

A Stochastic Block Model for Multilevel Networks

Saint-Clair Chabert-Liddell

Joint work with S. Donnet and P. Barbillon

UMR INRAe/AgroParisTech, MIA Paris
Séminaire des doctorants du LMO

3 June 2020

Outline

- 1 Modelling
- 2 Inference
- 3 Model Selection
- 4 Simulation Studies
- 5 Application to Real Dataset

Motivation Data Set

	$\overbrace{\hspace{2cm}}^{n_I}$		$\overbrace{\hspace{2cm}}^{n_O}$	
Individual 1 ⋮ Individual n_I	0 $X'_{ii'}$ 1	1	0 A_{ij} 0	- 1 1
Organization 1 ⋮ Organization n_O			1 $X^O_{jj'}$ 0	1 1
	Individual 1 ⋮ Individual n_I		Organization 1 ⋮ Organization n_O	

Data :

- X^I Interaction between individuals (advice ...)
- X^O Interaction between organizations (contract ...)
- A Affiliation of the individuals to the organizations
 $A_{ij} = 1$ if i is affiliated to j
 Only one affiliation per individual

Objectives

- Joint probabilistic model on $\mathbf{X} = \{X^I, X^O\}$ given A
- Evaluate the influence of the inter-organizational level on the inter-individual level

Outline

- 1 Modelling
- 2 Inference
- 3 Model Selection
- 4 Simulation Studies
- 5 Application to Real Dataset

Modelling of a Multilevel SBM



Stochastic Block Model (SBM)^a

^aSnijders and Nowicki, 1997

- Mixture model for graphs
- Latent variables on vertices
- Model heterogeneity of connection

Modelling of a Multilevel SBM

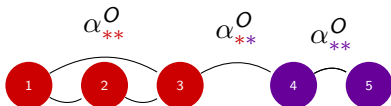


Inter-organizational Level

- n_O laboratories into Q_O clusters
- Latent variables are independent
- $Z_j^O = l \Leftrightarrow j \in l, \quad l \in \{1, \dots, Q_O\}$

$$\mathbb{P}(Z_j^O = l) = \pi_l^O$$

Modelling of a Multilevel SBM



Inter-organizational Level

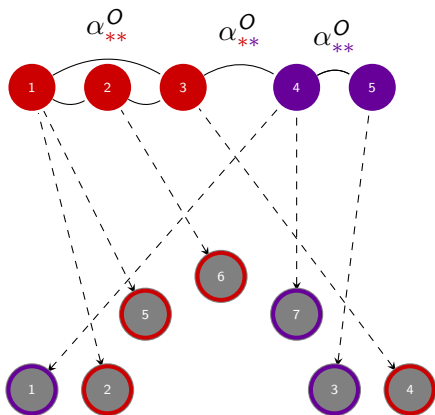
- n_O laboratories into Q_O clusters
- Latent variables are independent
- $Z_j^O = l \Leftrightarrow j \in l, \quad l \in \{1, \dots, Q_O\}$

$$\mathbb{P}(Z_j^O = l) = \pi_l^O$$

- Connectivity is independent given the latent variables

$$\mathbb{P}(X_{jj'}^O = 1 | Z_j^O = l, Z_{j'}^O = l') = \alpha_{ll'}^O$$

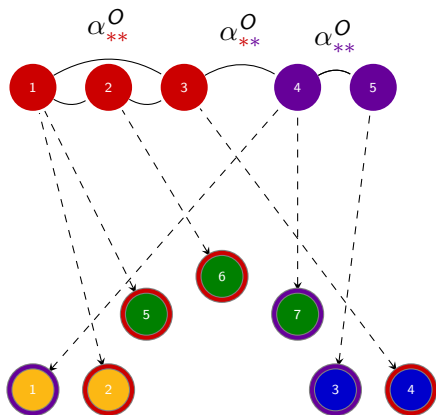
Modelling of a Multilevel SBM



Inter-individual Level

- n_I researchers into Q_I clusters
- A researcher's cluster depends on his laboratory's cluster

