

Statistical Learning with Complex Data



Pr. Charles BOUVEYRON

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

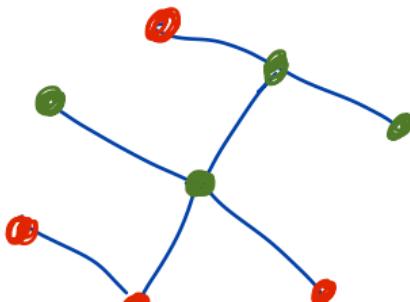
charles.bouveyron@univ-cotedazur.fr
 [@cbouveyron](https://twitter.com/cbouveyron)

The idea of clustering of networks could be:

Task 1)

group together the nodes of a network

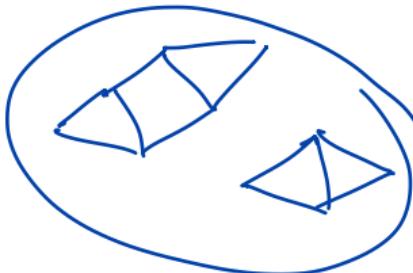
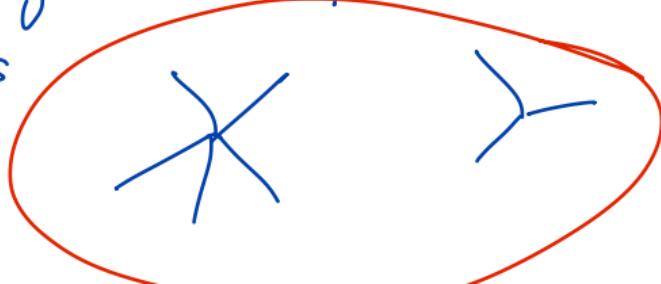
which have a similar role in the network



⇒ we will focus
on this problem

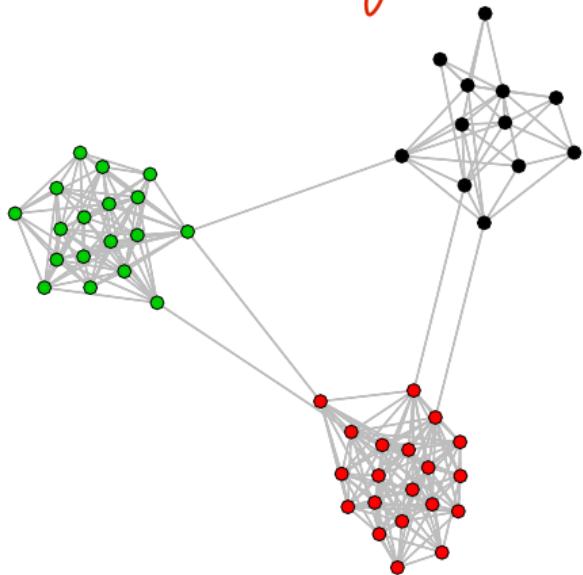
Task 2)

group together networks that have similar
structures

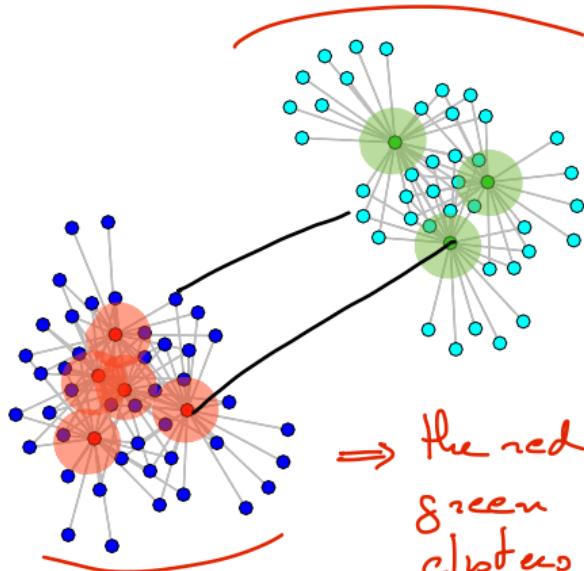


The clustering of networks

a Network made of communities



a network with stars.



