

Analyse de réseaux sociaux à l'aide de modèles de graphe aléatoire

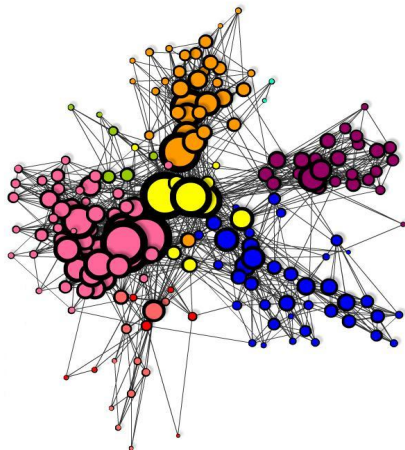
Pierre Latouche

Laboratoire Statistique et Génome

Journées MASHS, 23/06/2011



- ▶ **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.

▶ **Properties :**

- ▶ Sparsity
- ▶ Existence of a giant component
- ▶ Heterogeneity
- ▶ Preferential attachment
- ▶ Small world

↔ Topological structure (groups of vertices)

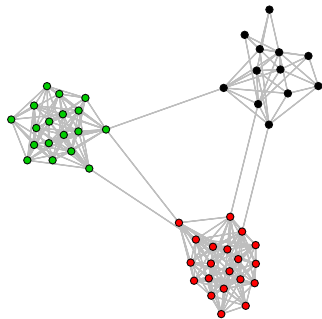
▶ **Properties :**

- ▶ Sparsity
- ▶ Existence of a giant component
- ▶ Heterogeneity
- ▶ Preferential attachment
- ▶ Small world

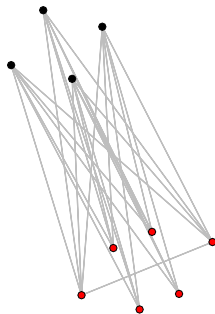
↪ Topological structure (groups of vertices)

- ▶ **Existing methods look for :**
 - ▶ Community structure
 - ▶ Disassortative mixing
 - ▶ Heterogeneous structure

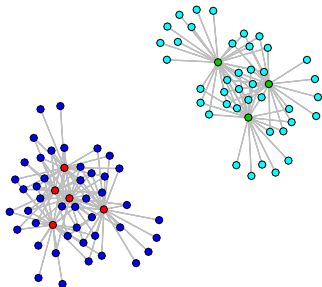
- ▶ **Existing methods look for :**
 - ▶ Community structure
 - ▶ Disassortative mixing
 - ▶ Heterogeneous structure



- ▶ **Existing methods look for :**
 - ▶ Community structure
 - ▶ Disassortative mixing
 - ▶ Heterogeneous structure



- ▶ **Existing methods look for :**
 - ▶ Community structure
 - ▶ Disassortative mixing
 - ▶ Heterogeneous structure

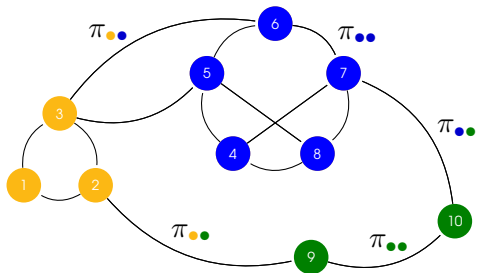


- ▶ Nowicki and Snijders (2001)
 - ▶ Earlier work : Govaert et al. (1977)
- ▶ \mathbf{Z}_i independent hidden variables :
 - ▶ $\mathbf{Z}_i \sim \mathcal{M}(1, \boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_Q))$
 - ▶ $Z_{iq} = 1$: vertex i belongs to class q
- ▶ $\mathbf{X} | \mathbf{Z}$ edges drawn independently :

$$X_{ij} | \{Z_{iq}Z_{jl} = 1\} \sim \mathcal{B}(\pi_{ql})$$

- ▶ A mixture model for graphs :

$$X_{ij} \sim \sum_{q=1}^Q \sum_{l=1}^Q \alpha_q \alpha_l \mathcal{B}(\pi_{ql})$$



