

UBC Spatial Stats Course VI

Jürgen Pilz

Institut für Statistik
Universität Klagenfurt
Universitätsstr. 65-67, 9020 Klagenfurt, Austria
juergen.pilz@uni-klu.ac.at

November 30, 2010 / UBC Vancouver

2nd Extension to Non-Gaussianity: GLGM

- Diggle and Ribeiro (2007):

GLGM = Generalized Linear Geostatistical Model

$$g(\mu_i) = \mathbf{c}_i^T \boldsymbol{\beta} + S(x_i), \text{ where}$$

$S(x_i)$ = random spatial effects

- R-Software: library(**geoR**), library(**geoRglm**)
- Bayesian Hierarchical GLM framework in Banerjee et al. (2004)
BHGLM

Going beyond GLM and BHGLM

- How to proceed with distributions not covered by the GLM-framework?
- How to deal with extreme events coming from a different random process as opposed to the „background“ process?

Modelling of heavy-tailed/**extreme-value** distributions:
These are distributions which have a much slower decay of probability in the tails

Remember: Exponential prob. decay for normal distribution

GEV: $Z = Z(x)$ has c.d.f.

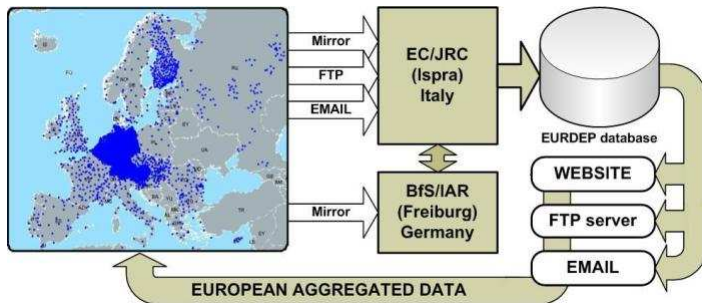
$$F(z; \mu, \sigma, \tau) = \exp \left(- \left[1 - \tau \left(\frac{z - \mu}{\sigma} \right) \right]_+^{-1/\tau} \right)$$

Main challenges in the EU-Project INTAMAP

- Development of new methodology for Bayesian spatial interpolation for non-Gaussian observations, with specific attention to extreme-value distributions
- Development of new methods for characterizing spatial dependence structures (based on copulas)
- Implementation (existing + new methods) in case of non-Gaussian, skewed and heavy-tailed observations:

library(intamap)

Library of R-functions



Briefly on Copulas

Copulas: distribution functions on the unit cube $[0, 1]^n$ with uniformly distributed margins, introduced by Abe Sklar.

Sklar's Theorem: Let H be an n -dimensional distribution function with margins F_1, \dots, F_n . Then there exists an n -dimensional copula C such that for all $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$

$$H(x_1, \dots, x_n) = C(F_1(x_1), \dots, F_n(x_n)).$$

If F_1, \dots, F_n are all continuous, then C is unique. Conversely, if C is an n -copula and F_1, \dots, F_n are d. f.s, then H is an n -dimensional d. f. with margins F_1, \dots, F_n , and

$$C(u_1, \dots, u_n) = H\left(F_1^{-1}(u_1), \dots, F_n^{-1}(u_n)\right).$$

Main Properties of Copulas

- Copulas describe the dependence between the quantiles of random variables
- Copulas are **invariant** under strictly increasing transformations of the marginals: Thus, frequently applied data transformations (e.g. square root and log transformations) do not change the copula.
- Random variables X_1, \dots, X_n are stochastically independent if and only if their copula is the product copula

$$\Pi^n(\mathbf{u}) = \prod_{i=1}^n u_i.$$

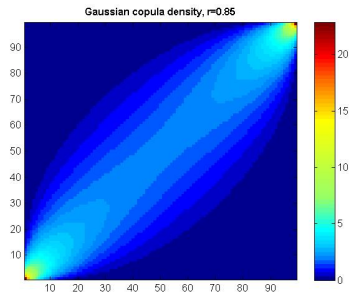
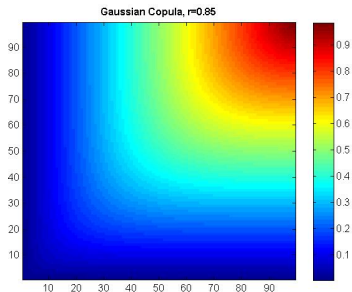
- Furthermore: (conditional) copula densities available

