

# **GeoModels Tutorial: simulation, estimation and prediction of global spatial data on the planet Earth using Gaussian random fields.**

Moreno Bevilacqua

## Introduction

This tutorial illustrates how to analyze spherical data (*i.e.* spatial data on the surface of a sphere) with the R package `GeoModels` (Bevilacqua and Morales-Oñate (2018)). Albeit the main focus of the tutorial is on Gaussian random fields (GRFs) defined over a sphere of arbitrary radius, any of the non GRFs implemented in the package `GeoModels` can be coupled with spherical data.

We denote with  $\mathbb{S}^2$  the sphere of  $\mathbb{R}^3$  of arbitrary radius  $R > 0$ , defined as

$$\mathbb{S}^2 = \{\mathbf{x} \in \mathbb{R}^3 : \|\mathbf{x}\| = R\}.$$

This tutorial shows how to simulate, estimate and predict GRFs on a sphere of arbitrary radius using a specific covariance model proposed in Alegría et al. (2020). which is a quite flexible for the reasons explained below. We start by loading the libraries needed for the analysis and set the radius of the sphere:

```
rm(list=ls())
library(GeoModels)
library(mapproj)
library(globe)
library(fields)
library(sphereplot)
set.seed(1891)
```

## Gaussian random fields defined on the sphere

We consider a zero mean and unit variance GRF,  $Y = \{Y(\mathbf{x}), \mathbf{x} \in \mathbb{S}^2\}$ , defined on a sphere of radius  $R > 0$ .

The great circle (GC) distance corresponds to the shortest path joining two points on the spherical surface and it is the natural distance to be used when analyzing data on a sphere. GC distance between two points  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{S}^2$  can be computed as:

$$d_{GC}(\mathbf{x}_i, \mathbf{x}_j) = R[\arccos(\mathbf{x}_i^T \mathbf{x}_j / R^2)].$$

Using the geometric definition of the dot product  $\mathbf{x}_i^T \mathbf{x}_j = \|\mathbf{x}_i\| \|\mathbf{x}_j\| \cos(\theta_{ij}) = R^2 \cos(\theta_{ij})$ , this implies that GC distance can be computed as  $d_{GC}(\mathbf{x}_i, \mathbf{x}_j) = R\theta_{ij}$  where  $\theta_{ij}$  is the angle between  $\mathbf{x}_i, \mathbf{x}_j$ .

A point  $\mathbf{x} \in \mathbb{S}^2$  can be parametrized through spherical coordinates using radians  $(lon^*, lat^*, R)$  where  $lon^* \in [-\pi, \pi]$  is the longitude and  $lat^* \in [-\pi/2, \pi/2]$  is the latitude. The connection between Euclidean and spherical coordinates is given by  $\mathbf{x} = (x_1, x_2, x_3)^T$  with  $x_1 = R\cos(lon^*)\cos(lat^*)$ ,  $x_2 = R\sin(lon^*)\cos(lat^*)$  and  $x_3 = R\sin(lat^*)$ .

Alternatively, one can use decimal degrees  $(lon, lat, R) = (180lon^*/\pi, 180lat^*/\pi, R)$  where  $lon \in [-180, 180]$  is the longitude and  $lat \in [-90, 90]$  is the latitude (the **GeoModels** package uses decimal degree).

To find the angle  $\theta_{ij}$  between  $\mathbf{x}_i, \mathbf{x}_j$ , we use the relation  $\theta_{ij} = \arccos(\mathbf{x}_i^T \mathbf{x}_j / R^2)$  and using spherical coordinates we obtain:

$$\theta_{ij} = [\arccos\{\sin lat_i^* \sin lat_j^* + \cos lat_i^* \cos lat_j^* \cos(lon_i^* - lon_j^*)\}]$$

Alternatively, one can compute the chordal (CH) distance, that is the segment below the arc joining any pair of points located over the spherical shell. It is given by:

$$d_{CH}(\mathbf{x}_i, \mathbf{x}_j) = 2R\sin(\theta_{ij}/2). \quad (1)$$

CH obviously underestimates the GC distance, and the approximation error increases with the size of the considered portion of the planet.