- ▶ **Log-likelihoods of the model :**

- ▶ Complete-data : $\log p(\mathbf{X}, \mathbf{Z} | \alpha, \mathbf{\Pi})$
- ▶ Observed-data : $\log p(\mathbf{X} | \alpha, \mathbf{\Pi}) = \log \{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \alpha, \mathbf{\Pi}) \}$
 $\hookrightarrow Q^N$ terms

- ▶ Expectation Maximization (EM) algorithm requires the knowledge of $p(\mathbf{Z} | \mathbf{X}, \alpha, \mathbf{\Pi})$

Problem

$p(\mathbf{Z} | \mathbf{X}, \alpha, \mathbf{\Pi})$ is not tractable (no conditional independence)

Criteria

Since $\log p(\mathbf{X} | \boldsymbol{\alpha}, \boldsymbol{\Pi})$ is not tractable, we *cannot* rely on:

- ▶ $AIC = \log p(\mathbf{X} | \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Pi}}) - K$
- ▶ $BIC = \log p(\mathbf{X} | \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Pi}}) - \frac{K}{2} \log \frac{N(N-1)}{2}$

Introduction

- Real networks

- Graph clustering

Variational Bayesian inference and complexity control for stochastic block models

- Bayesian framework

- VBEM algorithm

- Model selection

- Experiments

Overlapping stochastic block models

- The model

- Inference and model selection

- Experiments

▶ **Conjugate prior distributions :**

▶ $p(\boldsymbol{\alpha} | \mathbf{n}^0 = \{n_1^0, \dots, n_Q^0\}) = \text{Dir}(\boldsymbol{\alpha}; \mathbf{n}^0)$

▶ $p(\boldsymbol{\Pi} | \boldsymbol{\eta}^0 = (\eta_{ql}^0), \boldsymbol{\zeta}^0 = (\zeta_{ql}^0)) = \prod_{q \leq l} \text{Beta}(\pi_{ql}; \eta_{ql}^0, \zeta_{ql}^0)$

▶ **Non informative Jeffreys prior :**

▶ $n_q^0 = 1/2$

▶ $\eta_{ql}^0 = \zeta_{ql}^0 = 1/2$

- ▶ $p(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\Pi} | \mathbf{X})$ not tractable

Decomposition

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

where

$$\mathcal{L}(q) = \sum_{\mathbf{Z}} \int \int q(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\Pi}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\Pi})}{q(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\Pi})} \right\} d\boldsymbol{\alpha} d\boldsymbol{\Pi}$$

Factorization

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

E-step

- ▶ $q(\mathbf{Z}_i) = \mathcal{M}(\mathbf{Z}_i; 1, \boldsymbol{\tau}_i = \{\tau_{i1}, \dots, \tau_{iQ}\})$

M-step

- ▶ $q(\boldsymbol{\alpha}) = \text{Dir}(\boldsymbol{\alpha}; \mathbf{n})$
- ▶ $q(\boldsymbol{\Pi}) = \prod_{q \leq l} \text{Beta}(\pi_{ql}; \eta_{ql}, \zeta_{ql})$

A new model selection criterion : ILvb

Latouche et al. (2011a)

- ▶ $\log p(\mathbf{X} | Q) = \mathcal{L}(q) + \text{KL}(\dots)$
- ▶ After convergence, use $\mathcal{L}(q)$ as an approximation of $\log p(\mathbf{X} | Q)$

