

# GeoModels Tutorial: analysis of positive spatial data using Weibull random fields

Moreno Bevilacqua

## Introduction

In this tutorial we show how to analyze geo-referenced spatial data with positive support using Weibull random fields (RFs) (Bevilacqua et al., 2018) with the R package **GeoModels** (Bevilacqua and Morales-Oñate (2018)). The Weibull distribution is a flexible parametric model for positive data allowing both right and left 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(devtools)
install_github("vmoprojs/GeoModels")
require(GeoModels)
require(fields)
require(hypergeo)
model="Weibull" # model name in the GeoModels package
```

## Simulation of Weibull random fields

The definition of a Weibull 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  $\rho(\mathbf{h}) := \text{cor}(Z(\mathbf{s} + \mathbf{h}), Z(\mathbf{s}))$ .

Given  $Z_1, Z_2$ , two independent copies of  $Z$ , a RF  $U = \{U(\mathbf{s}), \mathbf{s} \in S\}$  with marginal distribution  $Weibull(\kappa, \nu(\kappa))$  can be derived by the transformation

$$U(\mathbf{s}) = \nu(\kappa) \left( \frac{1}{2} \sum_{k=1}^2 Z_k(\mathbf{s})^2 \right)^{1/\kappa}, \quad (1)$$

where  $\nu(\kappa) = \Gamma^{-1}(1 + 1/\kappa)$ ,  $\kappa > 0$  is a shape parameter and  $\Gamma(\cdot)$  is the gamma function. Under this specific parametrization,  $\mathbb{E}(U(\mathbf{s})) = 1$   $\text{var}(U(\mathbf{s})) = (\Gamma(1 + 2/\kappa) \nu^2(\kappa) - 1)$  and the correlation function is given by:

$$\rho_U(\mathbf{h}) = \frac{\nu^{-2}(\kappa)}{[\Gamma(1 + 2/\kappa) - \nu^{-2}(\kappa)]} [{}_2F_1(-1/\kappa, -1/\kappa; 1; \rho^2(\mathbf{h})) - 1]. \quad (2)$$

Here  ${}_2F_1(a, b; c; x)$  is the Gaussian hypergeometric function (Abramowitz and Stegun (1970)). In the **GeoModels** package it is computed using the function **hypergeo** of the **hypergeo** package (Hankin (2016)).

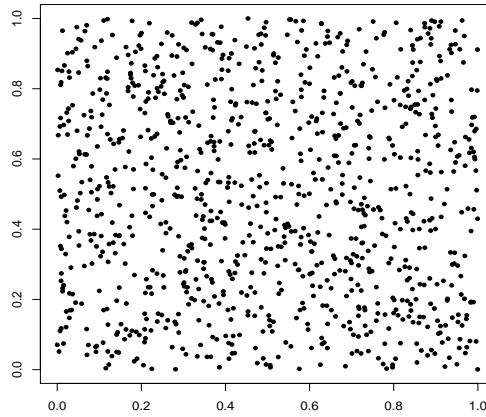
Then a non stationary version can be defined defined trough  $W = \{W(\mathbf{s}), \mathbf{s} \in S\}$  with

$$W(\mathbf{s}) = \mu(\mathbf{s})U(\mathbf{s}), \quad \mu(\mathbf{s}) > 0. \quad (3)$$

In this case  $\mathbb{E}(W(\mathbf{s})) = \mu(\mathbf{s})$ ,  $\text{var}(W(\mathbf{s})) = \mu(\mathbf{s})^2(\Gamma(1 + 2/\kappa)\nu^2(\kappa) - 1)$  and 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 Weibull RF we need to specify a regression mean parameters, a shape parameter and a parametric correlation model for  $\rho(\mathbf{h})$ . We first set the spatial coordinates

```
N=1000 # number of location sites
set.seed(24)
x = runif(N, 0, 1)
y = runif(N, 0, 1)
coords=cbind(x,y) # spatial coordinates
plot(coords,pch=20,xlab="",ylab="")
```



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 = -0.3; mean1=0.5 # regression parameters
```

where mean and mean1 are respectively  $\beta_1$  and  $\beta_2$ .

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

```
corrmodel = "Wend0"      ## correlation model and parameters
scale = 0.2
power2=4
```

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

Finally, we set the shape parameter of the Weibull RF and the nugget parameter

```
shape=2      # shape of the weibull RF
nugget=0     # nugget parameter
```

We are now ready to simulate a Weibull random field using the function `GeoSim`:

```
param=list(mean=mean,mean1=mean1,sill=1-nugget, nugget=nugget,
           scale=scale,power2=power2,shape=shape)
set.seed(312)
data = GeoSim(coordx=coords, corrmodel=corrmodel, model=model,
             param=param, X=X,sparse=TRUE)$data
```

The simulation is performed using Cholesky decomposition for the two Gaussian RFs involved. Note that the option `sparse=TRUE` allows to exploit specific algorithms for sparse matrices implemented in the `spam` package (Gerber et al. (2017)) when performing cholesky decomposition (Furrer and Sain (2010)).

