

Statistical Learning with Complex Data



Pr. Charles BOUVEYRON

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

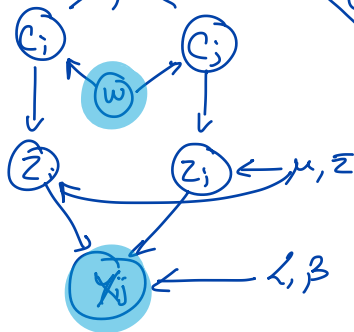
✉ charles.bouveyron@univ-cotedazur.fr
🐦 @cbouveyron

The latent position cluster model (LPCM)

Extension #2: mixture of experts LPCM

This model assumes that some covariates w may have an effect on the clustering

$$p(z_i) = \sum_{k=1}^K p_k(w_i) \mathcal{N}(z_i; \mu_k, \Sigma_k)$$



The prior probability of the groups depends on the individual covariate.

Note: This model of course comes with some complications regarding inference

The latent position cluster model (LPCM)

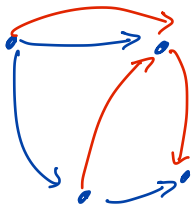
Extension #3: taking into account a dynamic

In order to model real-world networks where interactions may evolve along the time, it is interesting to model this. A way to do that is to assume that the cluster probabilities π evolve :

$$\text{logit}(\pi_k(t)) = \alpha_k(t) \sim \mathcal{N}(\alpha_k(t-1), \sigma^2)$$

Extension #4: dealing with multi-networks

this modeling is known as the State Space Model (SSM).



LSTM / LPCM to this situation $z_i \sim \sum \pi_k^c \mathcal{N}(\cdot)$

$$\text{logit}(P(X_{ij}^c | \theta)) = \alpha_c - \beta_c |z_i - z_j|$$

The Stochastic Block Model (SBM):

SBM is, at the moment, the most popular and efficient clustering model for networks. SBM has two main interests:

- 1) it is able to recover both communities and stars at the same time
- 2) the output of the model can be seen as a network summary (meta-network).

The stochastic block model (SBM)

The SBM model assumes:

- $C_i \sim \mathcal{H}(1; p)$

- $X_{ij} | C_{ik}=1, C_{je}=1$
 $\sim B(\pi_{ke})$

$C_{je}=1$
means that
 j belongs to
cluster e .

π_{ke} is the probability that people from cluster k connect with people from cluster e .

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

$\Rightarrow i$ belongs to cluster 3.

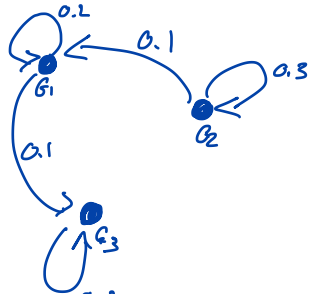
where $C_i = (C_{i1}, \dots, C_{iu})$
with $C_{ik} \in \{0, 1\}$.

and $p = (p_1, \dots, p_u)$
of the prior probabilities of
the groups.

Two remarks:

(i) The matrix Π can be seen as a network between the groups (it adjacency matrix, weighted)

$$\Pi = \begin{pmatrix} 0.2 & 0 & 0.1 \\ 0.1 & 0.3 & 0 \\ 0 & 0 & 0.3 \end{pmatrix}$$



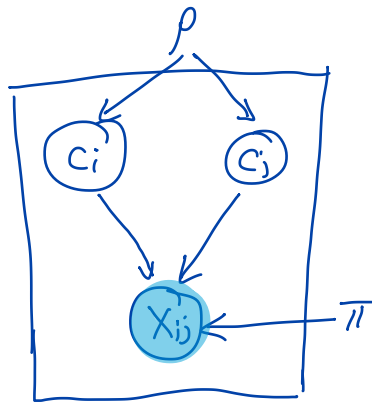
$$\Pi = \begin{pmatrix} 0.01 & 0.2 & 0 \\ 0.2 & 0.4 & 0 \\ 0.2 & 0 & 0.5 \end{pmatrix}$$

(ii) Π can also indicate if the groups are communities or stars



The stochastic block model (SBM)

The graphical model:



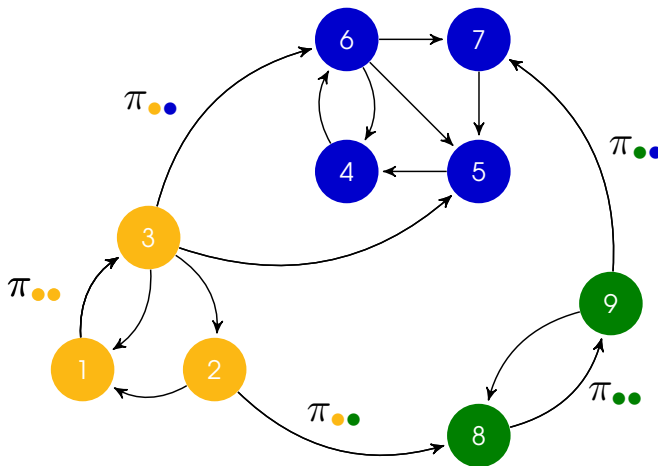
The inference of this model will have to estimate the model parameters, ρ and π , and the latent variables C .

→ Variational EM algorithm.

→ MCMC with a Bayesian version of the model.

The stochastic block model (SBM)

A simple example:



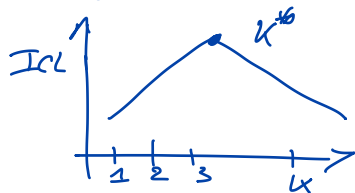
The stochastic block model (SBM)

Choosing the number of clusters:

As for other statistical models, we can rely here on model selection tools:

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

$$\text{ICL}(\mathcal{M}) = \text{Bic} - \sum_i \sum_k c_{ik} \log(c_{ik})$$



The mixed membership SBM (MMSBM)

The MMSBM extends the SBM as follows; in order to allow people to have different clusters depending on their roles in the network.

- $C_{i \rightarrow j} \sim \mathcal{O}(1; \rho_i)$ and $C_{i \leftarrow j} \sim \mathcal{O}(1; \rho_j)$
and $\rho_i \sim \text{Dir}(\alpha)$

- $X_{ij} \mid C_{i \rightarrow j}, C_{i \leftarrow j} \sim \mathcal{B}(\pi_{kl})$