We assume  $\text{cor}(Y(\mathbf{x}_i), Y(\mathbf{x}_j)) = \rho(d_{GC}(\mathbf{x}_i, \mathbf{x}_j))$  where  $\rho : [0, R\pi] \rightarrow \mathbb{R}$  is continuous with  $\rho(0) = 1$ . Under this setting, the correlation function is called *geodesically isotropic* (Gneiting, 2013; Porcu et al., 2016). We then consider the location and scale transformation

$$Z(\mathbf{x}) = \mu(\mathbf{x}) + \sigma Y(\mathbf{x})$$

where  $\mu(\mathbf{x}) \in \mathbb{R}$  is a spherical varying mean and  $\sigma > 0$ . The package **GeoModels** allows a mean regression specification, that is  $\mu(\mathbf{x}) = M(\mathbf{x})^T \boldsymbol{\beta}$  where  $M(\mathbf{x})$  is a vector of spherical covariates. Nevertheless, for the sake of simplicity this tutorial assumes a constant mean:  $\mu(\mathbf{x}) = \mu$ .

Then  $\mathbb{E}(Z(\mathbf{x})) = \mu$  and  $\text{cov}(Z(\mathbf{x}_i), Z(\mathbf{x}_j)) = \sigma^2 \rho(d_{GC}(\mathbf{x}_i, \mathbf{x}_j))$  where  $\sigma^2$  is the variance parameter and in order to obtain a realization from a spherical RF we need to specify a mean, a variance parameter and a valid (*i.e.* positive definite) geodesically isotropic parametric correlation function  $\rho_\gamma(\cdot)$ . Since a sphere of radius  $R = 1$  is a (rescaled) satisfactory approximation of planet Earth we set the radius equal to 1 hereafter

```
radius=1
```

We set the spherical coordinates in decimal degrees on the unit sphere using the function `pointsphere` of the package `sphereplot` (Robotham, 2013)

```
NN=700 ## number of location sites on the sphere
coords=pointsphere(NN,c(-180,180),c(-90,90),c(radius,radius))[,1:2]
```

We can visualize (Figure 1) the locations on the planet Earth using some functions of the package `globe` (Baddeley et al., 2017) with the following code:

```
globeearth(eye=place("newyorkcity"))
globepoints(loc=coords,pch=20,cex=0.4)
globeearth(eye=place("everest"))
globepoints(loc=coords,pch=20,cex=0.4)
```



Figure 1: Location sites on the planet Earth used in this tutorial

An issue when considering correlation models defined on the sphere is that classical models defined on the euclidean spaces are not necessarily suitable on the sphere (Gneiting, 2013). For instance, the popular Matérn model equipped with the great circle distance defined as:

$$\mathcal{M}_{\nu,\alpha}(d_{GC}) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{d_{GC}}{\alpha} \right)^\nu \mathcal{K}_\nu \left( \frac{d_{GC}}{\alpha} \right),$$

is not valid on the sphere when  $\nu > 0.5$ . Here,  $\mathcal{K}_\nu$  is a modified Bessel function of the second kind of order  $\nu$ .

A interesting covariance model defined on the unit sphere having the same features of the Matérn model is the the so called  $\mathcal{F}$  family of correlation functions proposed in Alegria

et al. (2020), defined as:

$$\mathcal{F}_{\tau,\alpha,\nu}(d_{GC}) = \frac{B(\alpha, \nu + \tau)}{B(\alpha, \nu)} {}_2F_1(\tau, \alpha, \alpha + \nu + \tau; \cos(d_{GC})), \quad 0 \leq d_{GC} \leq \pi \quad (2)$$

where  $\tau, \alpha$  and  $\nu$  are strictly positive parameters and where  ${}_2F_1$  is the Gauss Hypergeometric function

$${}_2F_1(a, b, c; x) = \sum_{k=0}^{\infty} \frac{(a)_k (b)_k}{(c)_k} \frac{x^k}{k!}, \quad |x| < 1,$$