## Estimation of Weibull random fields

The density of the bivariate random vector  $U(\mathbf{s}_i), U(\mathbf{s}_j)$  is given by (Bevilacqua et al. (2018)).

$$f_U(u_i, u_j) = \frac{\kappa^2(u_i u_j)^{\kappa-1}}{\nu^{2\kappa}(\kappa)(1 - \rho_{ij}^2)} \exp \left[ -\frac{u_i^\kappa + u_j^\kappa}{\nu^\kappa(\kappa)(1 - \rho_{ij}^2)} \right] I_0 \left( \frac{2|\rho|(u_i u_j)^{\kappa/2}}{\nu^\kappa(\kappa)(1 - \rho_{ij}^2)} \right). \quad (4)$$

where  $I_\alpha(x)$  denotes the modified Bessel function of the first kind of order  $\alpha$  and the bivariate densities of  $W$  can be derived from (4) as

$$f_W(w_i, w_j) = (\mu_i \mu_j)^{-1} f_U(w_i/\mu_i, w_j/\mu_j). \quad (5)$$

Given  $w(\mathbf{s}_1), \dots, w(\mathbf{s}_n)$ ,  $i = 1, \dots, n$  observations, then, the pairwise likelihood function is defined as:

$$pl(\boldsymbol{\theta}) = \sum_{i=1}^{N-1} \sum_{j=i+1}^N \log(f_W(w_i, w_j)) c_{ij}$$

where  $\boldsymbol{\theta} = (\beta_0, \beta_1, \kappa, \alpha, \delta)^T$ ,  $w_i = w(\mathbf{s}_i)$  for notation convenience and  $c_{ij}$  are non-negative weights, not depending on  $\boldsymbol{\theta}$ , specified as:

$$c_{ij} := \begin{cases} 1 & \|\mathbf{s}_i - \mathbf{s}_j\| < d \\ 0 & \text{otherwise} \end{cases}. \quad (6)$$

The 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. Pairwise likelihood estimation is performed with the function **GeoFit**:

```
start=list(mean=mean, mean1=mean1, scale=scale, shape=shape)
fixed=list(sill=1-nugget, nugget=nugget, power2=power2)
# Maximum composite-likelihood fitting of the Weibull random field:
fit = GeoFit(data=data, coordx=coords, corrmodel=corrmodel, model=model,
             X=X, optimizer="BFGS", start=start, fixed=fixed, maxdist=0.02)
```

The object **fit** include informations about the pairwise likelihood estimation

```
fit

#####
Maximum Composite-Likelihood Fitting of Weibull Random Fields
```

```

Setting: Marginal Composite-Likelihood
Model: Weibull
Type of the likelihood objects: Pairwise
Covariance model: Wend0
Optimizer: BFGS
Number of spatial coordinates: 1000
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: -706.47
Estimated parameters:
      mean      mean1      scale      shape
-0.3176    0.5195    0.1951    1.9678
#####

```

Note that the option `maxdist=0.02` set the compact support of the weight function (6) i.e.  $d = 0.02$ .

## 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

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

can be viewed as a realization of  $U$  a stationary RF with marginal distribution  $Weibull(\kappa, \nu(\kappa))$  with unit mean and correlation function given by (2).

The 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 a qq-plot:

```

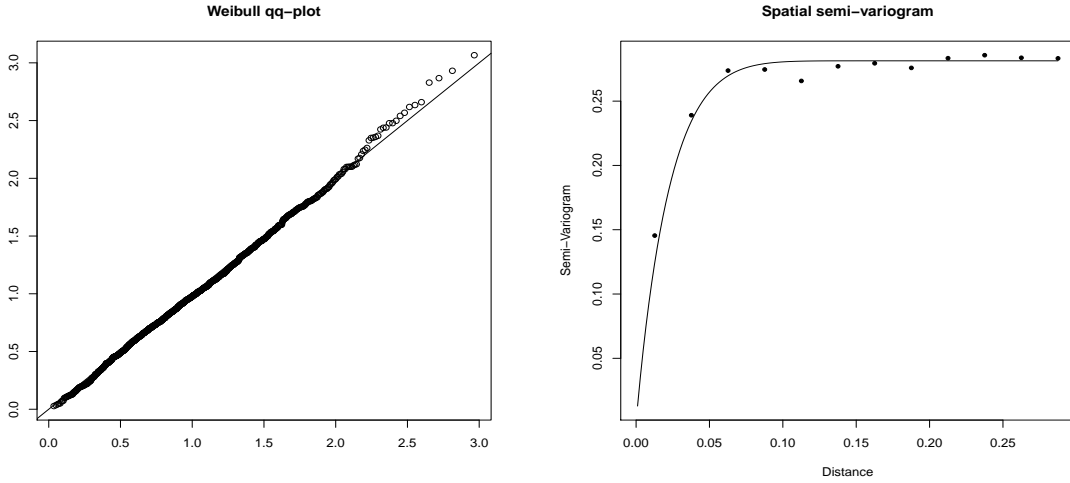
shape=fit$param["shape"]
probabilities = (1:N)/(N+1)
weibull.quantiles = qweibull(probabilities, shape=shape,
scale = 1/(gamma(1+1/shape)))