Modelling of a Multilevel SBM

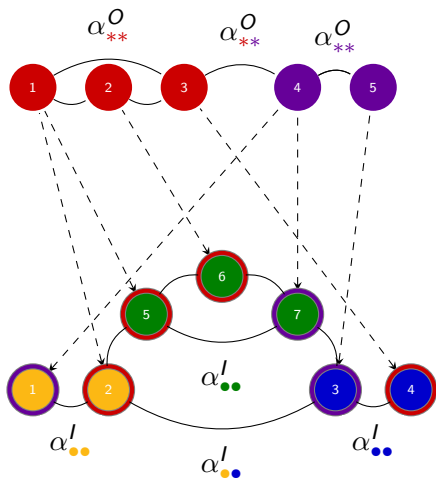


Inter-individual Level

- n_I researchers into Q_I clusters
- A researcher's cluster depends on his laboratory's cluster
- $Z_i^I = k \Leftrightarrow i \in k, k \in \{1, \dots, Q_I\}$

$$\mathbb{P}(Z_i^I = k | A_i = j, Z_j^O = l) = \gamma_{kl}$$

Modelling of a Multilevel SBM



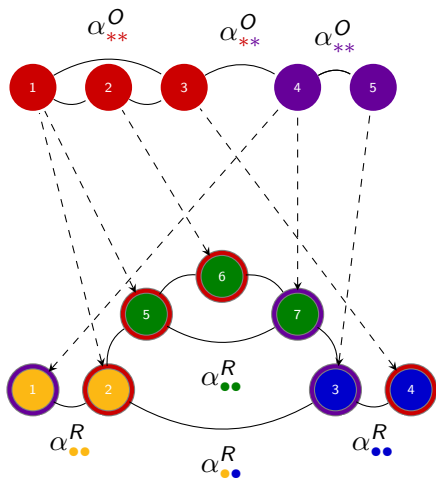
Inter-individual Level

- n_I researchers into Q_I clusters
- A researcher's cluster depends on his laboratory's cluster
- $Z_i^I = k \Leftrightarrow i \in k, k \in \{1, \dots, Q_I\}$
- Connectivity is independent given the latent variables

$$\mathbb{P}(Z_i^I = k | A_i = j, Z_j^O = l) = \gamma_{kl}$$

$$\mathbb{P}(X_{ij}^I = 1 | Z_i^I = k, Z_j^O = l) = \alpha_{kk'}^I$$

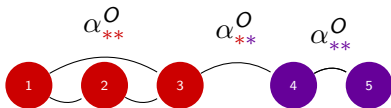
Independence Between Levels



- π^O is a probability vector
- Each column of γ as well
- If $\gamma_{kl} = \gamma_{k'l'} \quad \forall l, l'$

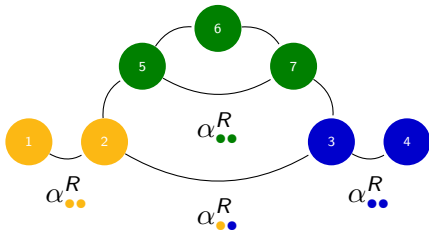
$$\mathcal{L}(X^I, X^O|A) = \mathcal{L}(X^I)\mathcal{L}(X^O)$$

Independence Between Levels



- π^O is a probability vector
- Each column of γ as well
- If $\gamma_{kl} = \gamma_{k'l'} \quad \forall l, l'$

$$\mathcal{L}(X^I, X^O|A) = \mathcal{L}(X^I)\mathcal{L}(X^O)$$



- Each level of the multilevel network is a SBM with $\pi^I = \gamma_{\cdot 1}$
- Organisational structure has no influence on the connectivity of individuals

Identifiability

Proposition

The multilevel model is identifiable up to label switching under the following assumptions:

- (i) All coefficients of $\alpha^O \cdot \pi^O$ are distinct
- (ii) All coefficients of $\alpha^I \cdot \gamma \cdot \pi^O$ are distinct
- (iii) $n_I \geq 2Q_I$
- (iv) $n_O \geq \max\{2Q_O, Q_I + Q_O - 1\}$
- (v) At least $2Q_I$ organizations contain one individual or more.

