# GeoModels Tutorial: simulation, estimation and prediction of bivariate spatial data using Gaussian random fields

Moreno Bevilacqua

## Introduction

In this tutorial we show how to analyze geo-referenced spatial bivariate data using Gaussian random fields (RFs) with the R package `GeoModels` (Bevilacqua and Morales-Oñate, 2018).

We first load the R libraries needed for the analysis and set the name of the model in the `GeoModels` package:

```
rm(list=ls())
require(devtools)
install_github("vmoprojs/GeoModels")
require(GeoModels)
require(fields)
model="Gaussian"  # model name in the  GeoModels package
set.seed(24)
```

## Simulation of a bivariate Gaussian random field

Let $\boldsymbol{Y}_{12} = \{\boldsymbol{Y}_{12}(\boldsymbol{s}) = (Y_1(\boldsymbol{s}), Y_2(\boldsymbol{s}))^T, \boldsymbol{s} \in A \subseteq \mathbb{R}^d\}$ be a standardized bivariate Gaussian random field where $Y_i = \{Y_i(\boldsymbol{s}), \boldsymbol{s} \in A \subseteq \mathbb{R}^d\}$ are two correlated univariate standard Gaussian random fields with $\mathbb{E}(Y_i(\boldsymbol{s})) = 0$, $\mathrm{var}(Y_i(\boldsymbol{s})) = 1$, $i = 1, 2$.

Under second order stationary assumption, the correlation function between $\mathbf{Y}_{12}(\boldsymbol{s}_l)$ and $\mathbf{Y}_{12}(\boldsymbol{s}_m)$, for any pair $\boldsymbol{s}_l, \boldsymbol{s}_m$ in the spatial domain, is represented by a mapping $\boldsymbol{R} : \mathbb{R}^d \to M_{2\times 2}$ defined through

$$\boldsymbol{R}(\boldsymbol{h}) = [R_{ij}(\boldsymbol{h})]^2_{i,j=1} = [\mathrm{cor}\,(Y_i(\boldsymbol{s}_l), Y_j(\boldsymbol{s}_m))]^2_{i,j=1}, \quad \boldsymbol{h} = \boldsymbol{s}_l - \boldsymbol{s}_m \in \mathbb{R}^d. \tag{1}$$

The function $\boldsymbol{R}(\boldsymbol{h})$ is called bivariate correlation function. Here, $M_{2\times 2}$ is the set of two dimensional squared, symmetric and positive definite matrices. The functions $R_{ii}(\boldsymbol{h})$ $i = 1, 2$ are the marginal correlation functions of the Gaussian random fields $Y_i$, $i = 1, 2$ while $R_{ij}(\boldsymbol{h})$ is called cross correlation function between $Y_i$ and $Y_j$ for $i, j = 1, 2$ and $i \neq j$ at spatial lag $\boldsymbol{h}$. The cross-covariance function is not in general symmetric, *i.e* $R_{12}(\boldsymbol{h}) \neq R_{21}(\boldsymbol{h})$ (Wackernagel, 2003). However, the majority of the existing multivariate parametric covariance models are symmetric, with some few exceptions (Genton and Kleiber, 2015). In this tutorial we assume $R_{12}(\boldsymbol{h}) = R_{21}(\boldsymbol{h})$.

Additionally, we assume that the cross-correlation function assume the following parametric form:

$$\boldsymbol{R}(\boldsymbol{h}) = [R_{ij}(\boldsymbol{h})]_{i,j=1}^2 = \left[\rho_{ij} K_{\psi_{ij}}(\boldsymbol{h})\right]_{i,j=1}^2, \quad \rho_{ii} = 1, \quad |\rho_{12}| < 1 \tag{2}$$

where $\rho_{12} = \rho_{21}$, is the so-called colocated correlation parameter and $K_{\psi_{ij}}(\boldsymbol{h})$ is a univariate parametric correlation model and $\psi_{ij}$ a parameter set. Other parametric models, as for instance the linear model of coregionalization are possible (Wackernagel, 2003), but in this tutorial we focus in the parametric model (2).

