

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 latent space model (LSM)

Adding covariates:

$$Y_{ij}$$

- a nb of years in common in a society, ...
- a type of relationship (categorical var)

$$Y_{ij} \in \{1, \dots, k\} \rightarrow \tilde{Y}_{ij} = (0, 0, 1, 0, 0) \Rightarrow \begin{matrix} \text{β is now} \\ \text{a vector.} \end{matrix}$$

\uparrow

$$Y_{ij} = 3.$$

Choice of the distance:

$$\|z_i - z_j\|^2$$

It could be an Euclidean distance $\| \cdot \|_2$ (the most natural), but it could be any other distance if you prefer ($\| \cdot \|_1, \dots$)

It is possible to compute afterward $P(Y_{ij}=1|\hat{\theta})$

$$\text{Logit } P(Y_{ij}=1|\theta) = \alpha - d(z_i, z_j)$$

$$\Leftrightarrow \log\left(\frac{p}{1-p}\right) = \alpha - d(z_i, z_j)$$

$$\Leftrightarrow p/(1-p) = \exp(\alpha - d(z_i, z_j))$$

$$\Leftrightarrow p = \exp(-) - p \exp(-)$$

$$\Leftrightarrow p(1 + \exp(-)) = \exp(-)$$

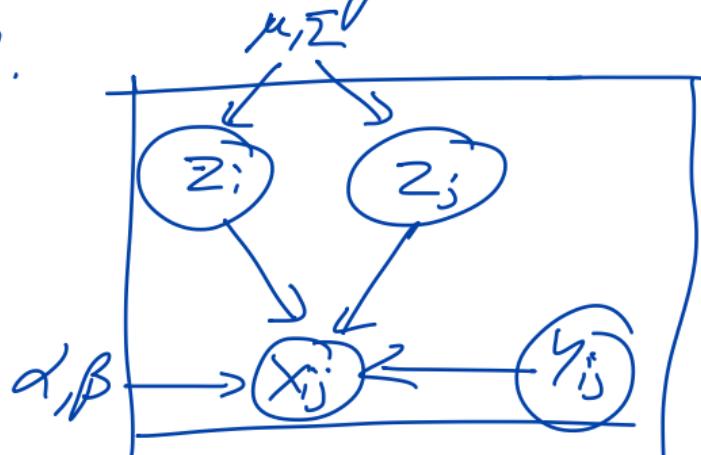
$$\Leftrightarrow P(Y_{ij}=1|\hat{\theta}) = \frac{\exp(\hat{\alpha} - d(\hat{z}_i, \hat{z}_j))}{1 + \exp(\hat{\alpha} - d(\hat{z}_i, \hat{z}_j))}$$

The latent space model (LSM)

LSM in R via the *latentnet* package: with implements the original approach of Hoff et al (2001) \Rightarrow MCMC for a Bayesian version of LSM.

$$= \text{LSR} + z_i \sim N(\mu, \Sigma)$$

and the *VBLPCR* package implements a VBEM algorithm for this model.



It is possible to extend the LSR model :

- by adding covariates on the nodes:

$$\text{Clogit} \left(P(Y_{ij} = 1 | \theta) \right) = \alpha + \beta X_{ij} - d(z_i, z_j)$$

where X_{ij} is a covariate about the pair of nodes

$$\text{Ex: } X_{ij} = |\text{age}_i - \text{age}_j| \quad \text{or} \quad X_{ij} = \prod \{ \text{pos}_i = \text{pos}_j \}$$

- by "playing" with the notion of distance:

- $d(z_i, z_j) = \|z_i - z_j\|^2$ is ok for symmetric graphs.

- $d_i(z_i, z_j) = |z_i| \cos \theta_{ij}$ where $\theta_{ij} = \frac{z_i^T z_j}{|z_i| |z_j|}$
is better for asymmetric networks

- by modifying the model :

$$\text{logit} \left(P(Y_{ij}=1 | \theta) \right) = \alpha + \beta X_{ij} - d(z_i, z_j) + \underbrace{\delta_i + \gamma_j}_{\text{a sender-receiver effect}}$$

where $\begin{cases} \delta_i \sim N(0, \sigma_s^2) \\ \gamma_j \sim N(0, \sigma_r^2) \end{cases}$

\leftarrow propensity to send connections
 \leftarrow to receive connections.

\Rightarrow this is of course for directed networks $\begin{pmatrix} Y_{ij}=1 \\ \cancel{Y_{ji}=1} \end{pmatrix}$

Rank: this is a quite highly parameterized model
 \Rightarrow this model has $(3n+2)$ parameters to estimate
and requires a lot of connections to estimated properly!