→ the red and
green clusters contain

hubs

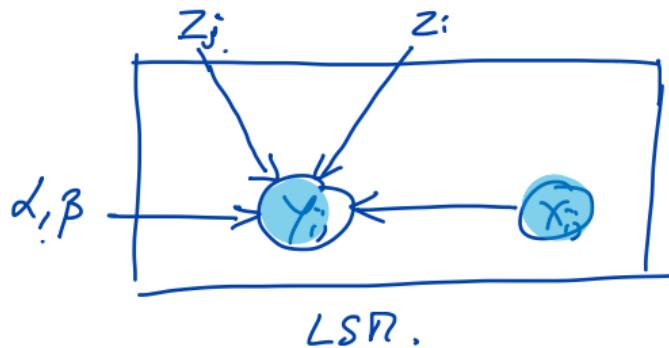
⇒ the two blue clusters contain
nodes that are "governed" by the hubs.

The latent position cluster model (LPCM)

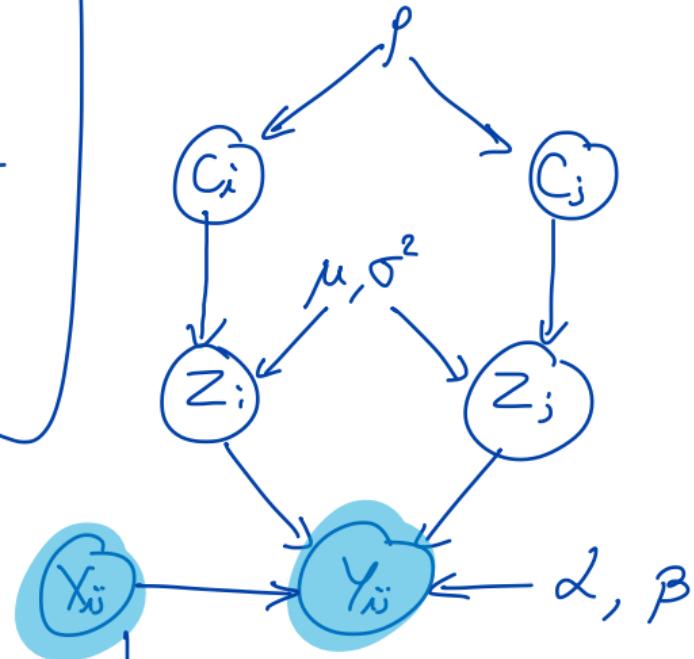
The LPCM extends LSM by adding a clustering structure:

$$\left\{ \begin{array}{l} \text{• } \text{Logit} \left(P(Y_{ij}=1|\theta) \right) = \alpha + \beta X_{ij} - d(z_i, z_j) \\ \text{• } c_i \sim \text{Cat}(1; \rho_h) \quad h=1, \dots, k \\ \text{• } z_i | c_{ih}=1 \sim N(\mu_h, \Sigma_h^2 I) \end{array} \right. \quad \left. \begin{array}{l} \text{This a Bayesian extension of LPR.} \\ \Rightarrow p(z_i, \theta) \sim \sum_{h=1}^k \rho_h N(\mu_h, \Sigma_h^2 I) \end{array} \right.$$

here the parameter $\rho = (\rho_1, \dots, \rho_k)$ is the vector of cluster proportions.



The inference of LPCR
should be done with
a NCPAC procedure or
a VBEN algorithm.



the graphical model for LPCR

Remark: the VBEN algo is working for large networks (VBLPCR package)

The latent position cluster model (LPCM)

Extension #1: adding a sender/receiver effect

Naturally, it is possible to add to LPCM the sender receiver effect.

