

Statistical Learning with Complex Data

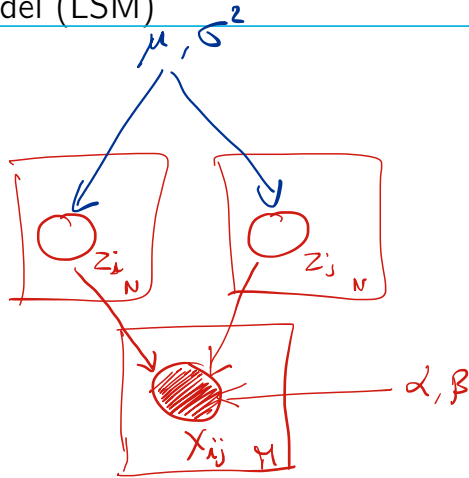


Pr. Charles BOUVEYRON

Professor of Statistics
Chair of the Institut 3IA Côte d'Azur
Université Côte d'Azur & Inria

✉ charles.bouveyron@univ-cotedazur.fr
🐦 @cbouveyron

The latent space model (LSM)



N : nb of nodes

M : — edges

\Rightarrow More efficient inference techniques relies on Bayesian version of the model $\Rightarrow z_i \sim N(\mu, \sigma^2 \mathbf{I}_r)$

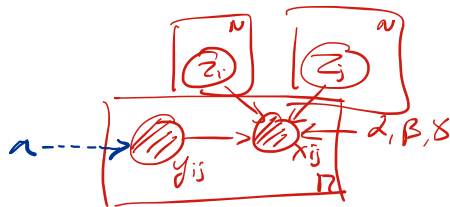
The latent space model (LSM)

Adding covariates:

$$\text{logit}(P(X_{ij}=1|G)) = \alpha - \beta \|z_i + z_j\|^2 + \gamma y_{ij}$$

where y_{ij} may indicate a difference between the nodes i and j

Ex: $y_{ij} = |\text{Age}_i - \text{Age}_j|$



Choice of the distance:

- Euclidean distance
- directional distance

$$\begin{aligned} d(z_i, z_j) &= \|z_i\| \cos(\hat{z_i z_j}) \\ &= \frac{z_i^T z_j}{\|z_j\|} \end{aligned}$$

More adapted to directed network with a strong asymmetry -

The latent space model (LSM)

A more tenable inference approach:

will rely on a Bayesian modeling of the LSM \Rightarrow MCMC procedure may be used to infer the model.

\Rightarrow "latentnet" package implements a MCMC algorithm for the LSM model.

Outline

1. Introduction
2. Characterization and manipulation of networks
3. The visualization of networks
4. Clustering of networks
5. Texts
6. Images

The clustering of networks

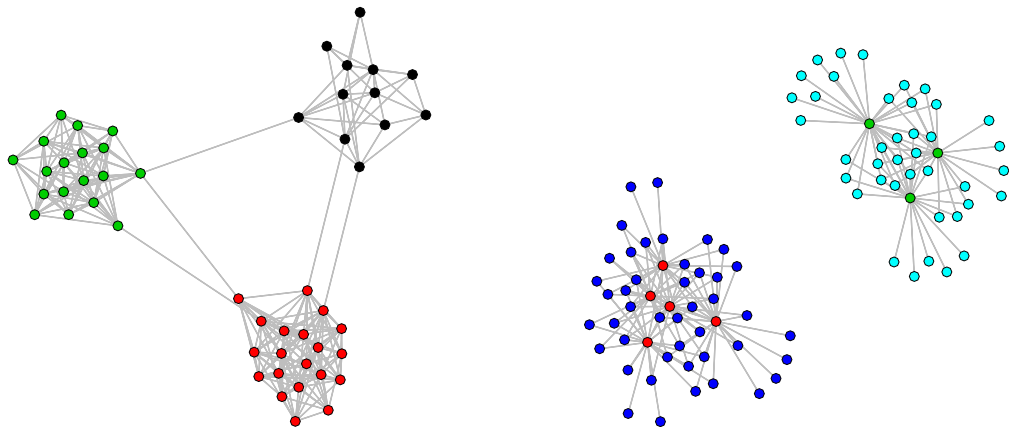
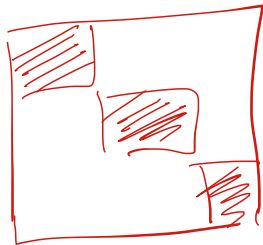


Figure: Clustering of communities vs. stars.