For each random field, we consider a location and scale transformation:

$$Z_i(\boldsymbol{s}) = \mu_i(\boldsymbol{s}) + \sigma_i Y_i(\boldsymbol{s}) + \tau_i \varepsilon_i(\boldsymbol{s}), \quad i = 1, 2.$$

Here $\{(\varepsilon_1(\boldsymbol{s}), \varepsilon_2(\boldsymbol{s}))^\top, \boldsymbol{s} \in A \subseteq \mathbb{R}^d\}$ is a bivariate Gaussian standard white noise such that $Cov(\varepsilon_1(\boldsymbol{s}), \varepsilon_2(\boldsymbol{s})) = \rho_{12}((\sigma_1^2 + \tau_1^2)^{\frac{1}{2}}(\sigma_2^2 + \tau_2^2)^{\frac{1}{2}} - \sigma_1 \sigma_2)$. In addition, $\sigma_i^2 > 0$, $i = 1, 2$ are the marginal variance parameters and $\tau_i^2 > 0$, $i = 1, 2$ are the marginal nugget parameters. Under these assumptions $\mathbb{E}(Z_i(\boldsymbol{s})) = \mu_i(\boldsymbol{s})$, $i = 1, 2$ and it can be shown that the entries of the cross-covariance function associated with $\boldsymbol{Z}_{12} = (Z_1(\boldsymbol{s}), Z_2(\boldsymbol{s}))^T$ i.e. $\boldsymbol{C_\theta}(\boldsymbol{h}) = [C_{ij;\boldsymbol{\theta}}(\boldsymbol{h})]_{i,j=1}^2 = [\text{cov}\,(Z_i(\boldsymbol{s}_l), Z_j(\boldsymbol{s}_m))]_{i,j=1}^2$, $\boldsymbol{h} = \boldsymbol{s}_l - \boldsymbol{s}_m \in \mathbb{R}^d$ are given by:

$$C_{ij;\boldsymbol{\theta}}(\boldsymbol{h}) = \begin{cases} \rho_{ij}(\sigma_i^2 + \tau_i^2)^{\frac{1}{2}}(\sigma_j^2 + \tau_j^2)^{\frac{1}{2}}, & \boldsymbol{h} = \boldsymbol{0}, \\ \sigma_i \sigma_j \rho_{ij} K_{\psi_{ij}}(\boldsymbol{h}), & otherwise. \end{cases} \tag{3}$$

where $\boldsymbol{\theta} = (\sigma_1^2, \sigma_2^2, \tau_1^2, \tau_2^2, \boldsymbol{\psi}_{11}^\top, \boldsymbol{\psi}_{12}^\top, \boldsymbol{\psi}_{22}^\top, \rho_{12})^\top$ is the vector parameter.

In the univariate setting, semi-variograms are often the main focus in geostatistics and are defined as the variance of contrasts. Similarly, in the bivariate setting, the semi-variogram matrix function can be defined as:

$$\boldsymbol{\Gamma}(\boldsymbol{h}) = [\gamma_{ij}(\boldsymbol{h})]_{i,j=1}^2 = 0.5\,[\text{cov}\,(Z_i(\boldsymbol{s}_l) - Z_i(\boldsymbol{s}_m), Z_j(\boldsymbol{s}_l) - Z_j(\boldsymbol{s}_m))]_{i,j=1}^2. \tag{4}$$

Under weakly stationarity and symmetry, the relation between the (cross) semi-variogram and the (cross) covariance is given by

$$\gamma_{ij;\boldsymbol{\theta}}(\boldsymbol{h}) = C_{ij;\boldsymbol{\theta}}(\boldsymbol{0}) - C_{ij;\boldsymbol{\theta}}(\boldsymbol{h}) \quad i, j = 1, 2. \tag{5}$$

The mapping $\boldsymbol{C_\theta}$ must be positive definite, which means that, for the bivariate random vector $\boldsymbol{Z}_N = (\boldsymbol{Z}_{1;N}^\top, \boldsymbol{Z}_{2;N}^\top)^\top$, where $\boldsymbol{Z}_{K;N} = (Z_K(\boldsymbol{s}_1), \dots Z_K(\boldsymbol{s}_N))^\top$, $i = 1, 2$, the $(2N) \times$

3

$(2N)$ associated covariance matrix $\boldsymbol{\Sigma_\theta} := [\boldsymbol{\Sigma}_{ij;\theta}]_{i,j=1}^2$ with $\boldsymbol{\Sigma}_{ij;\theta} = [C_{ij;\theta}(\boldsymbol{s}_l - \boldsymbol{s}_m)]_{l,m=1}^N$ is positive semidefinite.