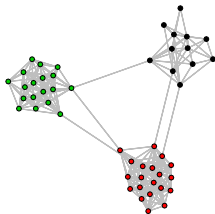
ILvb

$$IL_{vb} = \log \left\{ \frac{\Gamma(\sum_{q=1}^Q n_q^0) \prod_{q=1}^Q \Gamma(n_q)}{\Gamma(\sum_{q=1}^Q n_q) \prod_{q=1}^Q \Gamma(n_q^0)} \right\} \\ + \sum_{q \leq l}^Q \log \left\{ \frac{\Gamma(\eta_{ql}^0 + \zeta_{ql}^0) \Gamma(\eta_{ql}) \Gamma(\zeta_{ql})}{\Gamma(\eta_{ql} + \zeta_{ql}) \Gamma(\eta_{ql}^0) \Gamma(\zeta_{ql}^0)} \right\} - \sum_{i=1}^N \sum_{q=1}^Q \tau_{iq} \log \tau_{iq}$$

► **Two topological structures :**

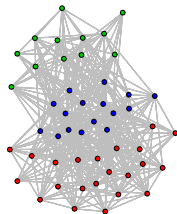
- Affiliation :

$$\mathbf{\Pi} = \begin{pmatrix} \lambda & \epsilon & \dots & \epsilon \\ \epsilon & \lambda & & \vdots \\ \vdots & & \ddots & \epsilon \\ \epsilon & \dots & \epsilon & \lambda \end{pmatrix}$$



- Affiliation and a class of hubs :

$$\mathbf{\Pi} = \begin{pmatrix} \lambda & \epsilon & \dots & \epsilon & \lambda \\ \epsilon & \lambda & & & \vdots \\ \vdots & & \ddots & & \vdots \\ \lambda & \dots & \dots & \dots & \lambda \end{pmatrix}$$



- ▶ $N = 50$
- ▶ $\lambda = 0.9$
- ▶ $\epsilon = 0.1$
- ▶ $\alpha_Q = 1/Q$
- ▶ $Q_{True} \in \{3, \dots, 7\}$
- ▶ $Q \in \{1, \dots, 7\}$
- ▶ 100 simulations
- ▶ 2 criteria : ICL (Biernacki et al. 2000, Daudin et al., 2008), lLvb

(a) $Q_{True} \setminus Q_{ICL}$

	2	3	4	5	6	7
3	0	100	0	0	0	0
4	0	0	100	0	0	0
5	0	0	23	77	0	0
6	0	1	28	59	12	0
7	0	8	49	42	1	0

(b) $Q_{True} \setminus Q_{ILvb}$

	2	3	4	5	6	7
3	0	100	0	0	0	0
4	0	0	100	0	0	0
5	0	0	0	99	1	0
6	0	0	4	23	73	0
7	0	2	14	44	27	13

Affiliation networks and a class of hubs

(c) $Q_{True} \setminus Q_{ICL}$

	2	3	4	5	6	7
3	0	100	0	0	0	0
4	0	0	100	0	0	0
5	0	0	12	88	0	0
6	0	0	19	59	22	0
7	0	3	29	56	12	0

(d) $Q_{True} \setminus Q_{ILvb}$

	2	3	4	5	6	7
3	0	100	0	0	0	0
4	0	0	100	0	0	0
5	0	0	2	98	0	0
6	0	0	1	29	70	0
7	0	0	3	34	45	18

- ▶ variational Bayes to approximate $p(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\Pi} | \mathbf{X})$
- ▶ computational cost : $O(Q^2 N^2)$
- ▶ model selection criterion : ILvb
- ▶ provides a relevant estimation of the number of classes
- ▶ implemented in a R package available on the CRAN :
mixer

Introduction

- Real networks

- Graph clustering

Variational Bayesian inference and complexity control for stochastic block models

- Bayesian framework

- VBEM algorithm

- Model selection

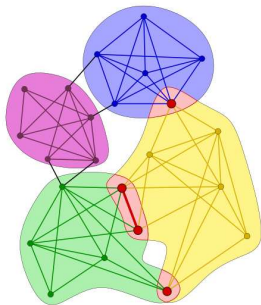
- Experiments

Overlapping stochastic block models

- The model

- Inference and model selection

- Experiments



Palla et al. (2005)

Problem

The stochastic block model (SBM) and most existing methods assume that each vertex belongs to a single class