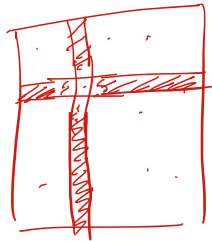
The clustering of networks

Difference between communities and stars:

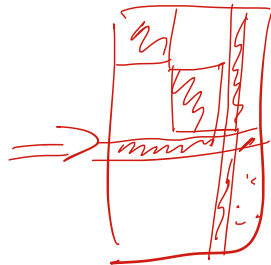
- Community : $P(X_{ij}=1 | C_i=C_j) \geq P(X_{ij}=1 | C_i \neq C_j)$
↳ assortative mixing
- Stars : $P(X_{ij}=1 | C_i \neq C_j) \geq P(X_{ij}=1 | C_i=C_j)$
↳ disassortative mixing



Communities



Stars



The latent position cluster model (LPCM)

The LPCM extends LSM by adding a clustering structure: (Hoff, Raftery, & Handcock, 2007)

$$\text{logit}(P(X_{ij} = 1 | \theta)) = \alpha - d(z_i, z_j) (+ \gamma y_{ij})$$

$C_i \sim \mathcal{C}(\tau)$, where C_i indicates the group membership of i .

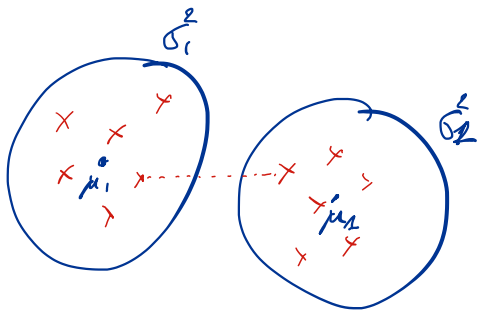
$C_i = (0, 0, 1, 0, 0) \Leftrightarrow i$ belongs to cluster #3.

$$Z_i | C_{ik} = 1 \sim N(\mu_k, \sigma_k^2 I_p)$$

$$\Leftrightarrow Z_i \sim \sum_{k=1}^K \pi_k N(\mu_k, \sigma_k^2 I_p)$$

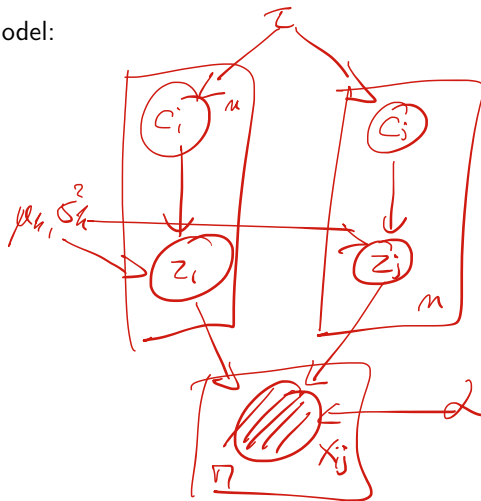
The latent position cluster model (LPCM)

The model:



The latent position cluster model (LPCM)

The graphical model:



The latent position cluster model (LPCM)

Inference:

- $\mathcal{M}CPC$ procedure \rightarrow latentnet package
- $\mathcal{V}BET$ procedure \rightarrow $\mathcal{V}B\mathcal{L}PCL$ —

⚠ $\mathcal{M}CPC$ is not adapted to large networks
(< 500 nodes)

\Rightarrow $\mathcal{V}BET$ may be used for networks smaller
than 20 000 nodes.

The latent position cluster model (LPCM)

In R, the *latentnet* and *VBLPCM* packages allow to use it:

The latent position cluster model (LPCM)

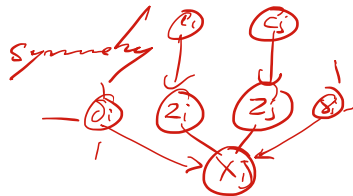
Extension #1: adding a sender/receiver effect

$$\text{Corr}(P(X_{ij} = 1 | \theta)) = \alpha - d(z_i, z_j) + \delta_i + \delta_j$$

$$\begin{cases} \delta_i \sim N(\mu_\delta, \sigma_\delta^2) \\ \delta_j \sim N(\mu_\delta, \sigma_\delta^2) \end{cases}$$

receiver effect
sender effect \approx prior probability for i to send a message to anyone.

\Rightarrow This model allows to break the symmetry of the LPCM / LSR model



The latent position cluster model (LPCM)

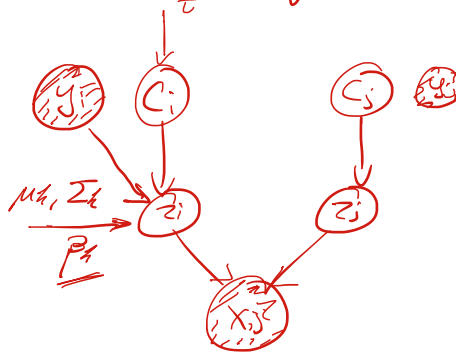
Extension #2: mixture of experts LPCM

incorporates covariate information

within the clustering model.

$$z_i \sim \sum_{k=1}^K \pi_k(y_i) \phi(z; \mu_k, \Sigma_k)$$

$$\text{where } \pi_k(y_i) = \frac{\exp(\beta_k^T y_i)}{\sum_{\ell=1}^K \exp(\beta_\ell^T y_i)}$$



\Rightarrow EM / VBEM procedures
to solve this problem.

