

GeoModels Tutorial: analysis of spatial precipitation anomalies using Gaussian and skew Gaussian random fields

Moreno Bevilacqua, Christian Caamaño-Carrillo

September 12, 2021

Introduction

In this tutorial we show how to analyze total precipitation anomalies registered at 7,352 location sites in the USA from 1895 to 1997. A detailed description of the data can be found in Kaufman et al. (2008). The yearly totals have been standardized by the long-run mean and standard deviation for each station from 1962. The size of dataset is large and this tutorial shows how to use the package `GeoModels` in order to perform estimation and prediction using Gaussian and SkewGaussian random fields. We first load the *R* libraries needed in this tutorial.

```
require(devtools)
install_github("vmoprojs/GeoModels")
require(GeoModels)
require(fields)
require(maps)
require(maptools)
require(mapdata)
require(geoR)
require(sn)
library(mapproj)
```

1 Preliminary data analysis

Precipitation anomalies data can be found in the `GeoModels` package. We first import the data:

```
data(anomalies)
head(anomalies)
      lon  lat      z
[1,] -85.25 31.57 -0.4586873
[2,] -87.42 32.23 -0.9253283
[3,] -85.87 32.98 -0.4370817
[4,] -88.13 33.13 -0.6026716
[5,] -86.50 31.32 -0.3519950
[6,] -85.85 33.58  0.5069722
```

and we select the coordinates (given in lon/lat format, decimal degree) and the anomalies data

```
loc=cbind(anomalies[,1],anomalies[,2])
z=cbind(anomalies[,3])
```

A colour map of the anomalies data can be obtained with the following code (see Figure 1).

```
quilt.plot(loc,z,xlab="long",ylab="lat")
map("usa", add = TRUE)
```

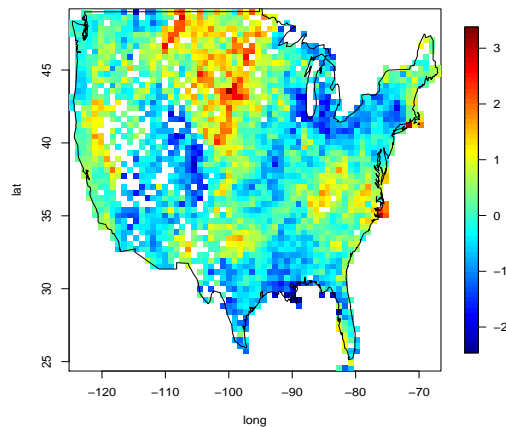


Figure 1: Coloured map of anomalies data.

We first project the spherical coordinates on a two dimensional euclidean space using a projection method. In this example we select a sinusoidal projection. The `GeoModels` package allows to handle spherical coordinates given in lon/lat format (decimal degree) and to work with geodesic or chordal distances. However in this example we work with projected coordinates.

```
P.sinusoidal <- mapproject(loc[,1],loc[,2],projection="sinusoidal")
loc<-cbind(P.sinusoidal$x,P.sinusoidal$y)*6371
maxdist=max(dist(loc))
```

Here 6371 is the radius of the earth in KM. The marginal distribution of the data (see the histogram in Figure 2 left part) suggests that the marginal Gaussian assumption seems quite reasonable. However a skew-Gaussian distribution could be more appropriate since

it can be appreciated a slight degree of asymmetry. Additionally the h -scatterplot suggests an elliptical dependence for the bivariate distributions (see Figure, 2 left part).

```
hist(z,main="Anomalies histogram")
GeoScatterplot(data=z,coordx=loc,maxdist=50,numbins=4)
```

