

# GeoModels Tutorial: simulation, estimation and prediction of positive spatial data using log-gaussian random fields

Moreno Bevilacqua  
Christian Caamaño

## Introduction

In this tutorial we show how to analyze geo-referenced spatial data with positive support using log-gaussian random fields (RFs) (Oliveira, 2006; Oliveira et al., 1997) with the R package **GeoModels** (Bevilacqua et al. (2023)). In particular, log-Gaussian processes have been widely used for the analysis of positive dependent data due to their well-known mathematical properties. The log-gaussian distribution is a flexible parametric model for positive data allowing right skewness.

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(GeoModels)  
require(fields)  
model="LogGaussian" # model name in the GeoModels package
```

## Simulation of log-Gaussian random fields

The definition of a log-Gaussian RF starts by considering a ‘parent’ Gaussian RF  $Z = \{Z(\mathbf{s}), \mathbf{s} \in S\}$ , where  $\mathbf{s}$  represents a location in the domain  $S$ . In this tutorial, we assume  $S = [0, 1]^2 \subseteq \mathbb{R}^2$  and that  $Z$  is stationary with zero mean, unit variance and correlation function  $\rho(\mathbf{h}) := \text{cor}(Z(\mathbf{s} + \mathbf{h}), Z(\mathbf{s}))$ .

Then a RF  $V = \{V(\mathbf{s}), \mathbf{s} \in S\}$  with marginal distribution  $\text{LogGaussian}(1, \sigma^2)$  can be derived by the transformation

$$V(\mathbf{s}) = \exp\{\sigma Z(\mathbf{s}) - \sigma^2/2\} \quad (1)$$

where  $\sigma > 0$  is a scale parameter. Under this specific parametrization,  $\mathbb{E}(V(\mathbf{s})) = 1$   $\text{var}(V(\mathbf{s})) = (\exp(\sigma^2) - 1)$  and the correlation function is given by:

$$\rho_V(\mathbf{h}) = \frac{\exp(\sigma^2 \rho(\mathbf{h})) - 1}{\exp(\sigma^2) - 1}. \quad (2)$$

Then a non stationary version can be defined through a multiplicative model as:

$$Y(\mathbf{s}) = \mu(\mathbf{s})V(\mathbf{s}), \quad \mu(\mathbf{s}) > 0 \quad (3)$$

with  $\mathbb{E}(Y(\mathbf{s})) = \mu(\mathbf{s})$ ,  $\text{var}(Y(\mathbf{s})) = \mu(\mathbf{s})^2(\exp(\sigma^2) - 1)$ . A spatial regression model can be obtained by assuming that  $\mu(\mathbf{s}) = e^{X(\mathbf{s})^T \boldsymbol{\beta}}$  where  $X(\mathbf{s})$  is a  $k$ -dimensional vector of covariates and  $\boldsymbol{\beta} = (\beta_0, \dots, \beta_k)^T$  is a  $k$ -dimensional vector of (unknown) parameters.

Thus, in order to obtain a realization from a log-Gaussian RF we need to specify a regression mean parameters, a scale parameter and a parametric correlation model for  $\rho(\mathbf{h})$ . We first set the spatial coordinates

```
N=1500 # number of location sites
set.seed(89)
x = runif(N, 0, 1)
y = runif(N, 0, 1)
coords=cbind(x,y)
```

Then we fix  $k = 2$  and we build the matrix covariates and fix the regression mean parameters

```
X=cbind(rep(1,N),runif(N)) # matrix covariates
mean = -0.2; mean1 =0.25 # regression parameters
nugget=0 # nugget parameter
sill=1
```

where `mean` and `mean1` are respectively  $\beta_1$  and  $\beta_2$ . Finally, we set the `sill` parameter of the log-Gaussian RF and the `nugget` parameter are respectively  $\sigma^2$  and  $\tau$ .

The names of the marginal parameters associated with the log-Gaussian model can be obtained with the function `NuisParam` (note that the option `num_betas` is the number of regression parameters involved):

```
NuisParam(model,num_betas=2)
[1] "mean" "mean1" "nugget" "sill"
```

For the correlation function we assume a special case of the isotropic Generalized Wendland class (Bevilacqua et al. (2019)) i.e the Askey model.

$$\rho(\mathbf{h}; \alpha, \delta) := \begin{cases} (1 - \|\mathbf{h}\|/\alpha)^\delta & \|\mathbf{h}\| < \alpha \\ 0 & \text{otherwise} \end{cases}.$$

Using asymptotic arguments Bevilacqua et al. (2019) show that this correlation model has the same features of the exponential correlation model. Additionally it is compactly supported an interesting feature from computational point of view. We set the Askey model

and the associated parameters. Note that the function `CorrParam` returns the names of the parameters associated for a given correlation model.