with  $(\cdot)_k$  being the Pochhammer symbol. Finally,  $B(\cdot, \cdot)$  is the Beta function. Alegría et al. (2020) show that (2) is  $k$  times differentiable at the origin if and only if  $\lfloor \nu/2 \rfloor > k$ , for any  $\tau, \alpha$ , and where  $\lfloor x \rfloor$  denotes the largest integer being smaller than  $x \in \mathbb{R}$ . From this point of view the  $\mathcal{F}$  class of correlation can be viewed as the analogous of the Matérn class on the sphere. In the **GeoModels** package the  $\mathcal{F}$  family is implemented as

$$\mathcal{F}_{1/\alpha, 1/\alpha+0.5, \nu}(d_{GC}) = \frac{B(1/\alpha + 0.5, \nu + 1/\alpha)}{B(1/\alpha + 0.5, \nu)} {}_2F_1(1/\alpha, 1/\alpha + 0.5, 2/\alpha + 0.5 + \nu; \cos(d_{GC})), \quad (3)$$

This kind of parametrization allows to solve problem of identifiability associated to (2) and the  $\alpha$  and  $\nu$  parameters are correlation range and smoothness parameters respectively as in the Matérn model.

We set the  $\mathcal{F}$  correlation model with associated parameters and the other parameters of the GRF. In the **GeoModels** package the  $\mathcal{F}$  model is implemented with the name of **Smoke** model.

```
corrmodel = "Smoke"      ## correlation model and parameters
scale=0.2
smooth=0.5
sill=1
nugget=0
mean=0
```

Here the **scale** parameter corresponds to  $\alpha$ , **sill** is the variance parameter  $\sigma^2$  and **smooth** is the smoothness parameter  $\nu$ . The choice of  $\nu = 0.5$  leads to a non mean square differentiable Gaussian random field. We are now ready to simulate a GRF on the unit sphere using the function **GeoSim**:

```
param=list(mean=mean, sill=sill, nugget=nugget, smooth=smooth,
```

```

    scale=scale)
data = GeoSim(coordx=coords, corrmodel=corrmodel, param=param,
               distance="Geod", radius=radius)$data

```

Note that the option `distance="Geod"` coupled with the location sites given in lon/lat format in decimal degrees (`coordx=coords`) and a valid correlation model on the sphere (`corrmodel`) with a given `radius`, allows to perform a simulation based on Cholesky decomposition of a GRF on  $\mathbb{S}^2$ .

## Estimation of Gaussian random fields on the sphere

Let  $\rho_{\alpha,\nu}(d_{GC}) = \mathcal{F}_{1/\alpha, 1/\alpha+0.5, \nu}(d_{GC})$ . Given a realization  $\mathbf{Z} = (Z(\mathbf{x}_1), Z(\mathbf{x}_2), \dots, Z(\mathbf{x}_N))^T$  from a GRF defined on  $\mathbb{S}^2$  with correlation (3) we can perform maximum likelihood estimation maximizing the log-likelihood function

$$l(\boldsymbol{\beta}) = -0.5\log(|\sigma^2 C_{\alpha,\nu}|) - 0.5 \frac{(\mathbf{Z} - \mu\mathbf{1})^T C_{\alpha,\nu}^{-1} (\mathbf{Z} - \mu\mathbf{1})}{\sigma^2}$$

with respect to  $\boldsymbol{\beta} = (\mu, \sigma^2, \alpha, \nu)^T$  where  $C_{\alpha,\nu} = [\rho_{\alpha,\nu}(d_{GC}(\mathbf{x}_i, \mathbf{x}_j))]_{i,j=1}^N$  is the correlation matrix.

