

# Penalized Gaussian Copula mixtures for mixed mode data

Panagou Foteini, Karlis Dimitrios

Department of Statistics  
Athens University of Business & Economics

[fwtpanagou@aueb.gr](mailto:fwtpanagou@aueb.gr), [karlis@aueb.gr](mailto:karlis@aueb.gr)

Athens, 26 Oct 2021

## Mixtures of Gaussian copulas for mixed mode data

For mixed data sets the joint probability is not always easy to be found. Gaussian copula offers flexibility to describe the dependencies between different types of variables.

Multivariate Gaussian copula is defined as

$$C^N(u_1, u_2, \dots, u_p; R) = \Phi_p(\Psi(u_1), \dots, \Psi(u_p); R)$$

where:  $\Phi_p$  is the p.m.f. of a standard p-variate normal distribution,  $\Psi(\cdot) = \Phi^{-1}(\cdot)$  is the inverse p.d.f. of a standard univariate normal distribution.

The correlation matrix  $R$  of  $k$  continuous &  $\ell$  discrete variables can be split into blocks based on the type of variables:

$$R = \begin{bmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{bmatrix}$$

where e.g.  $R_{11}$  is the correlation matrix of the  $k$  continuous variables.

Through Choleski factor and angles reparametrization,  $R = LL'$ ,  $R$  and  $\Theta$  are correlated:

$$L_{11} = 1, L_{ii} = \prod_{u=1}^{i-1} \sin \theta_{ui}, L_{ij} = \cos \theta_{ij} \prod_{u=1}^{i-1} \sin \theta_{uj}, \theta_{ij} = \cos^{-1} \left\{ \frac{L_{ij}}{\prod_{u=1}^{i-1} \sin \theta_{uj}} \right\}, i < j$$

$\theta_{ij} = 0$ ,  $i \geq j$  and angles are measured in radians. We require  $\theta_{ij} \in (0, \pi]$  so that the Choleski factor is unique.

$$\begin{aligned} \ell &= \sum_{g=1}^G \log \{ f(\mathbf{x}; \gamma_g, \Psi) \} = \sum_{g=1}^G \sum_{t=1}^k \log(f_t(\mathbf{x}; \gamma_{tg})) - \sum_{g=1}^G \sum_{t=1}^k \log \left\{ \phi(\Phi^{-1}(u_{tg})) \right\} \\ &+ \sum_{g=1}^G \log \left\{ \phi_k \left( \Psi(u_1), \dots, \Psi(u_k); \mathbf{0}, \mathbf{R}_{11g} \right) \right\} \\ &+ \sum_{g=1}^G \log \left\{ \sum_d \operatorname{sgn}(\mathbf{d}) \Phi_\ell \left( \Psi(u_{k+1}), \dots, \Psi(u_{k+l}); \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g \right) \right\} + \boxed{\lambda \sum_{g=1}^G \sum_{j=1}^{p(p-1)/2} \sin^2 \theta_{gj}} \end{aligned}$$

$\boldsymbol{\mu}_g = \mathbf{R}_{21g} \mathbf{R}_{11g}^{-1} [\Psi(u_1), \dots, \Psi(u_k)]'$ ,  $\boldsymbol{\Sigma}_g = \mathbf{R}_{22g} - \mathbf{R}_{21g} \mathbf{R}_{11g}^{-1} \mathbf{R}_{12g}$ ,  $u_i = F_i(x_i; \gamma_{gi})$ ,  $p = k + \ell$  and  $\lambda$  takes values into a grid ( $\lambda > 0$ ).

For  $\lambda = 0$  the approach is equivalent to the fully parametrized model without any constrain to the correlation matrix.

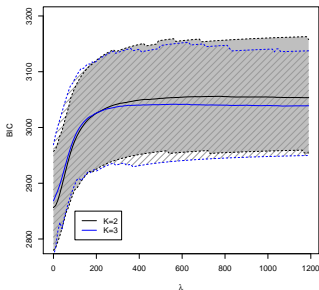
For  $\lambda \rightarrow \infty$  the log-likelihood shrinks to the independent model where all  $\theta$ 's are equal to  $\pi/2$ .

# Simulation Study

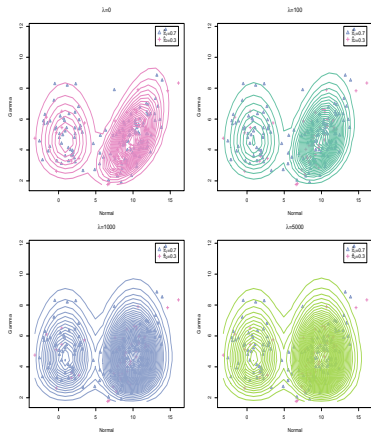
Assume a 4-variate mixed mode dataset (Normal, Gamma, Poisson & Bernoulli) of length  $n = 200$ ,  $G = 2$  components ( $\pi_1 = 0.7$ ).

Model	K=2	K=3
$\lambda = 0$	<b>2856.46</b>	2868.34
$\lambda = 10$	<b>2859.94</b>	2877.58
$\lambda = 50$	<b>2911.58</b>	2924.45
$\lambda = 100$	<b>2976.66</b>	2983.74
$\lambda = 500$	3052.28	<b>3040.95</b>
$\lambda = 1000$	3054.61	<b>3039.28</b>
$\lambda = 5000$	3043.43	<b>3028.10</b>

BIC vs  $\lambda$



BIC for various values of  $\lambda$



Contours for various values of  $\lambda$