Outline

- 1 Modelling
- 2 Inference**
- 3 Model Selection
- 4 Simulation Studies
- 5 Application to Real Dataset

Maximum Likelihood Inference

Objective Joint clustering of $\mathbf{Z} = \{Z^I, Z^O\}$ and estimates of $\theta = \{\pi^O, \gamma, \alpha^O, \alpha^I\}$

Method Maximum likelihood of the observed data

Idea Calculate the complete likelihood and integrate on the latent variables

Problem Intractable, sum of $Q_R^{n_R} \times Q_L^{n_L}$ terms

Solution EM algorithm

Problem $\mathcal{L}(\mathbf{Z}|\mathbf{X})$ also intractable

Solution Variational approach of the EM algorithm

Variational EM

Maximise a lower bound of the observed data likelihood

$$\begin{aligned} \ell_{\theta}(\mathbf{X}) &\geq \ell_{\theta}(\mathbf{X}) - KL(\mathcal{R}(\mathbf{Z}) \parallel \mathbb{P}_{\theta}(\mathbf{Z}|\mathbf{X})) \\ &= \mathbb{E}_{\mathcal{R}}[\ell_{\theta}(\mathbf{X}, \mathbf{Z})] + \mathcal{H}(\mathcal{R}(\mathbf{Z})) \\ &= \mathcal{I}_{\theta}(\mathcal{R}(\mathbf{Z})) \end{aligned}$$

$\mathcal{R}(\mathbf{Z})$ is a mean-field approximation of $\mathbf{Z}|\mathbf{X}$

\mathcal{H} is the entropy

VEM algorithm

2 steps iterative algorithm

VE Maximise $\mathcal{I}_{\theta}(\mathcal{R}(\mathbf{Z}))$ w.r.t. $\mathcal{R}(\mathbf{Z})$

M Maximise $\mathcal{I}_{\theta}(\mathcal{R}(\mathbf{Z}))$ w.r.t. θ

Parameters update

VE-Step : variational parameters

$$\widehat{\tau}_{jl}^O \propto \pi_l^O \prod_{i,k} \gamma_{kl}^{A_{ij} \widehat{\tau}_{ik}^I} \prod_{j' \neq j} \prod_{l'} \varphi(X_{jj'}^O, \alpha_{ll'}^O, \widehat{\tau}_{j'l'}^O)$$

$$\widehat{\tau}_{jl}^I \propto \prod_{j,l} \gamma_{kl}^{A_{ij} \widehat{\tau}_{jl}^O} \prod_{i' \neq i} \prod_{k'} \varphi(X_{i'i'}^I, \alpha_{kk'}^I, \widehat{\tau}_{i'k'}^I)$$

$$\tau_{ik}^I = \mathbb{P}_{\mathcal{R}}(Z_i^I = k) \quad \tau_{jl}^O = \mathbb{P}_{\mathcal{R}}(Z_j^O = l)$$

$$\varphi(X, \alpha, \tau) = (\alpha^X (1 - \alpha)^{1-X})^\tau$$

M-step : model parameters

$$\widehat{\pi}_l^O = \frac{1}{n_O} \sum_j \widehat{\tau}_{jl}^O$$

$$\widehat{\alpha}_{kk'}^I = \frac{\sum_{i' \neq i} \widehat{\tau}_{ik}^I \widehat{\tau}_{i'k'}^I X_{i'i'}^I}{\sum_{i' \neq i} \widehat{\tau}_{ik}^I \widehat{\tau}_{i'k'}^I}$$

$$\widehat{\alpha}_{ll'}^O = \frac{\sum_{j' \neq j} \widehat{\tau}_{jl}^O \widehat{\tau}_{j'l'}^O X_{jj'}^O}{\sum_{j' \neq j} \widehat{\tau}_{jl}^O \widehat{\tau}_{j'l'}^O}$$

$$\widehat{\gamma}_{kl} = \frac{\sum_{i,j} A_{ij} \widehat{\tau}_{ik}^I \widehat{\tau}_{jl}^O}{\sum_{i,j} A_{ij} \widehat{\tau}_{jl}^O}$$