```

```
plot(sort(weibull.quantiles), sort(c(res$data)),
     xlab="", ylab="", main="Weibull qq-plot")
abline(0,1)
```

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

```
### checking model residuals assumptions: covariance model
vario = GeoVariogram(data=res$data,
                    coordx=coords, maxdist=0.3) # empirical variogram
GeoCovariogram(res, show.vario=TRUE, vario=vario, pch=20)
```



## Prediction of Weibull random fields

The optimal linear prediction of Weibull RF is given by (Bevilacqua et al. (2018))

$$\widehat{W(s_0)} = \widehat{\mu(s_0)} \left( 1 + \sum_{i=1}^N \lambda_i [\widehat{U(s_i)} - 1] \right) \quad (8)$$

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

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

We first set the spatial locations to predict and the associated covariates:

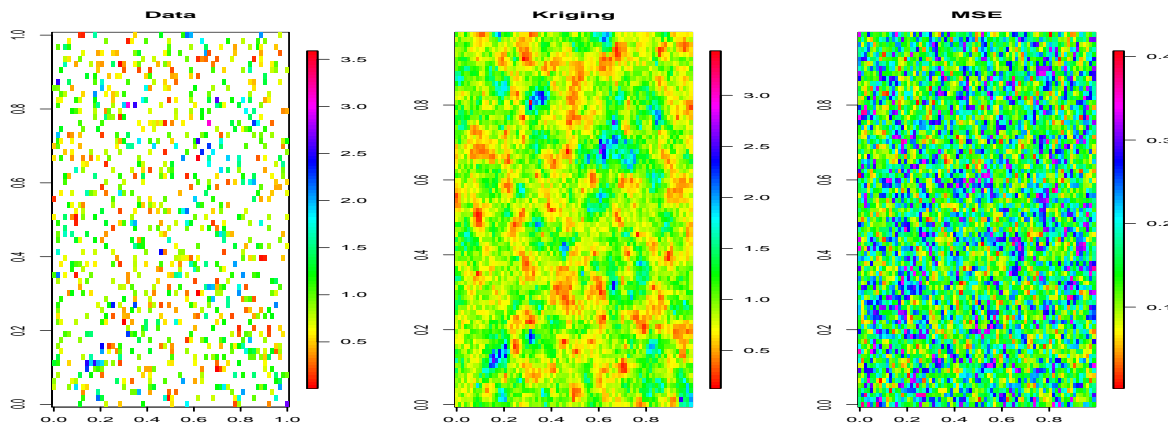
```
# locations to predict and associated covariates
xx=seq(0,1,0.013)
loc_to_pred=as.matrix(expand.grid(xx,xx))
Nloc=nrow(loc_to_pred)
Xloc=cbind(rep(1,Nloc),runif(Nloc))
```

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

```
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 predictions (and associated mean square error) with the following code:

```
colour = rainbow(100)
par(mfrow=c(1,3))
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")
```





## References

- Abramowitz, M. and I. A. Stegun (1970). *Handbook of Mathematical Functions*. New York: Dover.
- Bevilacqua, M., C. Caamano, and C. Gaetan (2018). On modelling positive continuous data with spatio-temporal dependence. *ArXiv e-prints*.
- 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. and V. Morales-Oñate (2018). *GeoModels: A Package for Geostatistical Gaussian and non Gaussian Data Analysis*. R package version 1.0.3-4.
- Furrer, R. and S. R. Sain (2010). spam: a sparse matrix R package with emphasis on mcmc methods for Gaussian Markov random fields. *Journal of Statistical Software* 36, 1–25.
- Gerber, F., K. Moesinger, and R. Furrer (2017). Extending R packages to support 64-bit compiled code: An illustration with spam64 and GIMMS NDVI3g data. *Computer & Geoscience* 104, 109–119.
- Hankin, R. K. S. (2016). *hypergeo: The Gauss Hypergeometric Function*. R package version 1.2-13.