We can use the `GeoFit` function to perform maximum likelihood estimation. We first choose the fixed parameters and the parameters that must be estimated.

```

fixed<-list(nugget=nugget)
start<-list(mean=mean, scale=scale, sill=sill, smooth=smooth)

```

We consider a quasi-Newton optimization as implemented in the `nlminb` function that allows box-constrained optimization using PORT routines. We first fix the lower and upper bound of the constrained optimization:

```

optimizer="nlminb"
I=Inf
upper<-list(mean=I, scale=I, sill=I, smooth=I)
lower<-list(mean=-I, scale=0, sill=0, smooth=0)

```

We are now ready to perform maximum likelihood estimation using the `GeoFit` function. This option can be time consuming depending on the number of point on the sphere.

```

fixed<-list(nugget=nugget)
start<-list(mean=mean, scale=scale, sill=sill, smooth=smooth)

```

```

fit_geo_ml <- GeoFit(data=data, coordx=coords, corrmodel=corrmodel,
                      likelihood="Full", type="Standard",
                      optimizer=optimizer, lower=lower, upper=upper,
                      start=start, fixed=fixed, distance="Geod", radius=radius)

```

Default maximization algorithm used in the `GeoFit` function is the Nelder-mead method as implemented in the R function `optim`.

Since maximum likelihood can be computationally expensive, in particular when the computation of the correlation function involves some special functions as in the Matérn or  $\mathcal{F}$  cases, we can use conditional pairwise likelihood estimation (Bevilacqua and Gaetan, 2015). obtained by maximizing the function

$$pl(\boldsymbol{\beta}) = \sum_{i=1}^N \sum_{j=1, j \neq i}^N \log(f_{Z_i, Z_j}(z_i, z_j; \boldsymbol{\beta})) - (\log(f_{Z_i}(z_i; \boldsymbol{\beta})) + \log(f_{Z_j}(z_j; \boldsymbol{\beta})))w_{ij} \quad (4)$$

with respect to  $\boldsymbol{\beta}$ . It can be performed with the following code

```

fit_geo_pl <- GeoFit(data=data, coordx=coords, corrmodel=corrmodel,
                      likelihood="Conditional", type="Pairwise", neighb=3,
                      optimizer=optimizer, lower=lower, upper=upper,
                      start=start, fixed=fixed, distance="Geod", radius=radius)

```

In (4)  $f_{Z_i, Z_j}(z_i, z_j; \boldsymbol{\beta})$  and  $f_{Z_j}(z_j; \boldsymbol{\beta})$  are the Gaussian bivariate and marginal density respectively and  $w_{ij}$  are non-negative weights, not depending on  $\boldsymbol{\beta}$ , specified as:

$$w_{ij}(m) = \begin{cases} 1 & \mathbf{x}_i \in N_m(\mathbf{x}_j) \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where  $N_m(\mathbf{s}_l)$  is the set of neighbours of order  $m = 1, 2, \dots$  of the point  $\mathbf{s}_l \in A$ .

The two estimates can be obtained from the objects `fit_geo_ml` and `fit_geo_pl` as:

```

unlist(fit_geo_ml$param)
  mean      scale      sill      smooth
-0.07685842  0.16698323  0.93005748  0.55309468

unlist(fit_geo_pl$param)
  mean      scale      sill      smooth
-0.09594479  0.20088214  0.93402952  0.48057181

```

Recall that the Matérn model is not a suitable covariance model on the sphere but it is still a valid model when coupled with the chordal distance. In the `GeoModels` package,

estimation can be performed using also the CH distance, or the Euclidean distance after a suitable projection of the spherical coordinates. Computation of CH distances involved in the covariance matrix can be obtained using (1), that is:

Then we can use conditional pairwise likelihood estimation using chordal distances (setting `distance="Chor"`) in order to estimate a Matérn covariance model:

```
fit_chor_pl <- GeoFit(data=data, coordx=coords, corrmodel="Matern",
                      likelihood="Conditional", type="Pairwise", neighb=3,
                      optimizer=optimizer, lower=lower, upper=upper,
                      start=start, fixed=fixed, distance="Chor", radius=radius)
```

We can also consider conditional pairwise likelihood estimation of the Matérn model using a suitable projection of the spherical points coupled with the Euclidean distance. We use the function `mapproj` of the package `mapproj` (McIlroy et al., 2017) in order to obtain projected coordinates using sinusoidal projection

```
prj=mapproject(coords[,1], coords[,2], projection="sinusoidal",
                orientation=c(90,0,0))
coords_prj=radius*cbind(prj$x, prj$y)
```

and we can visualize the projected points (Figure 2) with the following code:

```
sinusoidal.proj = map(database= "world", ylim=c(-90,90),
                      xlim=c(-180,180), col="grey80", fill=TRUE,
                      plot=FALSE, projection="sinusoidal")
map(sinusoidal.proj)
points(coords_prj, pch=20, cex=0.4)
```

Other types of projections can be chosen in the package `mapproj` such us, for instance, the mercator or cylindrical projections.