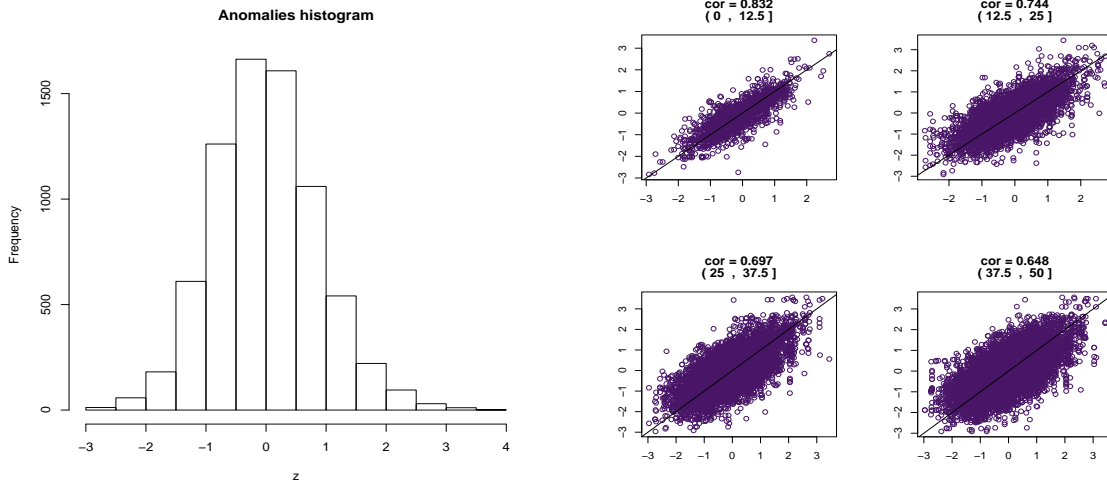


Figure 2: From left to right: histogram of anomalies data and associated h -scatterplot.

Finally, the empirical semivariogram in Figure 3 suggests the presence of a non-negligible nugget effect.

```
evariog=GeoVariogram(data=z,coordx=loc,maxdist=maxdist/4)
plot(evariog,ylim=c(0,1), pch=20,xlab="Km",ylab="Semi-variogram")
```

This preliminary graphical analysis suggest the a Gaussian or a skew-Gaussian random field with a covariance model including a nugget effect can be suitable models for the anomalies data.

2 Gaussian and skew-Gaussian random fields

We first consider a zero mean, unit variance and weakly stationary standard Gaussian random field $G = \{G(\mathbf{s}), \mathbf{s} \in S\}$, where \mathbf{s} represents a location in the domain $S \subset \mathbb{R}^2$ with isotropic exponential correlation model with a nugget effect that is:

$$\rho(\mathbf{h}) = \begin{cases} 1 & \|\mathbf{h}\| = 0 \\ (1 - \tau^2)e^{-\|\mathbf{h}\|/\alpha} & \|\mathbf{h}\| > 0 \end{cases}. \quad (1)$$

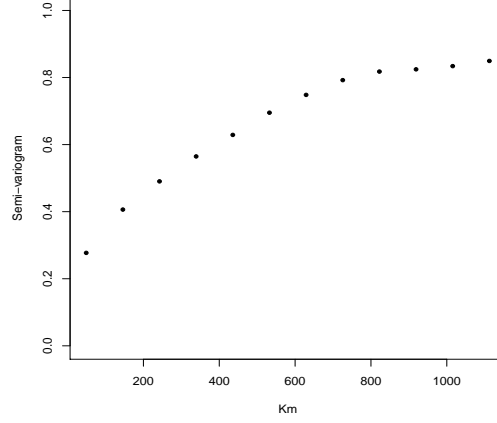


Figure 3: Empirical semi-variogram of anomalies data.

Here $\|\mathbf{h}\|$ is the Euclidean distance and $0 \leq \tau^2 < 1$ represents the nugget parameter. Additionally $\alpha > 0$ is a spatial dependence parameter.

We consider two random fields in our analysis. The first is a location-scale transformation of G that is a random field $Y = \{Y(\mathbf{s}), \mathbf{s} \in A\}$ defined as:

$$Y(\mathbf{s}) := \mu(\mathbf{s}) + \sigma G(\mathbf{s}) \quad (2)$$

with $\mathbb{E}(Y(\mathbf{s})) = \mu(\mathbf{s}) \in \mathbb{R}$ and $Var(Y(\mathbf{s})) = \sigma^2 > 0$.