In this general approach, the difficulty lies in deriving conditions on the model parameters that result in a valid (positive definite) multivariate covariance model. Note that

$$\rho_{12} = \frac{C_{ij;\theta}(\boldsymbol{0})}{\sqrt{C_{ii;\theta}(\boldsymbol{0})C_{jj;\theta}(\boldsymbol{0})}}$$

that is the colocated correlation parameters express the marginal correlation between the two marginal Gaussian random fields $Z_1$ and $Z_2$.

For instance Gneiting et al. (2010) proposed the model (3) with $K_\psi(\boldsymbol{h})$ equal to the Matérn isotropic correlation model:

$$\mathcal{M}_{\nu,\alpha}(\boldsymbol{h}) = \frac{1}{2^{\nu-1}\Gamma(\nu)} \left( \frac{||\boldsymbol{h}||}{\alpha} \right)^\nu K_\nu \left( \frac{||\boldsymbol{h}||}{\alpha} \right), \quad \alpha > 0, \nu > 0. \tag{6}$$

Putting together (3) and (6) we obtain the bivariate Matérn model

$$C_{ij;\theta}(\boldsymbol{h}) = \begin{cases} \rho_{ij}(\sigma_i^2 + \tau_i^2)^{\frac{1}{2}}(\sigma_j^2 + \tau_j^2)^{\frac{1}{2}}, & \boldsymbol{h} = \boldsymbol{0}, \\ \sigma_i\sigma_j\rho_{ij}\mathcal{M}_{\nu_{ij},\alpha_{ij}}(\boldsymbol{h}), & otherwise. \end{cases} \tag{7}$$

with $\boldsymbol{\theta} = (\sigma_1^2, \sigma_2^2, \nu_{11}, \nu_{12}, \nu_{22}, \alpha_{11}, \alpha_{12}, \alpha_{22}, \rho_{12})^T$. Gneiting et al. (2010) find a set of sufficient and necessary conditions on the colocated correlation parameter $\rho_{12}$ in order the model (7) to be valid.

If $\alpha_{11} = \alpha_{12} = \alpha_{22} > 0$ and $\nu_{11} = \nu_{12} = \nu_{22} > 0$ then the condition is $|\rho_{12}| < 1$ (in this case we obtain the so-called separable bivariate Matérn model). Otherwise the range of validity of $\rho_{12}$ is restricted to $|\rho_{12}| < a < 1$ with $a > 0$ and this kind of restriction on the upper and lower bound of the colocated parameter can be more or less severe depending on the scale and smoothness parameters (see Gneiting et al. (2010) for details).

Bivariate models of type (3) are implemented in the package `Geomodels` when $K_\psi(\boldsymbol{h})$ is a Matérn, a Generalized Wendland (Bevilacqua et al., 2019) and a Generalized Cauchy model (Gneiting and Schlather, 2004).

Suppose we want to simulate a realization of $\boldsymbol{Z}_{12}$ at $\boldsymbol{s}_1, \dots, \boldsymbol{s}_N$ location sites uniformly distributed in the unit square with $N = 500$ that is $\boldsymbol{z}_{500} = (\boldsymbol{z}_{1;500}^\top, \boldsymbol{z}_{2;500}^\top)^\top$, where $\boldsymbol{z}_{i;500} = (z_i(\boldsymbol{s}_1), \dots z_i(\boldsymbol{s}_{500}))^\top$, $i = 1, 2$. The total number of observations is given by $500 \times 2 = 1600$.

We first set the spatial coordinates:

```
NN=500 # number of spatial locations
x = runif(NN, 0, 1);
y = runif(NN, 0, 1)
coords=cbind(x,y)
```

We assume that the bivariate covariance function is given by a bivariate Matérn model in equation (7). We first we set the names of parameters of the bivariate covariance model in the package GeoModels.

```
corrmodel="Bi_Matern"
CorrParam("Bi_Matern")
 [1] "sill_1"    "sill_2"    "nugget_1"  "nugget_2"  "pcol" "scale_1"
 [7] "scale_12"  "scale_2"   "smooth_1"  "smooth_12" "smooth_2"
```