Then, we can estimate with conditional pairwise likelihood using Euclidean distances:

```
fit_eucl_pl <- GeoFit(data=data, coordx=coords_prj, corrmodel="Matern",
                        neighb=3, likelihood="Conditional", type="Pairwise",
                        optimizer=optimizer, lower=lower, upper=upper,
                        start=start, fixed=fixed, distance="Eucl", radius=radius)
```

The two estimates of the Matérn model using chordal and euclidean distances can be obtained from the objects `fit_chor_ml` and `fit_uecl_pl` as:

```
unlist(fit_chor_pl$param)
```

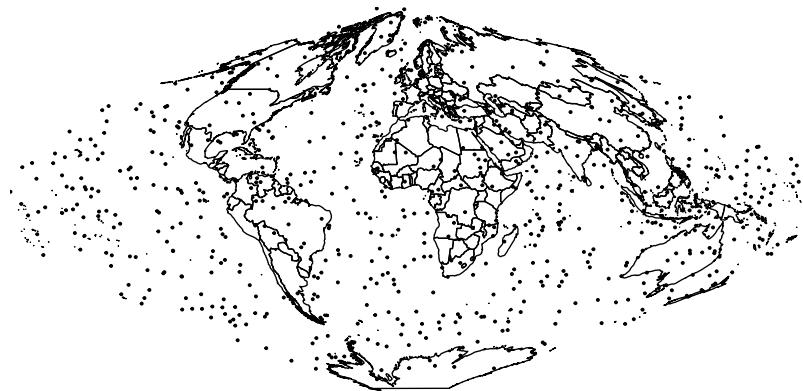


Figure 2: Sinusoidal projection

```

    mean      scale      sill      smooth
-0.09845857  0.15492058  0.95646508  0.46724768
unlist(fit_eucl_pl$param)
    mean      scale      sill      smooth
-0.1045473   0.1448448   0.9546144   0.4221286

```

Boxplots of the distances (Figure 4) can be useful to depict the different behaviour of the three type of distances in this specific example:

```

geod_ds=rdist.earth(coords,miles=F,R=1)
chor_ds=2*sin(geod_ds/2)
eucl_ds=dist(coords_prj)
boxplot(c(geod_ds),c(chor_ds),c(eucl_ds),
        names=c("Greatcircle","Chordal","Euclidean"))

```

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

$$\widehat{Y}(\mathbf{x}_i) = \frac{Z(\mathbf{x}_i) - \hat{\mu}}{(\hat{\sigma}^2)^{\frac{1}{2}}} \quad i = 1, \dots, N$$

can be viewed as a realization of a zero mean unit variance spherical GRF with correlation model (3). Using conditional pairwise likelihood estimates, the residuals can be computed using the `GeoResiduals` function:

```
res=GeoResiduals(fit_geo_pl)
```

Then the marginal distribution assumption on the residuals can be graphically checked with a qq-plot (Figure 3, left part) using the function `GeoQQ`.