- ▶ Nowicki and Snijders (2001)
- ▶ \mathbf{Z}_i independent hidden variables :

$$\mathbf{Z}_i \sim \mathcal{M}\left(1, \boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_Q)\right)$$

- ▶ Latouche et al. (2011b)
- ▶ Z_{iq} independent hidden variables :

$$\mathbf{z}_i \sim \prod_{q=1}^Q \mathcal{B}(Z_{iq}; \alpha_q) = \prod_{q=1}^Q \alpha_q^{Z_{iq}} (1 - \alpha_q)^{1-Z_{iq}}$$

- ▶ Latouche et al. (2011b)
- ▶ Z_{iq} independent hidden variables :

$$\mathbf{Z}_i \sim \prod_{q=1}^Q \mathcal{B}(Z_{iq}; \alpha_q) = \prod_{q=1}^Q \alpha_q^{Z_{iq}} (1 - \alpha_q)^{1-Z_{iq}}$$

- ▶ $\mathbf{X} | \mathbf{Z}$ edges drawn independently :

$$X_{ij} | \mathbf{Z}_i, \mathbf{Z}_j \sim \mathcal{B}(X_{ij}; g(a_{\mathbf{Z}_i, \mathbf{Z}_j}))$$

- ▶ $g(t) = 1 / (1 + \exp(-t))$ is the logistic function

$$a_{\mathbf{Z}_i, \mathbf{Z}_j} = \mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_j + \mathbf{Z}_i^\top \mathbf{U} + \mathbf{V}^\top \mathbf{Z}_j + W^*$$

- ▶ $\tilde{\mathbf{Z}}_i = (\mathbf{z}_i, 1)^\top$
- ▶ $\tilde{\mathbf{W}} = \begin{pmatrix} \mathbf{W} & \mathbf{U} \\ \mathbf{V}^\top & W^* \end{pmatrix}$
- ▶ $a_{\mathbf{z}_i, \mathbf{z}_j} = \tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{W}} \tilde{\mathbf{Z}}_j$
- ▶ Parameter set : $\{\alpha, \tilde{\mathbf{W}}\}$

► **Conjugate prior distributions :**

- $p(\boldsymbol{\alpha}) = \prod_{q=1}^Q \text{Beta}(\alpha_q; \eta_q^0, \zeta_q^0)$
- $p(\tilde{\mathbf{W}}^{\text{vec}}) = \mathcal{N}(\tilde{\mathbf{W}}^{\text{vec}}; \tilde{\mathbf{W}}_0^{\text{vec}}, \mathbf{S}_0)$

► The vec operator : if

$$\mathbf{A} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix},$$

then

$$\mathbf{A}^{\text{vec}} = \begin{pmatrix} A_{11} \\ A_{21} \\ A_{12} \\ A_{22} \end{pmatrix}$$

Problem

$p(\mathbf{Z}, \boldsymbol{\alpha}, \tilde{\mathbf{W}} \mid \mathbf{X})$ not tractable

Decomposition

$$\log p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p)$$

where

$$\mathcal{L}(q) = \sum_{\mathbf{Z}} \int q(\mathbf{Z}, \boldsymbol{\alpha}, \tilde{\mathbf{W}}) \log \left(\frac{p(\mathbf{X} | \mathbf{Z}, \tilde{\mathbf{W}}) p(\mathbf{Z} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) p(\tilde{\mathbf{W}})}{q(\mathbf{Z}, \boldsymbol{\alpha}, \tilde{\mathbf{W}})} \right) d\boldsymbol{\alpha} d\tilde{\mathbf{W}}$$

Lower bound

$$\log p(\mathbf{X}) \geq \mathcal{L}(q)$$

Problem

$\mathcal{L}(q)$ has a too complex form \leftrightarrow no variational Bayes EM algorithm ??