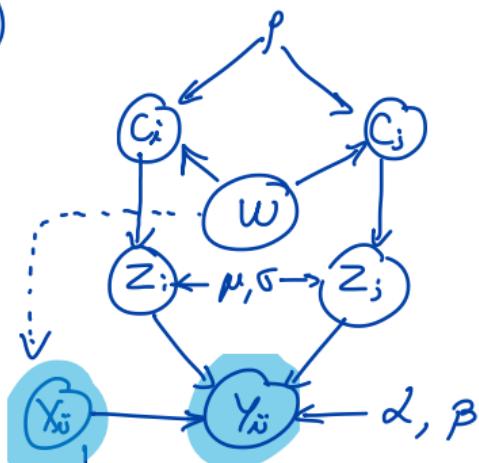
- $\text{Logit}(P(Y_{ij}=1|\theta)) = \alpha + \beta X_{ij} - d(z_i, z_j) + \delta_i + \gamma_j$
- $C_i \sim \mathcal{M}(1; \rho_h) \quad h=1, \dots, k$
- $z_i | C_{ih}=1 \sim N(\mu_h, \Sigma_h^2 I)$
- $\delta_i \sim N(0, \sigma_s^2)$ and $\gamma_j \sim N(0, \sigma_\gamma^2)$

The latent position cluster model (LPCM)

Extension #2: mixture of experts LPCM : this model assumes that the covariates W may have an effect on the clustering.

$$p(z_i) = \sum_{h=1}^k p_h(w_i) N(z_i; \mu_h, \Sigma_h^{-1})$$

Rank : this model may be useful in very specific situations and of course comes with difficulties for inference.



The stochastic block model (SBM)

The SBM model assumes:

$$C_i \sim \mathcal{M}(1, p)$$

$$\Leftrightarrow \text{cluster } \#3.$$

$$C_i = (0, 0, 1, 0, 0)$$

which indicates the cluster membership of the node i

$$Y_{ij} \mid C_{iq}, C_{je} = 1 \sim B(\pi_{qe}), q, e \in \{1, \dots, k\}$$

\uparrow $C_{iq} = 1$ means that
node i belongs to cluster q

$\pi_{qe} \in [0, 1]$ indicates
the probability to connect
for nodes coming from
cluster q and e
respectively.

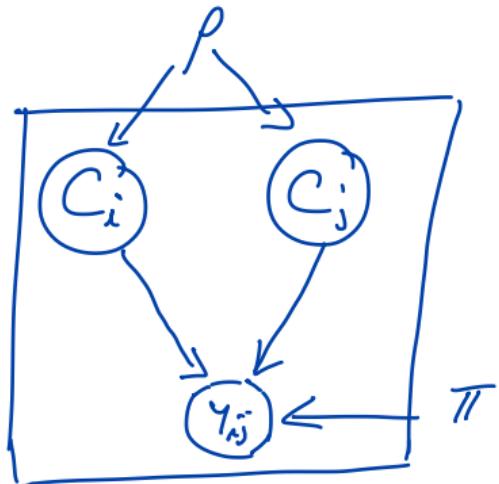
$X_{ij} = 1$
if $i \rightarrow j$ (in the directed case)

or if $i \leftrightarrow j$ (in the undirected case)

$X_{ij} = 0$ otherwise

The stochastic block model (SBM)

The graphical model:

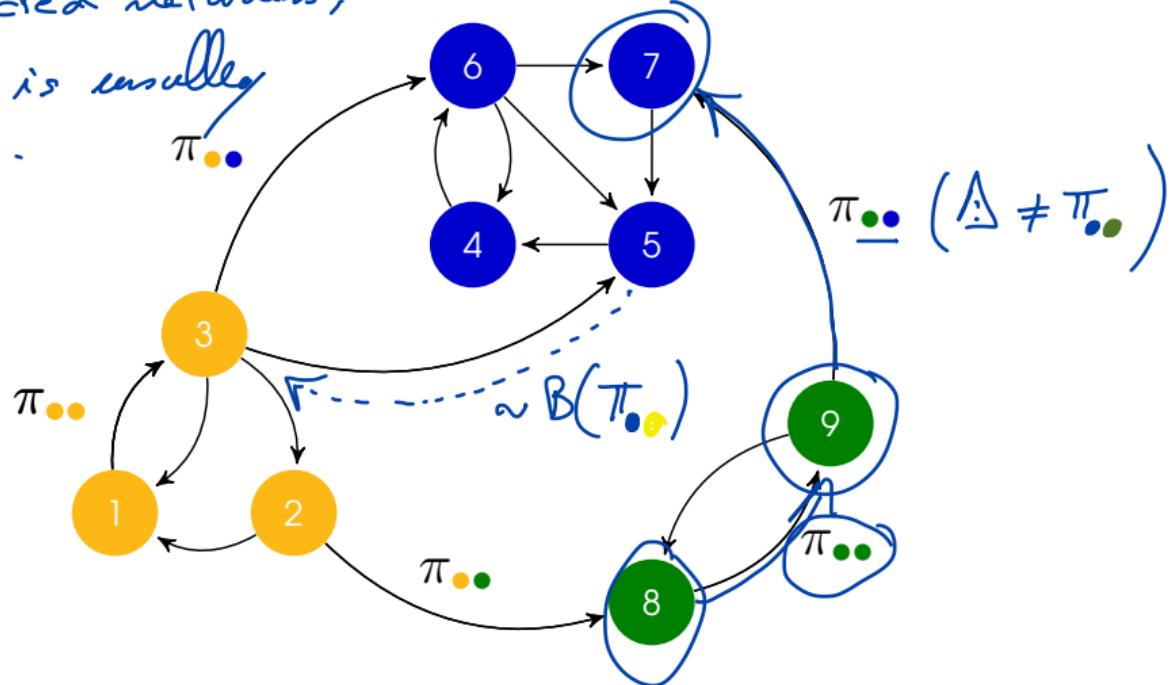


So, the inference of the SBM model will have to estimate from the data Y , the parameters p and π , plus the cluster memberships C .

The stochastic block model (SBM)

A simple example:

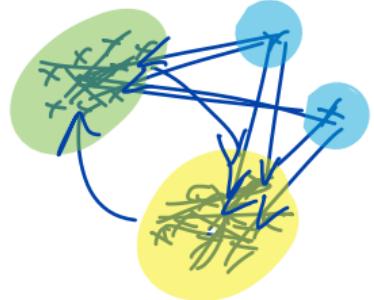
Rank: for directed networks,
the matrix Π is usually
not symmetric.



One of the most interesting feature of the SBM is that it can model:

- communities : $\pi_{qq} \geq \pi_{ql}$ for $l \neq q$
- stars/hubs : $\pi_{qq} \leq \pi_{ql}$ —

and a mixture of stars and communities:



$$\pi = \begin{pmatrix} 1 & 0.01 & 0.4 & 0.35 \\ 2 & 0.2 & 0.25 & 0.05 \\ 3 & 0.3 & 0.02 & 0.5 \end{pmatrix}$$

\rightarrow grp 1 is a group of hubs/influencers
 \rightarrow grp 2 and 3 are communities.

The stochastic block model (SBM)

Inference: SBP can be infer with :

- MCMC algorithms
- VEP algorithms

and other advanced techniques.

\Rightarrow VEP is implemented in the mixe package.

The stochastic block model (SBM)

Choosing the number of clusters: it can be choose thanks to
model selection criteria (AIC, BIC, ...)

↳ in practice for this kind of application, the
most popular technique is ICL
(integrated Classification likelihood)

$$ICL = BIC - \sum_i \sum_h \hat{c}_{ih} \log(\hat{c}_{ih})$$