The Gaussian Copula

Sklar's Theorem provides a simple way of constructing copulas from multivariate distributions: Suppose $H = \Phi_{\mathbf{0}, \Sigma}$ and $F_1 = \dots = F_n = \Phi$, then

$$C_{\Sigma}^G(u_1, \dots, u_n) = \Phi_{\mathbf{0}, \Sigma} \left(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_n) \right)$$

is called the Gaussian copula.



Spatial Modeling using Copulas

- Bardossy (2006) was the first to present a method for spatial modeling using copulas
- The relation between two locations separated by the vector \mathbf{h} is characterized by the bivariate distribution

$$P(Z(\mathbf{x}) \leq z_1, Z(\mathbf{x} + \mathbf{h}) \leq z_2) = C_{\mathbf{h}}(F_Z(z_1), F_Z(z_2))$$

The copula thus becomes a function of the separating vector \mathbf{h}

- Spatial copulas describe spatial dependence over the whole range of quantiles for a given separating vector \mathbf{h} , and not only the mean dependence as the variogram does.
- H. Kazianka made enormous progress in this direction

Incorporating Spatial Trend

- Natural extension to incorporate trend: parameterize F_Z
- If F_Z belongs to the class of exponential dispersion models, $ED(\mu, \sigma)$, the mean μ appears explicitly inside the analytical expression of the density f_Z and we can set

$$g(\mu_i) = \eta(\mathbf{c}_i) = \mathbf{c}_i^T \boldsymbol{\beta}, \quad i = 1, \dots, n$$

where g is the **link function**, $\eta(\mathbf{c}_i)$ is the response surface, \mathbf{c}_i is the vector of **covariates** corresponding to location \mathbf{x}_i and $\boldsymbol{\beta}$ is the vector of regression parameters.

- For other models, e.g. the GEV, parameterize the location parameter μ of the distribution in the same way.

Relation to Diggle's GLGM

- Copula-based spatial modeling approach is an alternative to the model-based approach of Diggle (1998, 2007)
- Their approach can be easily extended to include mixed linear effects:

$$g(\mu_i) = \mathbf{c}_i^T \boldsymbol{\beta} + S(\mathbf{x}_i), \quad i = 1, \dots, n$$

where $S(\cdot)$ is a stationary Gaussian process with mean zero, variance σ^2 and given correlation function $\rho(\cdot)$

- The main difference: in GLGM only μ_i is modeled while we build a complete multivariate distribution for the response variables $Z(\mathbf{x}_i); i = 1, \dots, n$

Parameter Estimation for Cont. Margins

- Parameter vector $\Theta = (\theta, \lambda, \eta)$, where
 θ = correlation function parameters, λ = copula parameters and
 η = parameters of the known fam. of univ. distributions F
- Likelihood of the data $Z = (Z(x_1), \dots, Z(x_n))^T$:

$$l(\Theta; Z) = c_{\theta, \lambda}(F_{\eta}(Z(x_1)), \dots, F_{\eta}(Z(x_n))) \prod_{i=1}^n f_{\eta}(Z(x_i))$$

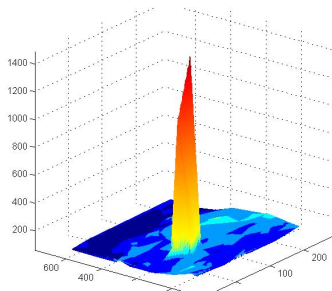
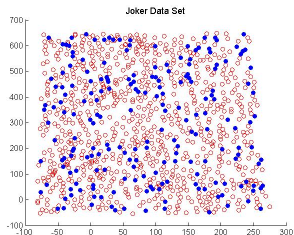
- For copulas different from the Gaussian c. the evaluation of the density is infeasible in higher dimensions. Here we proceed with **composite ml** using the bivariate copula densities, assuming that different pairs of observations can be treated as independent.