Outline

- 1 Modelling
- 2 Inference
- 3 Model Selection**
- 4 Simulation Studies
- 5 Application to Real Dataset

Model Selection for the number of clusters

Penalized criterion for choosing the number of clusters

$$\begin{aligned}
 ICL_{Multilevel}(Q_I, Q_O) = \max_{\theta} \ell_{\theta}(X^I, X^O, \hat{Z}^I, \hat{Z}^O | A) \\
 \underbrace{- \frac{1}{2} \frac{Q_I(Q_I + 1)}{2} \log \frac{n_I(n_I - 1)}{2}}_{\alpha^I} - \underbrace{\frac{Q_O(Q_I - 1)}{2} \log n_I}_{\gamma} \\
 \underbrace{- \frac{1}{2} \frac{Q_O(Q_O + 1)}{2} \log \frac{n_O(n_O - 1)}{2}}_{\alpha^O} - \underbrace{\frac{Q_O - 1}{2} \log n_O}_{\pi^O}
 \end{aligned}$$

- Step-wise procedure with relevant local initialization of VEM to optimise the ICL

Model selection for independence

- ICL can be used to state on the independence between levels
- New penalty term for γ

$$\text{pen}_\gamma = \frac{Q_I - 1}{2} \log n_I$$

- $ICL_{ind}(Q_I, Q_O) = ICL_{SBM}^I(Q_I) + ICL_{SBM}^O(Q_O)$
- We decide that levels are interdependent if

$$ICL_{ind}(Q_I, Q_O) < ICL_{Multilevel}(Q_I, Q_O)$$

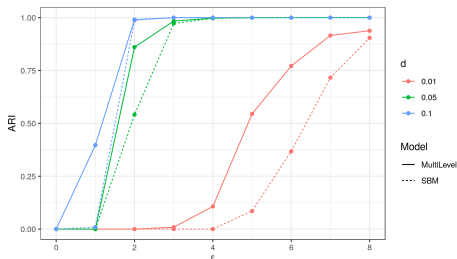
Outline

- 1 Modelling
- 2 Inference
- 3 Model Selection
- 4 Simulation Studies**
- 5 Application to Real Dataset

Simulation Studies

$$Q_O = 3 \quad Q_I = 3 \quad n_O = 20 * Q_I \quad n_I = 3 * n_O \quad \pi^O = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

$$\alpha^I = d * \begin{bmatrix} 1 + \epsilon & 1 & 1 \\ 1 & 1 + \epsilon & 1 \\ 1 & 1 & 1 + \epsilon \end{bmatrix} \quad \alpha^O = \begin{bmatrix} .5 & .1 & .1 \\ .1 & .5 & .1 \\ .1 & .1 & .5 \end{bmatrix} \quad \gamma = \begin{bmatrix} .8 & .1 & .1 \\ .1 & .8 & .1 \\ .1 & .1 & .8 \end{bmatrix}$$

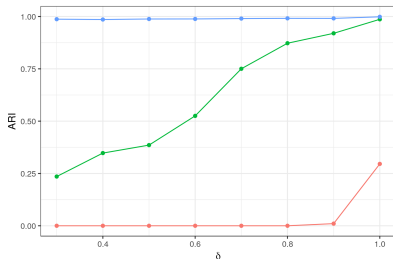
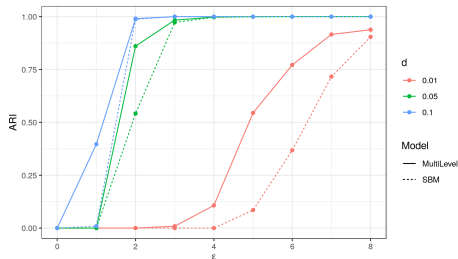