Decomposition

$$\mathcal{L}(\xi) = \mathcal{L}(q; \xi) + \text{KL}(q||p)$$

where

$$\mathcal{L}(q; \xi) = \sum_{\mathbf{Z}} \int q(\mathbf{Z}, \alpha, \tilde{\mathbf{W}}) \log \left(\frac{h(\mathbf{Z}, \tilde{\mathbf{W}}, \xi) p(\mathbf{Z} | \alpha) p(\alpha) p(\tilde{\mathbf{W}})}{q(\mathbf{Z}, \alpha, \tilde{\mathbf{W}})} \right)$$

Lower bound

$$\log p(\mathbf{X}) \geq \mathcal{L}(\xi) \geq \mathcal{L}(q; \xi)$$

Local optimization

- ▶ $\xi = \operatorname{argmax}_{\xi} \mathcal{L}(q; \xi)$

E-step

- ▶ $q(Z_{iq}) = \mathcal{B}(Z_{iq}; \tau_{iq})$

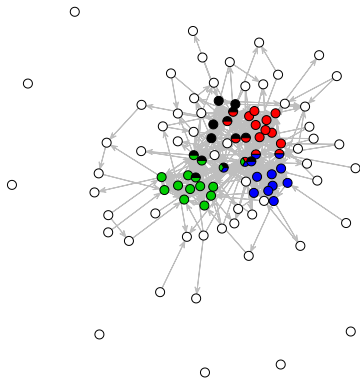
M-step

- ▶ $q(\alpha) = \prod_{q=1}^Q \operatorname{Beta}(\alpha_q; \eta_q^N, \zeta_q^N)$
- ▶ $q(\tilde{\mathbf{W}}^{vec}) = \mathcal{N}(\tilde{\mathbf{W}}^{vec}; \tilde{\mathbf{W}}_N^{vec}, \mathbf{S}_N)$

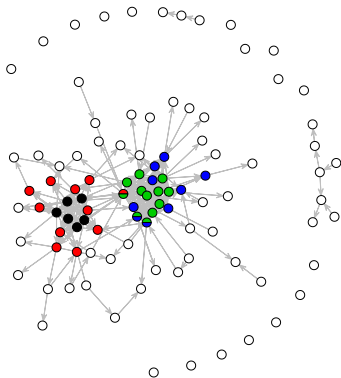
- ▶ After convergence, use $\mathcal{L}(q; \xi)$ as an approximation of $\log p(\mathbf{X} | Q)$

IL_{osbm}

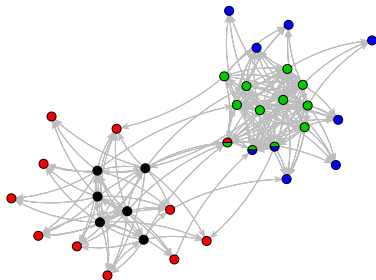
$$IL_{osbm} = \mathcal{L}(\hat{q}; \hat{\xi})$$



Example of an overlapping stochastic block model (OSBM) network with community structures.



Example of an overlapping stochastic block model (OSBM) network with community structures and stars.



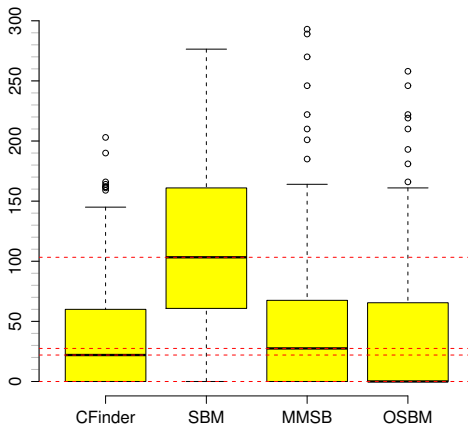
Example of an overlapping stochastic block model (OSBM) network with community structures and stars.

- ▶ $N = 100$
- ▶ $\lambda = 4$
- ▶ $\epsilon = 1$
- ▶ $W^* = -5.5$
- ▶ $\mathbf{U} = \mathbf{V} = (\epsilon \ \dots \ \epsilon)$
- ▶ $\alpha_q = 0.25$
- ▶ $Q = 4$
- ▶ 100 simulations
- ▶ 4 graph clustering methods :
 - ▶ CFinder (Palla et al. 2006)
 - ▶ Stochastic Block Model (SBM)
 - ▶ Mixed Membership Stochastic Block Model (MMSB) (Airoldi et al. 2008)
 - ▶ Overlapping Stochastic Block Model (OSBM)

How to compare the methods ?