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

The correlation model assumption can be graphically checked comparing the empirical and the estimated semivariogram functions using the `GeoVariogram` and `GeoCovariogram` functions (Figure 3, right part):

```
### checking model assumptions: variogram model
vario = GeoVariogram(data=res$data,coordx=coords,
                      maxdist=max(geod_ds)/2,distance="Geod",radius=radius)
GeoCovariogram(res,vario=vario,show.vario=TRUE,pch=20)
```

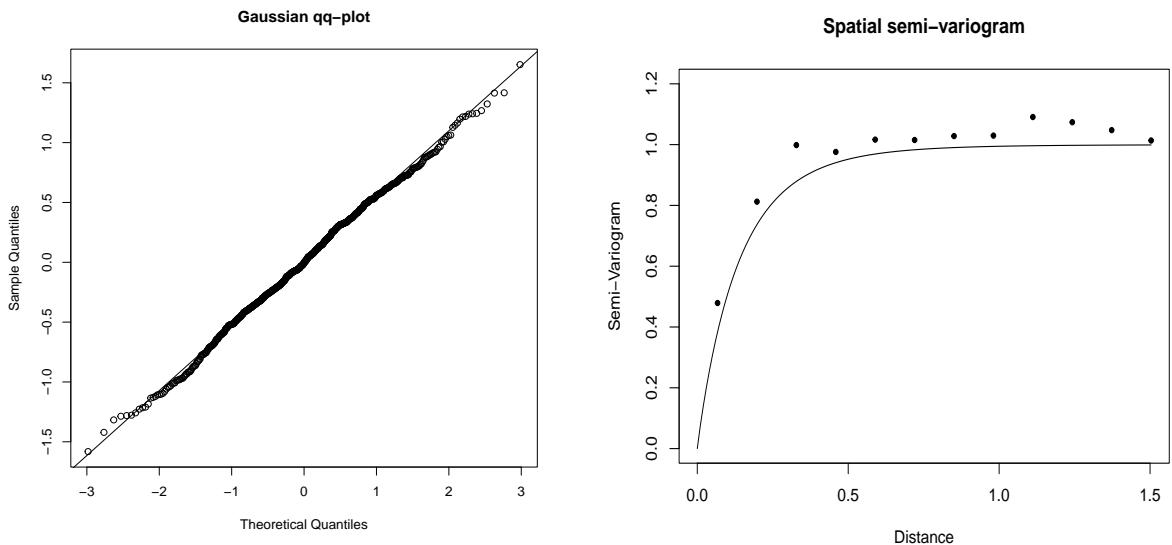


Figure 3: Left: QQ-plot for the model residuals. Right: empirical vs estimated semi-variogram function for the residuals

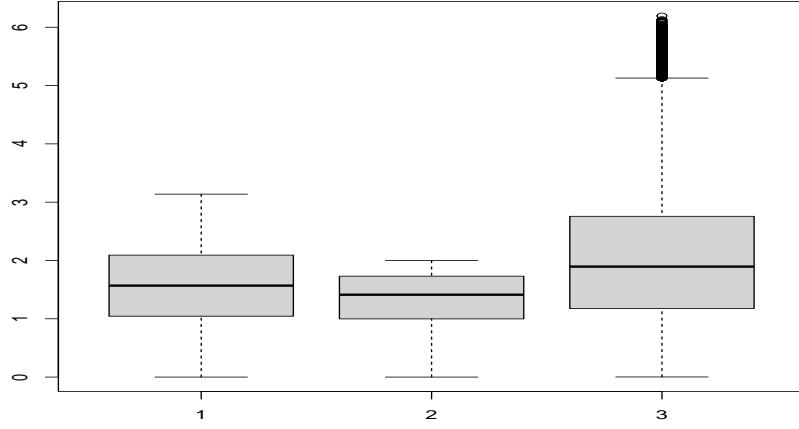


Figure 4: Boxplot of great circle, chordal and euclidean distances (from left to right)

