tilde{q}(\mathbf{Z}_i) = E_{\mathbf{Z} \setminus i, \alpha, \pi} [\ln p(\mathbf{X}, \mathbf{Z}, \alpha, \pi)] + \text{cste.}$
- 2 $\ln \tilde{q}(\alpha) = E_{\mathbf{Z}, \pi} [\ln p(\mathbf{X}, \mathbf{Z}, \alpha, \pi)] + \text{cste.}$
- 3 $\ln \tilde{q}(\pi) = E_{\mathbf{Z}, \alpha} [\ln p(\mathbf{X}, \mathbf{Z}, \alpha, \pi)] + \text{cste.}$

Variational Bayes E-step

$$q(\mathbf{Z}_i) = \mathcal{M}(\mathbf{Z}_i; \mathbf{1}, \tau_i = \{\tau_{i1}, \dots, \tau_{iQ}\}).$$

Variational Bayes M-step (1)

$$q(\alpha) = \text{Dir}(\alpha; \mathbf{n}),$$

where $n_q = n_q^0 + \sum_{i=1}^N \tau_{iq}$.

Variational Bayes M-step (2)

$$q(\pi) = \prod_{q,l} \text{Beta}(\pi_{ql} | \eta_{ql}, \zeta_{ql}),$$

where $\eta_{ql} = \eta_{ql}^0 + \sum_{i \neq j} X_{ij} \tau_{iq} \tau_{jl}$ and $\zeta_{ql} = \zeta_{ql}^0 + \sum_{i \neq j} (1 - X_{ij}) \tau_{iq} \tau_{jl}$.

- **The model evidence depends on Q .**
- **Bayes' rule leads to $p(Q|\mathbf{X}) \propto p(\mathbf{X}|Q)p(Q)$.**
- **If $p(Q)$ is broad, maximizing $p(Q|\mathbf{X})$ is equivalent to maximizing $p(\mathbf{X}|Q)$.**
- **Since $p(\mathbf{X}|Q)$ is intractable we propose to use the lower bound $\mathcal{L}(q(\cdot))$ and to add a term $\ln Q!$ to take the multimodality into account.**
- **First non-asymptotic criterion based on an approximation of the model evidence.**

1 Introduction

- Real networks
- Random graph models
- The MixNet model
- Maximum likelihood estimation

2 Bayesian View of MixNet

- Bayesian probabilistic model
- Variational inference
- Model selection

3 Applications

- Affiliation models
- Metabolic network of E. coli

- **Probability of intra-connection** : λ .
- **Probability of inter-connection** : ϵ .
- **Number of vertices** : $n = 50$.
- **For each graph model ($\lambda + \epsilon = 1$) and for each number of classes $Q_{True} \in \{2, 3, 4, 5\}$, we generated 100 graphs.**
- **5 initializations using spectral clustering techniques.**
- **Select the best number of estimated classes according to each criterion.**

Affiliation model (1)

		1	2	3	4	5	6
	2	0	100	0	0	0	0
a) Q_{True}/Q_{ICL}	3	0	0	100	0	0	0
	4	0	0	1	98	1	0
	5	0	0	10	61	29	0
		1	2	3	4	5	6
	2	0	100	0	0	0	0
b) Q_{True}/Q_{VB}	3	0	0	100	0	0	0
	4	0	0	0	98	2	0
	5	0	0	1	29	65	5

Table: $\lambda = 0.85$ and $\epsilon = 0.15$.

Affiliation model (2)

a) Q_{True}/Q_{ICL}

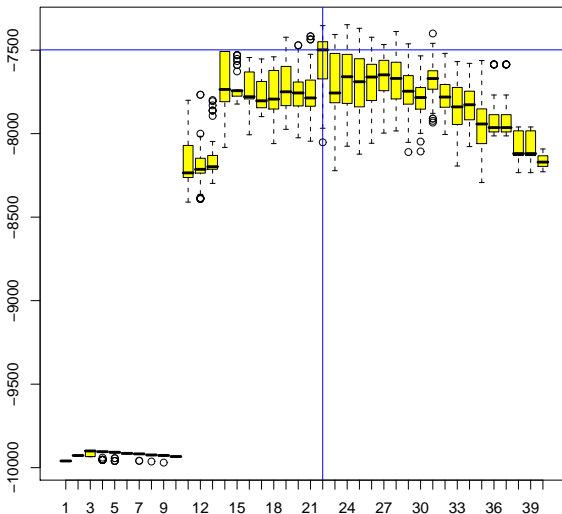
	1	2	3	4	5	6
2	0	100	0	0	0	0
3	0	0	100	0	0	0
4	0	0	14	86	0	0
5	0	17	36	44	3	0

b) Q_{True}/Q_{VB}

	1	2	3	4	5	6
2	0	100	0	0	0	0
3	0	0	100	0	0	0
4	0	0	5	94	1	0
5	0	4	18	43	29	6

Table: $\lambda = 0.8$ and $\epsilon = 0.2$.

Model selection



- **Flexibility of MixNet :**

- MixNet is a probabilistic model which captures features of real-networks,
- It considers classes of connectivity.

- **Bayesian framework**

- Estimate of the number of classes more robust,
- Can handle large graphs.

- **References :**

- Daudin J-J., Picard F., Robin S. (2008) , A mixture model for random graphs, *Statistic and Computing*
- Zanghi, H, Ambroise, C. and Miele, V. (to appear), Fast online Graph Clustering via Erdős-Rényi Mixture, *Pattern Recognition*
- J.M. Hofman and C.H. Wiggins (2008), A bayesian approach to network modularity, *Physical review letters*

- **Softwares :**

- MixNet, a C++ code (V. Miele)
<http://stat.genopole.cnrs.fr/software/mixnet>
- MixNet, a R wrapper of MixNet C++ code (available on demand)