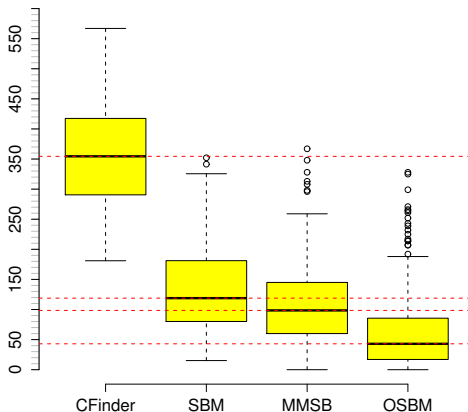
- ▶ CFinder and OSBM can deal with outliers ($\mathbf{Z}_i = \mathbf{0}$)
- ▶ SBM and MMSB are run with $Q + 1$ classes
↪ identify the class of outliers
- ▶ Compute $\mathbf{P} = \mathbf{Z}\mathbf{Z}^\top$ and $\hat{\mathbf{P}} = \hat{\mathbf{Z}}\hat{\mathbf{Z}}^\top$:
 - ▶ invariant to column permutations of \mathbf{Z} and $\hat{\mathbf{Z}}$
 - ▶ number of shared clusters between each pair of vertices
- ▶ Compute L_2 distance $d(\mathbf{P}, \hat{\mathbf{P}})$

Community structures



L_2 distance $d(\mathbf{P}, \hat{\mathbf{P}})$ over the 100 samples of networks with community structures for CFinder, SBM, MMSB and OSBM.

Community structures and stars



L_2 distance $d(\mathbf{P}, \hat{\mathbf{P}})$ over the 100 samples of networks with community structures for CFinder, SBM, MMSB and OSBM.

The French blogosphere network

	UMP	UDF	liberal	PS	analysts	others
cluster 1	30 + 3	0 + 1	0	0	0 + 1	0
cluster 2	2 + 3	29 + 1	0	0	1 + 3	0
cluster 3	0	0	24	0	1 + 1	0
cluster 4	0	0 + 2	0	40	0 + 4	1
outliers	5	1	1	17	5	30

Classification of the blogs into $Q = 4$ clusters using OSBM. 196 vertices, 2864 edges.

- ▶ A new random graph model : the overlapping stochastic block model (OSBM)
- ▶ Frequentist and Bayesian inference procedures
- ▶ Computational cost : $O(Q^4 N^2)$
- ▶ New model selection criterion : l_{osbm}
- ▶ R package **OSBM** soon available on the CRAN

-  K. Nowicki, T.A.B. Snijders. *Estimation and prediction for stochastic blockstructures*, Journal of the American Statistical Association, 2001.
-  G. Palla, I. Derenyi, I. Farkas, T. Vicsek. *Uncovering the overlapping community structure of complex networks in nature and society*, Nature, 2005.
-  J. Daudin, F. Picard, S. Robin. *A mixture model for random graphs*, Statistics and Computing, 2008.
-  E.M. Airoldi, D.M. Blei, S.E. Fienberg, E.P. Xing. *Mixed membership stochastic blockmodels*, Journal of Machine Learning Research, 2008.
-  P. Latouche, E. Birmelé, C. Ambroise. *Bayesian methods for graph clustering*, Advances in data handling and business intelligence, 2009.
-  P. Latouche, E. Birmelé, C. Ambroise. *Variational Bayesian inference and complexity control for stochastic block models*, Statistical Computing, 2011.
-  P. Latouche, E. Birmelé, C. Ambroise. *Overlapping stochastic block models with application to the French political blogosphere*, Annals of Applied Statistics, 2011.