Bayesian Spatial Modeling

- Use a proper prior for Θ to assure that the posterior distribution, $p(\Theta | Z)$, is proper too. Otherwise, one has to prove the propriety of the posterior, which can be a difficult task
- For computation of the posterior d. of $\Theta = (\theta, \lambda, \eta)$, we use an MH-Algorithm
- Let $\hat{\Theta}^{(i)} = (\hat{\theta}^{(i)}, \hat{\lambda}^{(i)}, \hat{\eta}^{(i)})$ denote the i-th sample from the posterior. The approximation to the predictive density is

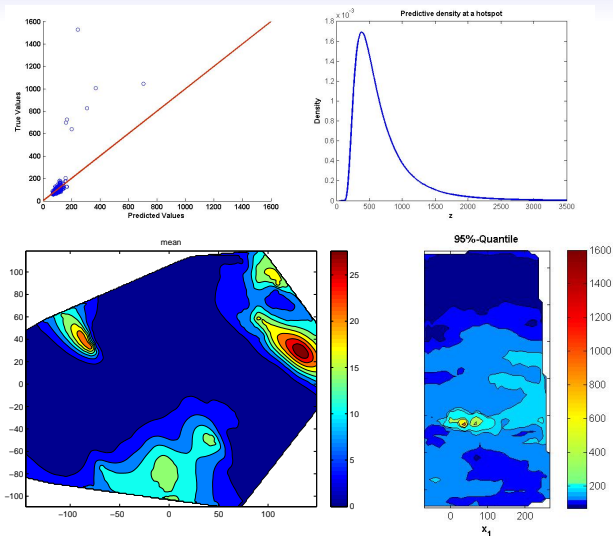
$$p(Z(x_0) | Z) \approx \frac{1}{N} \sum_{i=1}^N p(Z(x_0) | Z, \hat{\Theta}^{(i)})$$

SIC2004 Joker Data Set



- Extreme value data set from the Spatial Interpolation Contest 2004 (Dubois)
- 200 training data and 808 test data. The training data include two very large observations (1070.4 and 1499) and have a sample **skewness** of **9.92**
- We assume a **GEV** d. as the marginal d.
- The marginal location parameters are specified as
$$\mu_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i}.$$
- Fit a Gaussian copula with correlation function = mixture of a Gaussian and an exponential model. Geometric anisotropy is also considered.
- RMSE=64.47, MAE=15.83, ME=-2.56, Pearson r=0.72
Better results than with trans-Gaussian kriging using Box-Cox or Log-Log transformation and other classical geostatistical applications.

Spatial Interpolation of SIC2004 Joker Data



Conclusions

- Slight to moderate deviations from Gaussianity can be modelled with (Bayesian) Trans-Gaussian Kriging, using Box-cox-transf., GLM-framework of Diggle & Ribeiro (geoRglm) or Aston's **library(psgp)** or using Gaussian copulas (copula package)
- Heavy-tailed/Extreme value distributions we propose to model with function **copulaEstimation** in **library(intamap)**
- Code for estimation of parameters and prediction at unobserved locations in the case of the Gaussian spatial copula model is part of this library: **spatialPredict**
- Automatic choice of marginal distribution, correlation function model, anisotropy and starting values by using certain heuristics.
- Copula package not yet ready for **real**-time monitoring of extreme events
(it took 10 years to have current version of geoRglm!)

Bibliography

- Kazianka, H. & Pilz, J. (2010) Spatial interpolation using copula-based geostatistical models. In: GeoENV VII Geostatistics for Environmental Applications (P. Atkinson and C. Lloyd, eds), Springer, Berlin 2010, 307-319
- Kazianka, H. & Pilz, J. (2010) Copula-based geostatistical modeling of continuous and discrete data including covariates. Stochastic Environmental Research and Risk Assessment 24 (2010), 661–673
- <http://www.intamap.org>