The previous names of parameters are associated to $\sigma_1^2$, $\sigma_2^2$, $\tau_1^2$, $\tau_2^2$ $\rho_{12}$, $\alpha_{11}$, $\alpha_{12}$, $\alpha_{22}$ and $\nu_{11}$, $\nu_{12}$, $\nu_{22}$ respectively. Note that the function CorrParam is useful since it returns the names of the correlation parameters for a given bivariate correlation model. Then we set the covariance parameters. In this example, we assume a common value (equal to 0.5) for the smoothness parameters a negative colocated correlation parameter and, for simplicity, two constant mean parameters $\mathbb{E}(Z_i(\boldsymbol{s})) = \mu_i$ with $\mu_1 = 2$ and $\mu_2 = -1$. However, the package GeoModels allows a mean regression specification for both mean functions. Additionally, we assume zero nugget parameters.

```
mean_1 = 2; mean_2= -1
nugget_1 =0;nugget_2=0
sill_1 =0.5; sill_2 =1;
scale_1=0.2/3; scale_2=0.15/3; scale_12=0.5*(scale_2+scale_1)
smooth_1=smooth_2=smooth_12=0.5
pcol=-0.4
param= list(nugget_1=nugget_1,nugget_2=nugget_2,
            sill_1=sill_1,sill_2=sill_2,
            mean_1=mean_1,mean_2=mean_2,
            smooth_1=smooth_1, smooth_2=smooth_2,smooth_12=smooth_12,
            scale_1=scale_1, scale_2=scale_2,scale_12=scale_12,
            pcol=pcol)
```

We are now ready to simulate the bivariate random fields using the function GeoSim:

```
ss1 = GeoSim(coordx=coords, corrmodel=corrmodel,
```

```
                    model=model,param=param)$data
dim(ss1)
[1]    2 500
```

The simulation is performed using Cholesky decomposition of the covariance matrix. The covariance matrix $\mathbf{\Sigma_\theta}$ can be obtained using the function `GeoCovmatrix` with the following code:

```
cc = GeoCovmatrix(coordx=coords, corrmodel=corrmodel,
                  model=model,param=param)
cc$covmatrix[1:4,1:4]
             [,1]          [,2]          [,3]          [,4]
[1,] 5.000000e-01 5.249200e-02 6.399256e-06 2.568249e-05
[2,] 5.249200e-02 5.000000e-01 1.639101e-05 1.020845e-04
[3,] 6.399256e-06 1.639101e-05 5.000000e-01 3.089991e-02
[4,] 2.568249e-05 1.020845e-04 3.089991e-02 5.000000e-01
```

## Estimation of a bivariate Gaussian random field

Given $\boldsymbol{z}_N = (\boldsymbol{z}_{1;N}^\top, \boldsymbol{z}_{2;N}^\top)^\top$, ($N = 500$ in this example), a realization of a bivariate Gaussian random field with bivariate Matérn covariance function, the estimation can be performed with maximum likelihood or weighted pairwise likelihood.

Maximum likelihood involves the maximization of the log-likelihood function:

$$l_N(\boldsymbol{\psi}) = -\frac{1}{2}\log|\mathbf{\Sigma_\theta}| - \frac{1}{2}(\boldsymbol{z}_N - \boldsymbol{\mu}_{12})^\top[\mathbf{\Sigma_\theta}]^{-1}(\boldsymbol{z}_N - \boldsymbol{\mu}_{12}). \tag{8}$$

where $\boldsymbol{\psi} = (\mu_1, \mu_2, \boldsymbol{\theta}^T)^T$. Here $\boldsymbol{\mu}_{12} = (\mathbf{1}\mu_1, \mathbf{1}\mu_2)^T$ and $\mathbf{1}$ is the unit vector of length $N$. To perform maximum likelihood estimation, we first set the parameters that we want to estimate and the parameters that we want to fix with the following two lists.

```
fixed=list(nugget_1=nugget_1,nugget_2=nugget_2,
           smooth_1=smooth_1, smooth_2=smooth_2,smooth_12=smooth_12)
start=list(mean_1=mean_1,mean_2=mean_2,sill_1=sill_1,sill_2=sill_2,
           scale_1=scale_1,scale_2=scale_2,scale_12=scale_12, pcol=pcol)
```