## Prediction of Gaussian random fields on the sphere

For a given location  $\mathbf{x}_0 \in \mathbb{S}^2$ , the optimal prediction in this example is computed by:

$$\widehat{Z}(\mathbf{x}_0) = \hat{\mu} + \sum_{i=1}^N \hat{\lambda}_i [Z(\mathbf{x}_i) - \hat{\mu}] \quad (6)$$

where  $\hat{\mu}$  is a mean estimation and the vector of weights  $\hat{\boldsymbol{\lambda}} = (\hat{\lambda}_1, \dots, \hat{\lambda}_N)'$  is given by  $\hat{\boldsymbol{\lambda}} = \widehat{C}_{\hat{\alpha}, \hat{\nu}}^{-1} \hat{\mathbf{c}}_{\hat{\alpha}, \hat{\nu}}$  where  $\widehat{C}_{\hat{\alpha}, \hat{\nu}}$  is the estimated correlation matrix and  $\hat{\mathbf{c}}_{\hat{\alpha}, \hat{\nu}}$  is the estimated correlation between the locations and the point to predict. Moreover, the associated mean squared error is computed through

$$MSE(\widehat{Z}(\mathbf{x}_0)) = \hat{\sigma}^2 (1 - \hat{\mathbf{c}}_{\hat{\alpha}, \hat{\nu}}^T \widehat{C}_{\hat{\alpha}, \hat{\nu}}^{-1} \hat{\mathbf{c}}_{\hat{\alpha}, \hat{\nu}}). \quad (7)$$

Optimal prediction and associated MSE can be computed using the `GeoKrig` function. Suppose we want to predict the GRF around the north pole area. First we need to specify the spherical location to predict (in lon/lat format):

```
## location to predict given in lon/lat format
loc_to_pred=as.matrix(expand.grid(seq(-180,180,2.5), seq(45,90,0.15)))
```

To speed up the computation of (6) and (7) the optimal local predictor using 100 neighbours and associated MSE (using the maximum conditional pairwise likelihood estimated parameters) can be performed with the following code:

```

param_est<-as.list(c(fixed,fit_geo_pl$param))
pr_geo=GeoKrigloc(data=data,loc=loc_to_pred,coordx=coords,corrmodel=corrmodel,
    radius=radius,distance="Geod",param=param_est,mse=TRUE,neighb=100)

```

The matrix of locations to predict, predictions and associated MSE can be obtained as:

```

predictions=cbind(loc_to_pred,pr_geo$pred,pr_geo$mse)
head(predictions)
  Var1  Var2
[1,] -180.0 45 -0.2590250 0.5359307
[2,] -177.5 45 -0.3447504 0.5179728
[3,] -175.0 45 -0.4610383 0.5490006
[4,] -172.5 45 -0.5853897 0.5756557
[5,] -170.0 45 -0.6889051 0.5721893
[6,] -167.5 45 -0.7494889 0.5335922

```

and a graphical visualization (Figure 5) of the prediction and associated MSE can be obtained with the following code:

```

globeearth(eye=place("newyorkcity"));
globepoints(loc=loc_to_pred,pch=20,
col = terrain.colors(length(pr_geo$pred),alpha=0.03)[rank(pr_geo$pred)])

globeearth(eye=place("newyorkcity"));
globepoints(loc=loc_to_pred,pch=20,
col = cm.colors(length(pr_geo$mse),alpha=0.03)[rank(pr_geo$mse)])

```

## References

Alegría, A., F. Cuevas, P. Diggle, and E. Porcu (2020). A family of covariance functions for random fields on spheres. <https://math.au.dk/forskning/publikationer/institutts-serier/publication/publid/1126/>.

Baddeley, A., T. Lawrence, and E. Rubak (2017). *globe: Plot 2D and 3D Views of the Earth, Including Major Coastline.* R package version 1.2-0.

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



Figure 5: From left to right: prediction and associated MSE in the north pole area

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.

Douglas Nychka, Reinhard Furrer, John Paige, and Stephan Sain (2015). *fields: Tools for spatial data*. R package version 9.0.

Gneiting, T. (2013). Strictly and non-strictly positive definite functions on spheres. *Bernoulli* 19, 1327–1349.

McIlroy, D., R. Brownrigg, T. Minka, and R. Bivand. (2017). *mapproj: Map Projections*. R package version 1.2-5.

Porcu, E., M. Bevilacqua, and M. Genton (2016). Spatio-temporal covariance and cross covariance functions of the great circle distance on a sphere. *Journal of the American Statistical Association* 111, 888–898.

Robotham, A. (2013). *sphereplot: Spherical plotting*. R package version 1.5.