

# GeoModels Tutorial: analysis 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(12)
```

## Simulation of a bivariate Gaussian random field

Let  $\mathbf{Z}_{12} = \{\mathbf{Z}_{12}(\mathbf{s}) = (Z_1(\mathbf{s}), Z_2(\mathbf{s}))^T, \mathbf{s} \in A \subseteq \mathbb{R}^d\}$  be a bivariate Gaussian random field. where  $Z_i = \{Z_i(\mathbf{s}), \mathbf{s} \in A \subseteq \mathbb{R}^d\}$ ,  $i = 1, 2$  are two univariate Gaussian random fields

In this tutorial we assume that the vector of the means are constant *i.e.*  $\mathbb{E}(\mathbf{Z}_{12}(\mathbf{s})) = \boldsymbol{\mu}$  with  $\boldsymbol{\mu} = (\mu_1, \mu_2)^T$  and  $\mathbb{E}(Z_i(\mathbf{s})) = \mu_i$ .

Under second order stationary assumption, the variances of the two random fields are bounded *i.e.*  $\mathbb{V}(Z_i(\mathbf{s})) = \sigma_i^2$ ,  $i = 1, 2$ , and the covariance function between  $\mathbf{Z}_{12}(\mathbf{s}_l)$  and  $\mathbf{Z}_{12}(\mathbf{s}_m)$ , for any pair  $\mathbf{s}_l, \mathbf{s}_m$  in the spatial domain, is represented by a mapping  $\mathbf{C} : \mathbb{R}^d \rightarrow M_{2 \times 2}$  defined through

$$\mathbf{C}(\mathbf{h}) = [C_{ij}(\mathbf{h})]_{i,j=1}^2 = [\text{cov}(Z_i(\mathbf{s}_l), Z_j(\mathbf{s}_m))]_{i,j=1}^2, \quad \mathbf{h} = \mathbf{s}_l - \mathbf{s}_m \in \mathbb{R}^d. \quad (1)$$

The function  $\mathbf{C}(\mathbf{h})$  is called bivariate covariance function. Here,  $M_{2 \times 2}$  is the set of two dimensional squared, symmetric and positive definite matrices. The functions  $C_{ii}(\mathbf{h})$   $i = 1, 2$  are the marginal covariance functions of the Gaussian random fields  $Z_i$ ,  $i = 1, 2$  while  $C_{ij}(\mathbf{h})$  is called cross covariance function between  $Z_i$  and  $Z_j$  for  $i, j = 1, 2$  and  $i \neq j$  at spatial lag  $\mathbf{h}$ .

The cross-covariance function is not in general symmetric, *i.e.*  $C_{12}(\mathbf{h}) \neq C_{21}(\mathbf{h})$  (Wackernagel, 2003). However, the majority of the existing multivariate parametric covariance models are symmetric, with some few exceptions (Genton and Kleiber, 2015).

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:

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

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

$$\gamma_{ij}(\mathbf{h}) = C_{ij}(\mathbf{0}) - C_{ij}(\mathbf{h}) \quad i, j = 1, 2. \quad (3)$$

The mapping  $\mathbf{C}$  must be positive definite, which means that, for the bivariate random vector  $\mathbf{Z}_N = (\mathbf{Z}_{1;N}^\top, \mathbf{Z}_{2;N}^\top)^\top$ , where  $\mathbf{Z}_{k;N} = (Z_k(\mathbf{s}_1), \dots, Z_k(\mathbf{s}_N))^\top$ ,  $k = 1, 2$ , the  $(2N) \times (2N)$  associated covariance matrix  $\mathbf{\Sigma} := [\mathbf{\Sigma}_{ij}]_{i,j=1}^2$  with  $\mathbf{\Sigma}_{ij} = [C_{ij}(\mathbf{s}_l - \mathbf{s}_m)]_{l,m=1}^n$  is positive semidefinite.

We shall assume throughout that the mapping  $\mathbf{C}$  comes from a parametric family of bivariate covariances  $\{\mathbf{C}_\theta(\cdot), \theta \in \Theta \subseteq R^p\}$ , with  $\theta$  an arbitrary parametric space.

A useful general symmetric parametric class is obtained through the following bivariate covariance function:

$$C_{ij;\theta}(\mathbf{h}) = \begin{cases} \rho_{ij}(\sigma_i^2 + \tau_i^2)^{\frac{1}{2}}(\sigma_j^2 + \tau_j^2)^{\frac{1}{2}}, & \mathbf{h} = \mathbf{0}, \\ \rho_{ij}\sigma_i\sigma_j R_{\psi_{ij}}(\mathbf{h}), & \text{otherwise.} \end{cases} \quad (4)$$

where  $\rho_{ii} = 1$ ,  $i = 1, 2$  and  $\theta = (\sigma_1^2, \sigma_2^2, \tau_1^2, \tau_2^2, \psi_{11}^\top, \psi_{12}^\top, \psi_{22}^\top, \rho_{12})^\top$  and  $R_\psi(\mathbf{h})$  is a univariate parametric correlation model.

In this general approach, the difficulty lies in deriving conditions on the model parameters that result in a valid multivariate covariance model. Here  $\sigma_i^2 > 0$ ,  $i = 1, 2$  are the marginal variance parameters,  $\tau_i^2 > 0$ ,  $i = 1, 2$  are the marginal nugget parameters and  $\rho_{12}$ , is the so-called colocated correlation parameter. Note that

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

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 (4) with  $R(\mathbf{h})$  equal to the Matérn isotropic correlation model:

$$\mathcal{M}_{\nu,\alpha}(\mathbf{h}) = \frac{1}{2^{\nu-1}\Gamma(\nu)} \left( \frac{\|\mathbf{h}\|}{\alpha} \right)^\nu K_\nu \left( \frac{\|\mathbf{h}\|}{\alpha} \right), \quad \alpha > 0, \nu > 0. \quad (5)$$

Putting together (4) and (5) and assuming zero nuggets for simplicity of notation, we obtain the bivariate Matérn model

$$\mathbf{C}_\theta(\mathbf{h}) = [\rho_{ij}\sigma_i\sigma_j\mathcal{M}_{\nu_{ij},\alpha_{ij}}(\mathbf{h})]_{i,j=1}^2, \rho_{ii} = 1. \quad (6)$$

with  $\theta = (\sigma_1^2, \sigma_2^2, \nu_{11}, \nu_{12}, \nu_{22}, \alpha_{11}, \alpha_{12}, \alpha_{22}, \rho_{12})^T$ . Gneiting et al. (2010) find a set of sufficient and necessary conditions on the colocated correlation parameter  $\rho_{12}$  in order the model (6) to be valid.

If  $\alpha_{11} = \alpha_{12} = \alpha_{22}$  and  $\nu_{11} = \nu_{12} = \nu_{22}$  then  $|\rho_{12}| < 1$  (in this case we obtain the so-called separable bivariate Matérn model). Otherwise the range of validity of  $\rho_{12}$  is restricted to  $|\rho_{12}| < a < 1$  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 (Gneiting et al., 2010).

Bivariate models of type (4) are implemented in the package **Geomodels** when  $R_\psi$  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  $\mathbf{Z}_{12}$  at  $\mathbf{s}_1, \dots, \mathbf{s}_N$  location sites uniformly distributed in the unit square with  $N = 800$  that is  $\mathbf{z}_{800} = (\mathbf{z}_{1;800}^\top, \mathbf{z}_{2;800}^\top)^\top$ , where  $\mathbf{z}_{k;800} = (z_k(\mathbf{s}_1), \dots, z_k(\mathbf{s}_{800}))^\top$ ,  $k = 1, 2$ . The total number of observations is given by  $800 \times 2 = 1600$ .

We first set the spatial coordinates:

```

NN=800 # 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 (6). We first we set the name 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 in the package `GeoModels`. Then we set the covariance parameters. In this case we assume a common value (equal to 0.5) for the smoothness parameter, a negative colocated correlation parameter and zero means for both marginal Gaussian random field.

```
mean_1 = 0; mean_2= 0
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] 2 800
```

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

```
cc = GeoCovmatrix(coordx=coords, corrmodel=corrmodel,
                model=model,param=param)
cc$covmatrix[1:5,1:5]
      [,1]      [,2]      [,3]      [,4]
[1,] 5.000000e-01 3.959941e-06 2.004530e-07 1.746098e-02
[2,] 3.959941e-06 5.000000e-01 1.579141e-05 1.064041e-04
[3,] 2.004530e-07 1.579141e-05 5.000000e-01 1.094790e-06
[4,] 1.746098e-02 1.064041e-04 1.094790e-06 5.000000e-01
```

## Estimation of a bivariate Gaussian random field

Given  $\mathbf{z}_N$  ( $N = 800$  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:

$$l_N(\boldsymbol{\theta}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}_N(\boldsymbol{\theta})| - \frac{1}{2} (\mathbf{z}_N - \boldsymbol{\mu}_{12})^\top [\boldsymbol{\Sigma}_N(\boldsymbol{\theta})]^{-1} (\mathbf{z}_N - \boldsymbol{\mu}_{12}). \quad (7)$$

where  $\boldsymbol{\theta} = (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \nu_{11}, \nu_{12}, \nu_{22}, \alpha_{11}, \alpha_{12}, \alpha_{22}, \rho_{12})^T$ . 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 two lists.

```
fixed=list(nugget_1=nugget_1,nugget_2=nugget_2,
          mean_1=mean_1,mean_2=mean_2,
          smooth_1=smooth_1, smooth_2=smooth_2,smooth_12=smooth_12)
start=list(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: 800
Number of dependent temporal realisations: 1
Type of the random field: bivariate
Number of estimated parameters: 6
Type of convergence: Successful
```

```

Maximum log-Likelihood value: -1268.00
AIC : 2548
BIC : 2580
Estimated parameters:
      pcol  scale_1  scale_12  scale_2  sill_1  sill_2
-0.35548  0.05294  0.04776  0.04679  0.43802  0.89745
#####

```

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:

$$pl(\boldsymbol{\theta}) = \sum_{(i,j,l,m) \in \Lambda} l_{ijmn}(\boldsymbol{\theta}) w_{ijmn}, \quad (8)$$

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

$$w_{ijlm} = \begin{cases} 1, & \|\mathbf{s}_l - \mathbf{s}_m\| \leq d_{ij} \\ 0, & \text{otherwise} \end{cases}, \quad (9)$$

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 (9) 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: 800
Number of dependent temporal realisations: 1
Type of the random field: bivariate
Number of estimated parameters: 6
Type of convergence: Successful
Maximum log-Composite-Likelihood value: -88758.97
Estimated parameters:
      pcol  scale_1  scale_12  scale_2  sill_1  sill_2
-0.34556  0.05196  0.04685  0.05480  0.42733  0.99172
#####

```

In order to check the adequacy of the estimated bivariate covariance model we can compare the empirical semi-variogram estimation of (2) ( *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 (3).

We first compute the semivariograms using the function `GeoVariogram` with the option `bivariate=TRUE`. Then the function `GeoCovariogram` allows to graphically compare the empirical semivariogram estimation with the estimated one (see Figure 1) using maximum weighted composite likelihood estimation (object `fit_pl`)

```

vario = GeoVariogram(data=ss1, coordx=coords, bivariate=TRUE,
                     maxdist=c(0.6,0.6,0.6))
GeoCovariogram(fit_pl, vario=vario, show.vario=TRUE, pch=20)

```



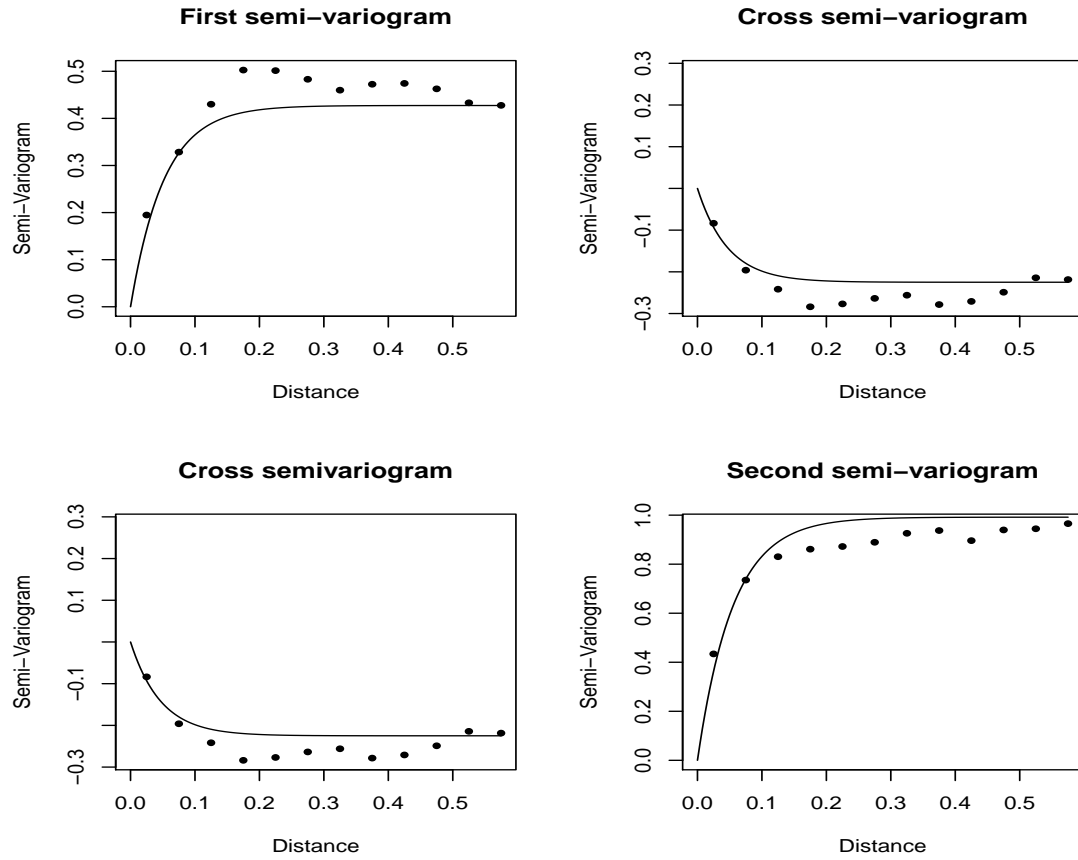


Figure 1: Empirical marginal and cross semi-variograms versus the estimated ones using maximum weighted composite likelihood estimation.

## Prediction of bivariate Gaussian random fields

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

$$\widehat{Y}_k(s_0) = \mu_k + \mathbf{c}^T(\boldsymbol{\theta})\boldsymbol{\Sigma}_N^{-1}(\boldsymbol{\theta})(\mathbf{z}_N - \boldsymbol{\mu}_{12}), \quad k = 1, 2 \quad (10)$$

where  $\mathbf{c} = (\text{cov}(Z_1(s_0), Z_1(s_1)), \dots, \text{cov}(Z_1(s_0), Z_1(s_N)), \dots, \text{cov}(Z_2(s_0), Z_2(s_N)))^T$ .

Optimal prediction (plugging-in the estimated parameters in (10)), 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(0,1,0.015)
loc_to_pred=as.matrix(expand.grid(xx,xx))
```

Then optimal prediction (10), using the estimated parameters using maximum likelihood (object `fit`), 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$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 2) 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")
```

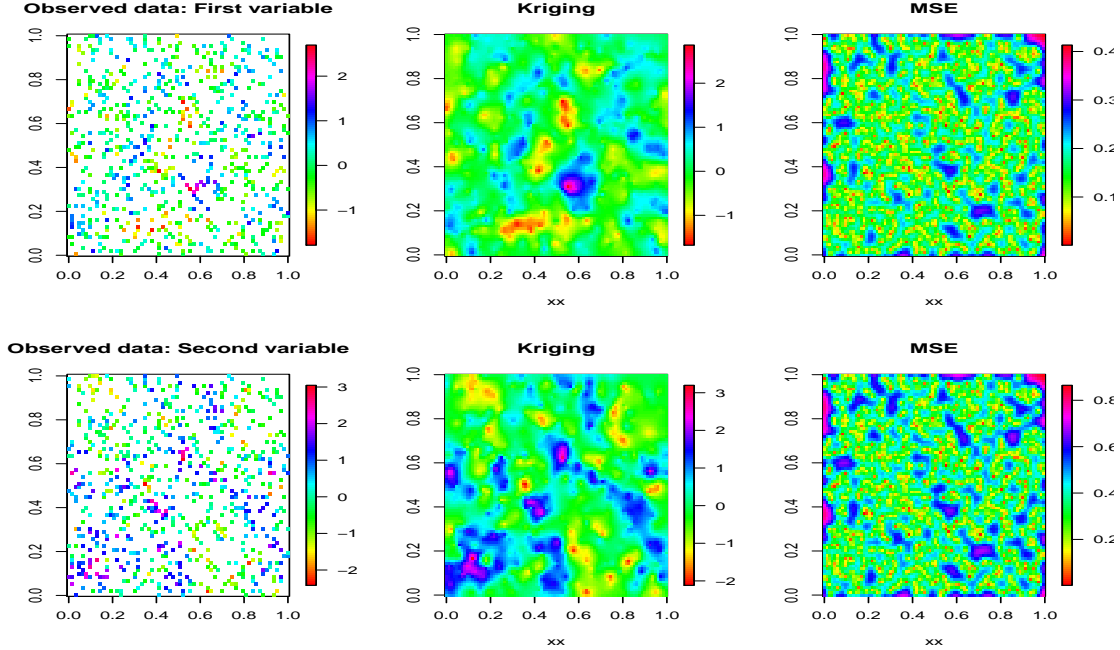


Figure 2: Observed spatial data, kriging prediction and associated mean square error for the two components of the bivariate random field with bivariate Matérn covariance model.

```
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="")
```

## 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. *Statist. Sci.* *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.