We are now ready to perform maximum likelihood estimation using the function `GeoFit`:

```
fit = GeoFit(data=ss1,coordx=coords, corrmodel=corrmodel,
             likelihood="Full",type="Standard",optimizer="BFGS",
```

```
        start=start,fixed=fixed)
```

The object `fit` include informations about the maximum likelihood estimation.

```
fit
#######################################################################
Maximum Likelihood Fitting of Gaussian Random Fields
Setting: Full Likelihood
Model: Gaussian
Type of the likelihood objects: Standard
Covariance model: Bi_Matern
Optimizer: BFGS
Number of spatial coordinates: 500
Number of dependent temporal realisations: 1
Type of the random field: bivariate
Number of estimated parameters: 8
Type of convergence: Successful
Maximum log-Likelihood value: -892.30
AIC : 1801
BIC : 1840
Estimated parameters:
   mean_1     mean_2        pcol    scale_1   scale_12    scale_2     sill_1
 1.88077   -1.12728   -0.41080    0.06293    0.04439    0.04051    0.48158
   sill_2
 0.94915
#######################################################################
```

Maximum likelihood can be computationally demanding if $N$ is large. An alternative method of estimation that can be useful in this case is the weighted composite likelihood method proposed in Bevilacqua et al. (2016). The weighted composite likelihood method involves the maximization of the function:

$$pl(\boldsymbol{\psi}) = \sum_{(i,j,l,m)\in\Lambda} l_{ijlm}(\boldsymbol{\psi})w_{ijlm}, \tag{9}$$

where $\Lambda$ is a specific index set (see Bevilacqua et al. (2016)), $l_{ijlm}(\boldsymbol{\psi})$ is the log-likelihood of the bivariate Gaussian random vector $[Z_i(\mathbf{s}_l), Z_j(\mathbf{s}_m)]^T$ and $w_{ijlm}$ are positive suitable weights specified as:

$$w_{ijlm} = \begin{cases} 1, & \|\boldsymbol{s}_l - \boldsymbol{s}_m\| \leq d_{ij} \\ 0, & otherwise \end{cases}, \tag{10}$$

Here $d_{ij} > 0$ are the compact support of the weight function that are arbitrary fixed. Specifically $d_{11} > 0$ is the compact support for the first component $d_{22} > 0$ is the compact support for the second component and $d_{12} = d_{21}$ is the compact support for the cross cases.

To compute weighted composite likelihood estimation we use the function `GeoFit`

```
fit_pl = GeoFit(data=ss1,coordx=coords, corrmodel=corrmodel,
 maxdist=c(0.1,0.1,0.1),likelihood="Marginal",type="Pairwise",
 optimizer="BFGS", start=start, fixed=fixed)
```

Note that the option `maxdist=c(0.1,0.1,0.1)` set the (arbitrary) compact supports of the weight function (10) i.e. $d_{11} = 0.1$, $d_{21} = d_{12} = 0.1$ and $d_{22} = 0.1$ respectively. A suitable choice of the compact supports of the weights allows to improve both the statistical and computational efficiency (Bevilacqua and Gaetan (2015))

The object `fit_pl` include informations about the weighted composite likelihood estimation.

```
fit_pl
##################################################################
Maximum Composite-Likelihood Fitting of Gaussian Random Fields
Setting: Marginal Composite-Likelihood
Model: Gaussian
Type of the likelihood objects: Pairwise
Covariance model: Bi_Matern
Optimizer: BFGS
Number of spatial coordinates: 500
Number of dependent temporal realisations: 1
Type of the random field: bivariate
Number of estimated parameters: 8
Type of convergence: Successful
Maximum log-Composite-Likelihood value: -33923.28
Estimated parameters:
  mean_1    mean_2       pcol    scale_1   scale_12    scale_2     sill_1
 1.85144  -1.09028   -0.39890    0.05728    0.03561    0.04246    0.45302
  sill_2
```
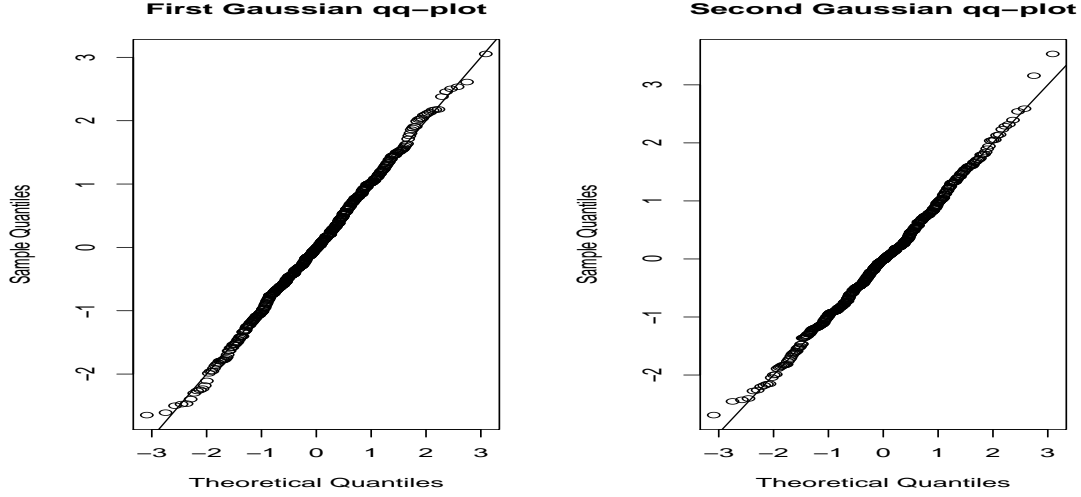
Figure 1: Gaussian qq-plot of the bivariate residuals.

```
 0.94970
##################################################################
```

Given the estimation of the mean and variance parameters, the estimated residuals

$$\widehat{y}_i(\boldsymbol{s}_j) = \frac{z_i(\boldsymbol{s}_j) - \hat{\mu}_i}{(\widehat{\sigma}_i^2)^{\frac{1}{2}}} \quad j = 1, \ldots, N, \quad i = 1, 2$$

can be viewed as a realization of a zero mean unit variance bivariate GRF. Using pairwise likelihood estimates, the residuals can be computed using the `GeoResiduals` function:

```
res=GeoResiduals(fit_pl)
```

Then the marginal distribution assumption on the bivariate residuals can be graphically checked with a Gaussian qq-plot (Figure 1) using the function `GeoQQ`.

```
### checking model assumptions: marginal distribution
GeoQQ(res)
```

In order to check the adequacy of the estimated bivariate covariance model we can compare the empirical semi-variogram estimation of (4) ( *i.e.* the marginal variograms and the cross-variogram) with the estimated ones plugging-in the estimated parameters in the bivariate covariance model and using relation (5).

We first compute the semivariograms using the function `GeoVariogram` with the option `bivariate=TRUE`. Then the function `GeoCovariogram` allows to graphically compare the
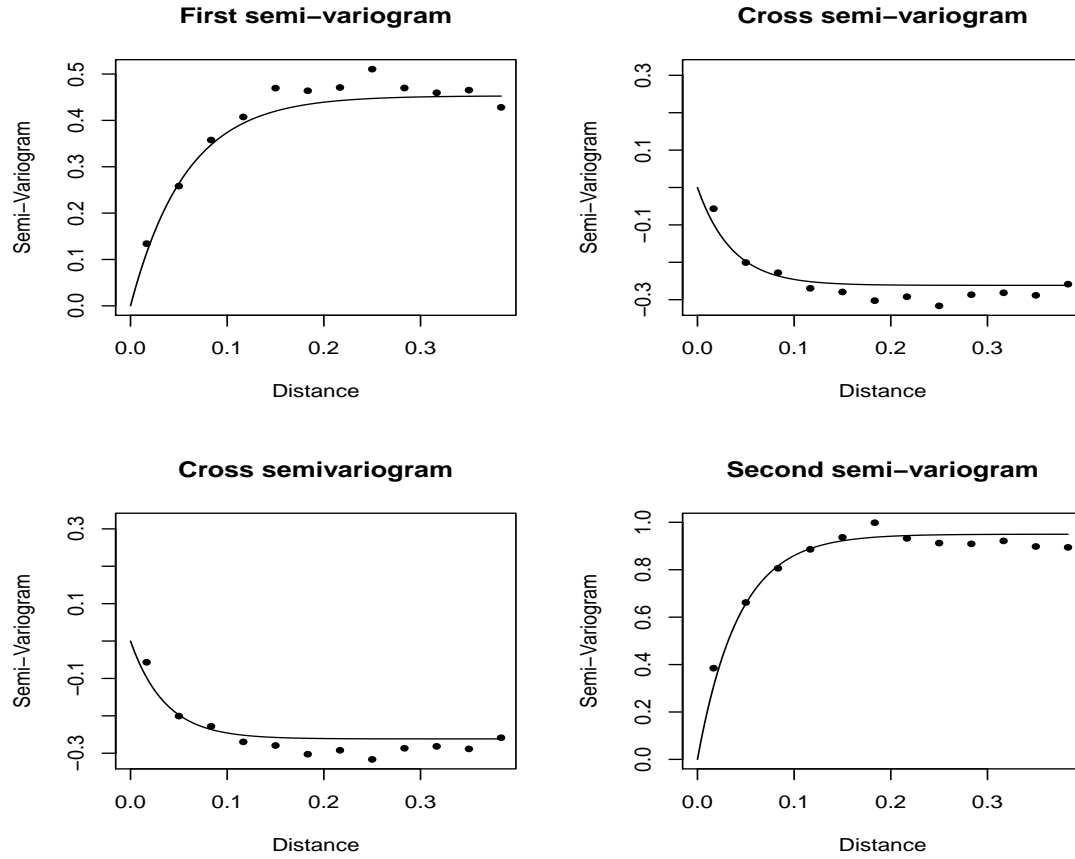
Figure 2: Empirical marginal and cross semi-variograms versus the estimated ones using maximum weighted composite likelihood estimation.

empirical semivariogram estimation with the estimated one (see Figure 2) using maximum weighted composite likelihood estimation (object `fit_pl`)

```
vario = GeoVariogram(data=ss1,coordx=coords, bivariate=TRUE,
               maxdist=c(0.4,0.4,0.4))
GeoCovariogram(fit_pl,vario=vario,show.vario=TRUE,pch=20)
```

## Prediction of bivariate Gaussian random fields

For a given spatial location $(s_0)$ the optimal prediction for a component of a bivariate Gaussian RF (co-kriging) is computed as:

$$\widehat{Z}_i(\boldsymbol{s}_0) = \mu_i + \boldsymbol{c}_{\boldsymbol{\theta}}^T \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1}(\boldsymbol{z}_N - \boldsymbol{\mu}_{12}), \qquad i = 1, 2 \tag{11}$$

10

where $c_{\theta} = (\text{cov}(Z_1(s_0), Z_1(s_1)), \ldots, \text{cov}(Z_1(s_0), Z_1(s_N)), \ldots, \text{cov}(Z_2(s_0), Z_2(s_N)))^T$.

Optimal prediction (plugging-in the estimated parameters in (11)), can be performed using the `GeoKrig` function. We need just to specify the spatial locations to be predict. In this example, we consider a spatial regular grid on the unit square:

```
xx=seq(0,1,0.012)
loc_to_pred=as.matrix(expand.grid(xx,xx))
```

Then optimal prediction (11), using the estimated parameters with pairwise maximum likelihood (object `fit_pl` ), can be performed using the `GeoKrig` function for the first and the second random field with the following code:

```
param_est=as.list(c(fit_pl$param,fixed))
pr1 = GeoKrig(data=ss1,coordx=coords,  corrmodel=corrmodel,which=1,
              model=model,mse=TRUE,loc=loc_to_pred,param=param_est)
pr2 = GeoKrig(data=ss1,coordx=coords,  corrmodel=corrmodel,which=2,
              model=model,mse=TRUE,loc=loc_to_pred,param=param_est)
```

Note that the option `which` allows to set the component of the bivariate Gaussian field to be predicted. A kriging map with associate mean square error (Figure 3) for the two random fields can be obtained with the following code:

```
par(mfrow=c(2,3))
colour <- rainbow(100)
quilt.plot(coords[,1],coords[,2],ss1[1,],col=colour,
                    main ="Observed␣data:First␣variable")
image.plot(xx, xx, matrix(pr1$pred,ncol=length(xx)),col=colour,
                    main = paste("Kriging"),ylab="")
image.plot(xx, xx, matrix(pr1$mse,ncol=length(xx)),col=colour,
                      main = paste("MSE"),ylab="")
quilt.plot(coords[,1],coords[,2],ss1[2,],col=colour,
                      main ="Observed␣data:Second␣variable")
image.plot(xx, xx, matrix(pr2$pred,ncol=length(xx)),col=colour,
                    main = paste("Kriging"),ylab="")
image.plot(xx, xx, matrix(pr2$mse,ncol=length(xx)),col=colour,
                    main = paste("MSE"),ylab="")
```
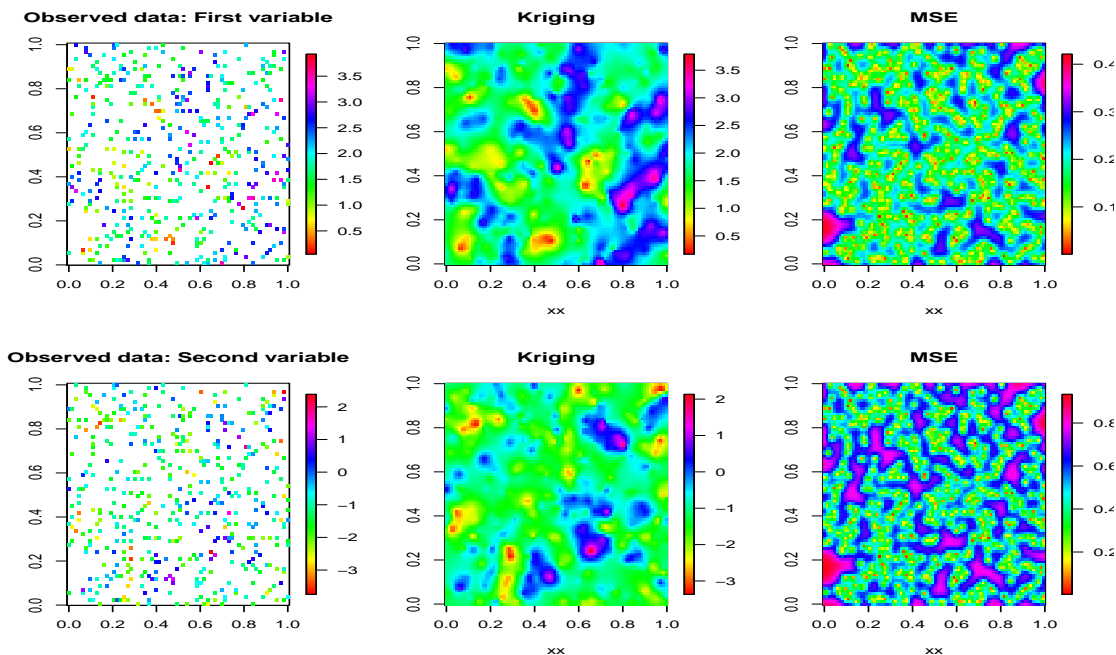
Figure 3: Observed spatial data, kriging prediction and associated mean square error for the two components of the bivariate random field with bivariate Matérn covariance model.

# References

Bevilacqua, M., A. Alegria, D. Velandia, and E. Porcu (2016). Composite likelihood inference for multivariate gaussian random fields. *Journal of Agricultural Biological and Environmental Statistics 21(3)*, 1236–1249.

Bevilacqua, M., T. Faouzi, R. Furrer, and E. Porcu (2019). Estimation and prediction using generalized Wendland functions under fixed domain asymptotics. *The Annals of Statistics 47(2)*, 828–856.

Bevilacqua, M. and C. Gaetan (2015). Comparing composite likelihood methods based on pairs for spatial Gaussian random fields. *Statistics and Computing 25*, 877–892.

Bevilacqua, M. and V. Morales-Oñate (2018). *GeoModels: A Package for Geostatistical Gaussian and non Gaussian Data Analysis*. R package version 1.0.3-4.

Genton, M. G. and W. Kleiber (2015, 05). Cross-covariance functions for multivariate geostatistics. *Statistical Science 30*(2), 147–163.

Gneiting, T., W. Kleiber, and M. Schlather (2010). Matérn Cross-Covariance functions for multivariate random fields. *Journal of the American Statistical Association 105*, 1167–1177.

Gneiting, T. and M. Schlather (2004). Stochastic models that separate fractal dimension and the hurst effects. *SIAM Rev. 46*, 269–282.

Wackernagel, H. (2003). *Multivariate Geostatistics: An Introduction with Applications* (3rd ed.). New York: Springer.