The latent position cluster model (LPCM)

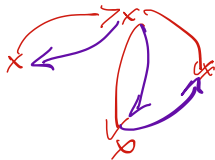
Extension #3: taking into account a dynamic

$$\left(X_{ij}^{(t)} \right)_{t=0 \dots T}$$

$$\longrightarrow \pi_k^{(t)}$$

$$\longrightarrow \pi_k^{(w_s)} \text{ where } w_s \text{ is a cluster of times.}$$

Extension #4: dealing with multi-networks



\rightarrow LPP 2 net package which implements an extension of the LPCM model to this situation.

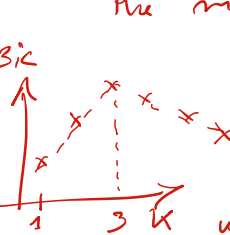
The latent position cluster model (LPCM)

Extension #1: adding a sender/receiver effect

The latent position cluster model (LPCM)

Choosing the number of clusters: Model selection is a theory which is based on penalized-likelihood criteria.

The most famous criterion is BIC:


$$Bic(\alpha) = \log \mathcal{L}(\hat{\theta}^n) - \frac{\gamma(\alpha)}{2} \log(n)$$

where $\gamma(\alpha)$ is the nb of free parameters in the model.

in the network context, the nb of edges

For instance: $LPCM = 1\alpha + K\mu + K\sigma^2 + (K-1)\pi \Rightarrow 3K$

The stochastic block model (SBM)

The SBM model ~~allows~~ allows to deal with:

- networks which are directed or not,
- — that contain communities but also stars.

It assumes:

- C_i indicates the group membership of node i

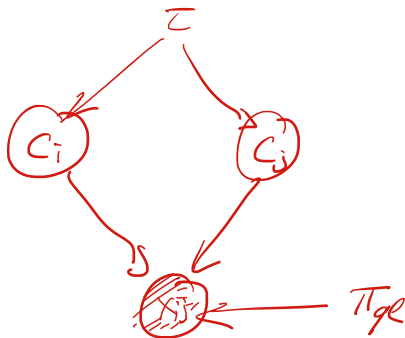
$$C_i \sim \mathcal{H}(\overline{\mathbf{C}}) \rightarrow C_i = (0, 0, 1, 0)$$

- the connection between i and j :

$$X_{ij} | C_{ig} C_{jl} = 1 \sim \mathcal{B}(\pi_{gl})$$

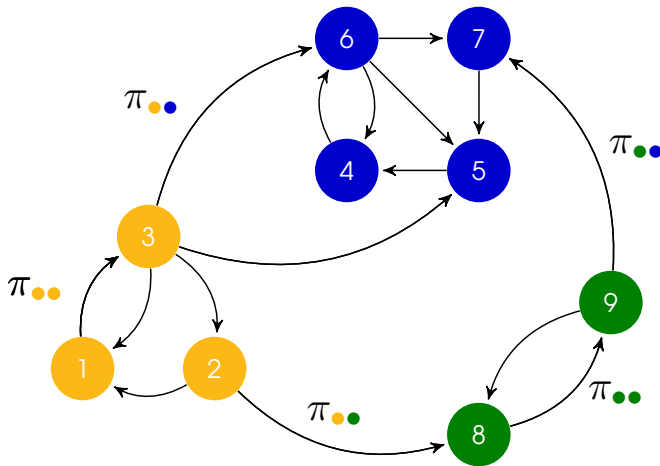
The stochastic block model (SBM)

The graphical model:



The stochastic block model (SBM)

A simple example:



The stochastic block model (SBM)

Inference: MCMC or V(B)EM procedures
are used for fitting a SBM model
from data.

Rank:

if cluster # q is a hub, what are the values
 π_{ql} , $\forall l = 1 \dots K$?

$$\pi_{ql} \geq \pi_{qq} \quad \forall l \neq q$$

(and π_{eq})

The stochastic block model (SBM)

Choosing the number of clusters:

→ BIC

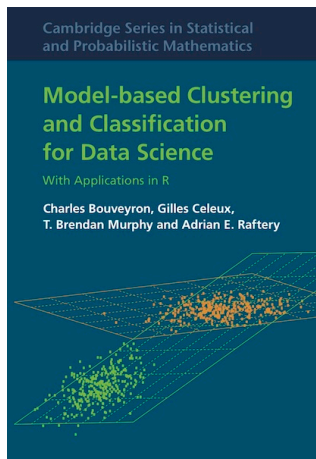
→ ICC

References



M. Salter-Townshend, A. White, I. Gollini and T. B. Murphy, *Review of Statistical Network Analysis: Models, Algorithms, and Software*, Statistical Analysis and Data Mining, Vol. 5(4), pp. 243–264, 2012.

References (more seriously ;-)



(Chapter 10 is devoted to network analysis!)