```
corrmodel = "Wend0" ## correlation model and parameters
CorrParam("Wend0")
[1] "power2" "scale"
scale = 0.25
power2 = 4
```

Here the `scale` parameter corresponds to  $\alpha$ , the compact support of the correlation model.

We are now ready to simulate a log-Gaussian random field using the function `GeoSim`:

```
param=list(mean=mean,mean1=mean1,sill=sill, nugget=nugget,
  scale=scale ,power2=power2)
data = GeoSim(coordx=coords , corrmodel=corrmodel , model=model ,
  param=param , X=X)$data
```

The simulation is performed using Cholesky decomposition.

## Estimation of log-Gaussian random fields

The density of the bivariate random vector  $(V(\mathbf{s}_i), V(\mathbf{s}_j))^T$  is given by

$$f_{\mathbf{V}}(v_i, v_j) = \frac{1}{2\pi\sigma^2 v_i v_j \sqrt{1 - \rho^2(\mathbf{h})}} e^{-\frac{1}{2(1-\rho^2(\mathbf{h}))} \left[ \left( \frac{\log(v_i)}{\sigma} + \frac{\sigma}{2} \right)^2 + \left( \frac{\log(v_j)}{\sigma} + \frac{\sigma}{2} \right)^2 - 2\rho(\mathbf{h}) \left( \frac{\log(v_i)}{\sigma} + \frac{\sigma}{2} \right) \left( \frac{\log(v_j)}{\sigma} + \frac{\sigma}{2} \right) \right]}, \quad (4)$$

and the bivariate densities of  $Y$  can be derived from (4) as

$$f_Y(y_i, y_j) = (\mu_i \mu_j)^{-1} f_V(y_i/\mu_i, y_j/\mu_j). \quad (5)$$

Given a realization  $y(\mathbf{s}_1), \dots, y(\mathbf{s}_n)$  of  $Y$ , then, the conditional pairwise likelihood function is defined as:

$$pl(\boldsymbol{\theta}) = \sum_{i=1}^N \sum_{j=1, j \neq i}^N \log(f_Y(y(\mathbf{s}_i), y(\mathbf{s}_j)) / f_Y(y(\mathbf{s}_j))) w_{ij}$$

where  $w_{ij}$  are non-negative weights, not depending on  $\boldsymbol{\theta}$ . An efficient way to specify the weights from computational and efficient viewpoint is based on neighborhoods:

$$w_{ij}(k) = \begin{cases} 1 & \mathbf{s}_i \in N_k(\mathbf{s}_j) \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

Here  $N_k(s_l)$  is the set of the neighbors of order  $k = 1, 2, \dots$  of the point  $s_l$  and in this case  $\boldsymbol{\theta} = (\beta_0, \beta_1, \sigma, \alpha, \delta)^T$ . The conditional pairwise likelihood estimator  $\hat{\boldsymbol{\theta}}_{pl}$  is obtained maximizing (5) with respect to  $\boldsymbol{\theta}$ . In the `GeoModels` package we can choose the fixed parameters and the parameters that must be estimated. Conditional pairwise likelihood estimation is performed with the function `GeoFit`:

```
start=list(mean=mean,mean1=mean1,scale=scale,sill=sill)
fixed=list(nugget=nugget ,power2=power2)

# Maximum composite-likelihood fitting of the loggaussian random field:
fit = GeoFit2(data=data,coordx=coords, corrmodel=corrmodel,
model=model,X=X,likelihood="Conditional",type="Pairwise",
start=start,fixed=fixed,neighb=1)
```

The option `neighb=1` sets the  $k$  value (i.e. the order of neighborhood) in the weight function (6). The object `fit` include informations about the conditional pairwise likelihood estimation