50 simulations on each dot

Simulation Studies

$$Q_O = 3 \quad Q_I = 3 \quad n_O = 20 * Q_I \quad n_I = 3 * n_O \quad \pi^O = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

$$\alpha^I = d * \begin{bmatrix} 1+2 & 1 & 1 \\ 1 & 1+2 & 1 \\ 1 & 1 & 1+2 \end{bmatrix} \alpha^O = \begin{bmatrix} .5 & .1 & .1 \\ .1 & .5 & .1 \\ .1 & .1 & .5 \end{bmatrix} \quad \gamma = \begin{bmatrix} \delta & 1 - \frac{\delta}{2} & 1 - \frac{\delta}{2} \\ 1 - \frac{\delta}{2} & \delta & 1 - \frac{\delta}{2} \\ 1 - \frac{\delta}{2} & 1 - \frac{\delta}{2} & \delta \end{bmatrix}$$



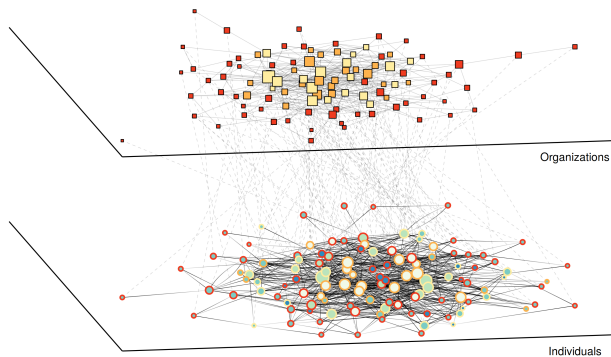
50 simulations on each dot

Outline

- 1 Modelling
- 2 Inference
- 3 Model Selection
- 4 Simulation Studies
- 5 Application to Real Dataset**

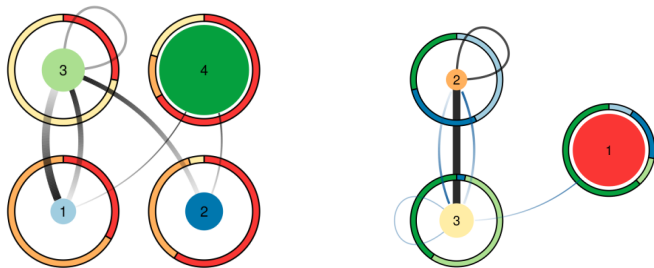
Application to a Television Program Trade Fair Dataset⁴

128 individuals with directed interactions (advice) and 109 organizations with undirected interactions (deal).



⁴Brailly, 2016

Dataset analysis



- 4 blocks of individuals and 3 blocks of organizations
- Levels are interdependent
- The structure of connection between individuals do not replicate the structure of connections between organizations

- Preprint available on arXiv: <https://arxiv.org/abs/1910.10512>
- R package available at <https://chabert-liddell.github.io/MLVSBM/>
 - Simulates and infers multilevel networks
 - Includes handling of missing data on X^I and X^O
 - Prediction on missing dyads, missing links and spurious links
 - Works with multi-affiliation datasets

References

- Biernacki, Christophe, Gilles Celeux, and Gérard Govaert (2000). "Assessing a mixture model for clustering with the integrated completed likelihood". In: *IEEE transactions on pattern analysis and machine intelligence* 22.7, pp. 719–725.
- Brailly, Julien (2016). "Dynamics of networks in trade fairs—A multilevel relational approach to the cooperation among competitors". In: *Journal of Economic Geography* 16.6, pp. 1279–1301.
- Celisse, Alain, Jean-Jacques Daudin, Laurent Pierre, et al. (2012). "Consistency of maximum-likelihood and variational estimators in the stochastic block model". In: *Electronic Journal of Statistics* 6, pp. 1847–1899.
- Daudin, J-J, Franck Picard, and Stéphane Robin (2008). "A mixture model for random graphs". In: *Statistics and computing* 18.2, pp. 173–183.
- Snijders, Tom A.B. and Krzysztof Nowicki (Jan. 1997). "Estimation and Prediction for Stochastic Blockmodels for Graphs with Latent Block Structure". In: *Journal of Classification* 14.1, pp. 75–100. ISSN: