# 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)

$\text{logit}\left(P(X_{ij}|\Theta)\right) = \alpha + \beta Y_{ij} - d(z_i, z_j)$

The covariate $Y_{ij}$ can be used to provide extra information to the model on the pairs of nodes. For instance:

- $Y_{ij}$ is the nb of years in common in a club/society between $i$ and $j$.
- a type of relationship (categorial var) $\sim$

Choice of the distance: $Y_{ij} \in \{1, ..., K\} \implies Y_{ij} = (0, 0, 1, 0, 0) \implies \beta$ is a vector

$\uparrow Y_{ij} = 3$

Another way to extend this model is to play with the definition of the distance within the latent space

- $d(z_i, z_j) = \| z_i - z_j \|_2 \quad$ or $\| z_i - z_j \|_2^2$
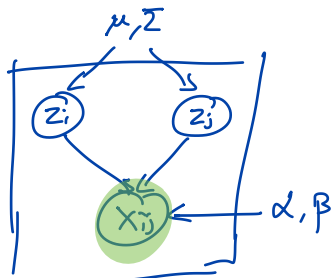- $d(z_i, z_j) = \| z_i - z_j \|_1 \quad$ (Manhattan distance)

34

**Modifying the model**: A specific and interesting case is the situation of directed networks, in which there are the roles of **sender** and **receiver**. It is naturally interesting to model this. A way to do that:

$$\text{logit}\left(P(X_{ij} = 1 \mid \theta)\right) = \alpha + \beta Y_{ij} - d(z_i, z_j) + \underline{\delta_i + \gamma_j}$$

a sender-receiver effect

$$\text{where} \begin{cases} d_i \sim N(0, \sigma_s^2) & \leftarrow \text{the prior for the propensity to send messages} \\ \gamma_j \sim N(0, \sigma_r^2) & \leftarrow \text{the prior for receiving messages.} \end{cases}$$

**Rmk**: this model is highly parametrized: it has $(3n + 2)$ parameters to estimate.

Exercise: draw the graphical model for this LSM version.

$\mu, \Sigma$

$z_i'$     $z_j'$

$X_{ij}'$     $\alpha, \beta$

Bayesian LSM.

$\mu, \Sigma$

$z_i'$     $z_j'$

$\sigma_s^2 \rightarrow$   $\delta_i'$     $X_{ij}'$     $\alpha, \beta$

$\sigma_r^2 \rightarrow$   $\gamma_j'$

# Outline
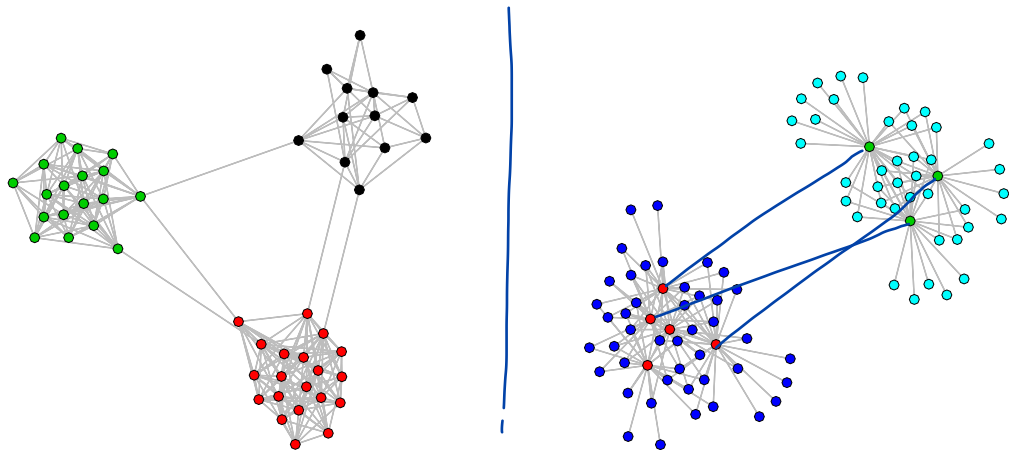
**40**

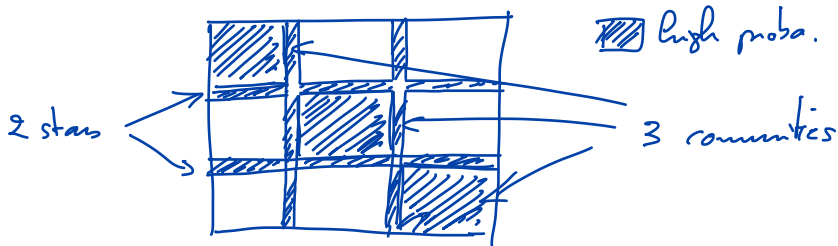# The clustering of networks



Figure: Clustering of communities vs. stars.

# The clustering of networks

Difference between communities and stars:

- in communities, people have a higher probability connection within the community than with other communities

- stars are people that connect less within the group than outside the group.
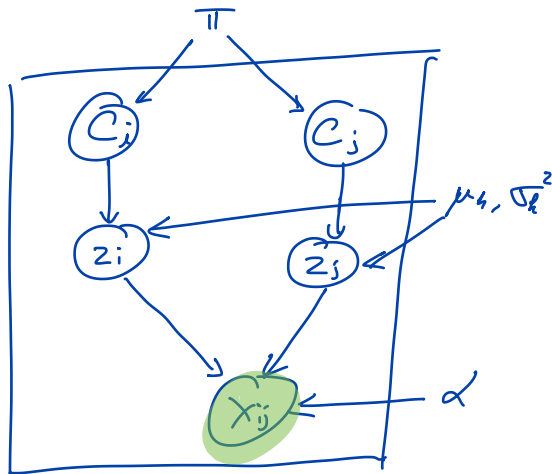


▨ high proba.

2 stars

3 communities

# The latent position cluster model (LPCM)

The LPCM extends LSM by adding a clustering structure:

- $\text{logit} \left( P(X_{ij} = 1 | \theta) = \alpha - d(z_i, z_j) \right.$

- $C_i \sim \mathcal{M}(1; \pi)$ where $\pi_h$ is the prior probability for cluster $k$, $h \in \{1, \dots, K\}$

- $Z_i \mid C_{ih} = 1 \sim N\left(\mu_h, \sigma_h^2 I\right)$

  $\Rightarrow \quad Z_i \sim \sum_{h=1}^{K} \pi_h \, N\left(\mu_h, \sigma_h^2 I\right).$

# The latent position cluster model (LPCM)

The model:



$\Rightarrow$ the inference of this Bayesian model has to be done using MCMC or advanced inference strategies (VBER)

44