```
fit
#####
Maximum Composite-Likelihood Fitting of Log Gaussian Random Fields
Setting: Conditional Composite-Likelihood
Model: LogGaussian
Type of the likelihood objects: Pairwise
Covariance model: Wend0
Optimizer: Nelder-Mead
Number of spatial coordinates: 1500
Number of dependent temporal realisations: 1
Type of the random field: univariate
Number of estimated parameters: 4
Type of convergence: Successful
Maximum log-Composite-Likelihood value: -144.02
Estimated parameters:
      mean      mean1      scale      sill
-0.3310    0.2640    0.2794    1.0579
#####
```

## Checking model assumptions

Given the estimation of the mean  $\widehat{\mu}(\mathbf{s}) = e^{X_1(\mathbf{s})\hat{\beta}_1 + X_2(\mathbf{s})\hat{\beta}_2}$ , the estimated residuals

$$\widehat{v}(\mathbf{s}_i) = y(\mathbf{s}_i) / \widehat{\mu}(\mathbf{s}_i) \quad i = 1, \dots, N \quad (7)$$

can be viewed as a realization of  $V$  a stationary RF with marginal distribution  $Loggaussian(1, \sigma^2)$  with unit mean and correlation function given by (2).

The estimated residuals can be computed using the `GeoResiduals` function:

```
res=GeoResiduals(fit) # computing residuals
```

Then the agreement of the marginal distribution assumption on the residuals with the theoretical model can be graphically checked with the `GeoQQ` function (Figure 1 left part):

```
GeoQQ(res)
```

The covariance model assumption can be checked comparing the empirical and the estimated semivariogram using the `GeoVariogram` and `GeoCovariogram` functions (Figure 1 right part). In particular the function `GeoVariogram` compute the empirical semivariogram:

```
### checking model residuals assumptions: semivariogram model
vario = GeoVariogram(data=res$data,
                     coordx=coords, maxdist=0.3) # empirical semivariogram
GeoCovariogram(res, show.vario=TRUE, vario=vario, pch=20, ylim=c(0,2))
```

## Prediction of log-Gaussian random fields

The optimal linear prediction of log-Gaussian RF at a location  $\mathbf{s}_0$  is given by (Bevilacqua et al. (2020)):

$$\widehat{Y}(\mathbf{s}_0) = \widehat{\mu}(\mathbf{s}_0) \left( 1 + \sum_{i=1}^N \lambda_i [\widehat{V}(\mathbf{s}_i) - 1] \right) \quad (8)$$

where the vector of weights  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_N)'$  is given by  $\boldsymbol{\lambda} = R^{-1}\mathbf{c}$ .

Here  $\mathbf{c} = (\text{cor}(V(\mathbf{s}_0), V(\mathbf{s}_1)), \dots, \text{cor}(V(\mathbf{s}_0), V(\mathbf{s}_n)))'$  and  $R = [\text{cor}(V(\mathbf{s}_i), V(\mathbf{s}_j))]_{i,j=1}^N$  is the (estimated) correlation matrix associated to (2).

We first set the spatial locations to predict and the associated covariates. In this example, we choose a regular fine grid in order to construct a prediction map.

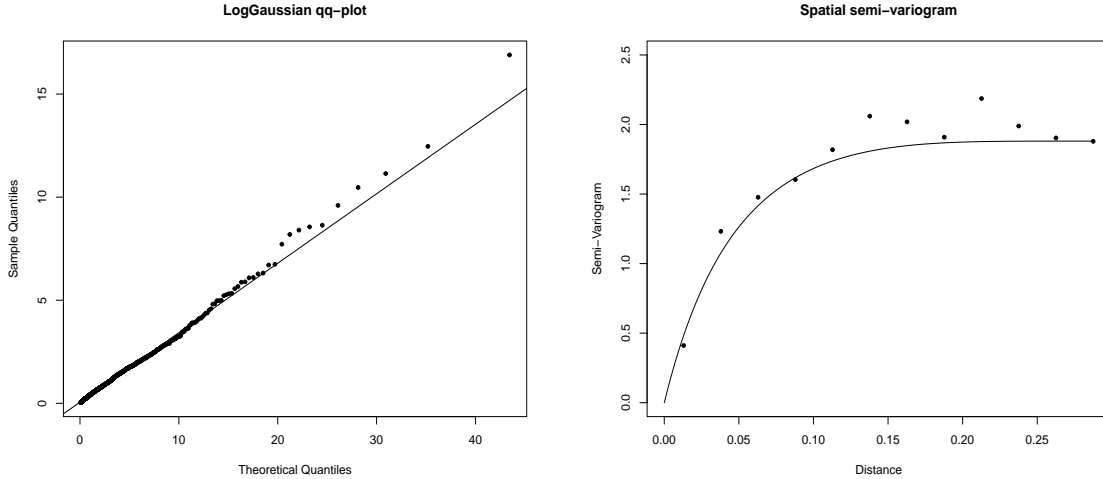


Figure 1: Left: QQ-plot for the log-Gaussian model residuals. Right: empirical vs estimated semi-variogram function for the residuals

```
# locations to predict and associated covariates
xx=seq(0,1,0.013)
loc_to_pred=as.matrix(expand.grid(xx,xx))
Nloc=nrow(loc_to_pred)
Xloc=cbind(rep(1,Nloc),runif(Nloc))
```

Then the optimal linear prediction (8), using the estimated parameters, can be performed using the `GeoKrig` function (computation can be time consuming):

```
param_est=as.list(c(fit$param,fixed))
pr=GeoKrig(data=data, coordx=coords,loc=loc_to_pred, X=X,Xloc=Xloc,
           corrmodel=corrmodel,model=model,mse=TRUE,
           sparse=TRUE,param=param_est)
```

and we can compare the map of simulated data with the kriging prediction (and associated mean square error) with the following code (see Figure 2):

```
colour = terrain.colors (20)
quilt.plot(x, y, data,col=colour,main="Data")
map=matrix(pr$pred,ncol=length(xx))
image.plot(xx, xx, map,col=colour,xlab="",ylab="",main="Kriging")
map_mse=matrix(pr$mse,ncol=length(xx))
image.plot(xx, xx, map_mse,col=colour,xlab="",ylab="",main="MSE")
```

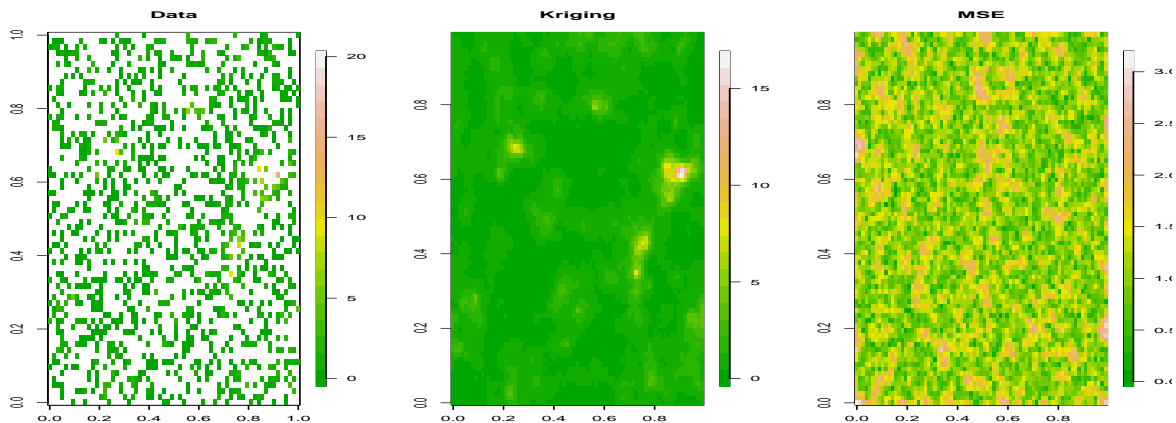


Figure 2: From left to right: colored map of observed data, kriging prediction and associated mean squared error

## References

- Bevilacqua, M., C. Caamaño-Carrillo, and C. Gaetan (2020). On modeling positive continuous data with spatiotemporal dependence. *Environmetrics* 31(7), e2632.
- 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, 828–856.
- Bevilacqua, M., V. Morales-Oñate, and C. Caamaño-Carrillo (2023). *GeoModels: A Package for Geostatistical Gaussian and non Gaussian Data Analysis*. R package version 1.1.0.
- Oliveira, V. D. (2006). On optimal point and block prediction in log-gaussian random fields. *Scandinavian Journal of Statistics* 33, 523–540.
- Oliveira, V. D., B. Kedem, and D. A. Short (1997). Bayesian prediction of transformed gaussian random fields. *Journal of the American Statistical Association* 92, 1422–1433.