The second is a location-scale transformation of the skew Gaussian random field proposed in Zhang and El-Shaarawi (2010) that is:

$$U_\eta(\mathbf{s}) = \mu(\mathbf{s}) + \sigma \left(\frac{\eta}{\sigma} |G_1(\mathbf{s})| + G_2(\mathbf{s}) \right) \quad (3)$$

with $\mathbb{E}(U_\eta(\mathbf{s})) = \mu(\mathbf{s}) + \eta\sqrt{2/\pi}$ and $Var(U_\eta(\mathbf{s})) = \sigma^2 + \eta^2(1 - 2/\pi)$ where $\eta \in \mathbb{R}$ is the asymmetry parameter, $\sigma > 0$ and G_i $i = 1, 2$ are two independent copies of a process G . More precisely, G_1 is a Gaussian random field with correlation (1) (assuming zero nugget) and G_2 is an independent Gaussian random field with correlation (1). Note that if $\eta = 0$ the Gaussian random field in (2) is obtained. As a consequence (3) is a generalization of (2). The correlation function of the Skew-Gaussian random field is given by (Zhang and El-Shaarawi, 2010)

$$\rho_{U_\eta}(\mathbf{h}) = \frac{2\eta^2}{\pi\sigma^2 + \eta^2(\pi - 2)} \left((1 - \rho_1^2(\mathbf{h}))^{1/2} + \rho_1(\mathbf{h}) \arcsin(\rho_1(\mathbf{h})) - 1 \right) + \frac{\sigma^2 \rho(\mathbf{h})}{\sigma^2 + \eta^2(1 - 2/\pi)}. \quad (4)$$

where $\rho_1(\mathbf{h})$ is the correlation function in (1) with $\tau^2 = 0$.

For both random fields we assume a constant mean $\mu(\mathbf{s}) = \mu$. However, the `GeoModels` package allows to specify a model regression for the spatial mean.

To obtain the names of the correlation parameters of the correlation model (1) and the names of the nuisance parameters of the Gaussian and Skew-Gaussian models, two useful functions are `CorrParam` and `NuisParam`:

```
CorrParam("Matern")
[1] "scale" "smooth"
NuisParam("Gaussian")
[1] "mean" "nugget" "sill"
NuisParam("SkewGaussian")
[1] "mean" "nugget" "sill" "skew"
```

For the exponential model we consider a Matern model where `scale` is the α parameter in equation (1) and `smooth` is a smoothness parameter that will be fixed equal to 0.5. Finally `nugget` is the τ^2 parameter, `sill` is the σ^2 parameter, `skew` is the η parameter and `mean` is the constant mean parameter μ .

Estimation of Anomalies data

Given a realization $\mathbf{Y} = (y(\mathbf{s}_1), y(\mathbf{s}_2), \dots, y(\mathbf{s}_N))^T$ from a Gaussian random field with correlation (1), the estimation of the parameters can be performed using maximum likelihood method that is maximizing the Gaussian multivariate pdf

$$f_{\mathbf{Y}}(y_1, \dots, y_N; \boldsymbol{\theta}_Y) = (2\pi)^{-N/2} |\sigma^2 R|^{-1/2} \exp \left\{ -\frac{(\mathbf{Y} - \mu \mathbf{1}_N)^T R^{-1} (\mathbf{Y} - \mu \mathbf{1}_N)}{2\sigma^2} \right\} \quad (5)$$

with respect to $\boldsymbol{\theta}_Y = (\mu, \sigma^2, \alpha, \beta, \tau^2)^T$. Here $R = [\rho(\mathbf{s}_i - \mathbf{s}_j)]_{i,j=1}^N$ is the correlation matrix.

However, since maximum likelihood method is computationally expensive for large datasets, in this example we focus on composite likelihood estimation method based on pairs. Specifically we focus on the conditional pairwise likelihood method as described in Bevilacqua and Gaetan (2015). This method involve only the bivariate and marginal pdfs associated with the bivariate random vector $\mathbf{Y}_{ij} = (Y(\mathbf{s}_i), Y(\mathbf{s}_j))^T$.

Similarly, given a realization $\mathbf{U}_\eta = (u_\eta(\mathbf{s}_1), u_\eta(\mathbf{s}_2), \dots, u_\eta(\mathbf{s}_N))^T$, from a skew-Gaussian random field the conditional pairwise likelihood method involves the bivariate and marginal

pdfs of the bivariate random vector $\mathbf{U}_{\eta;ij} = (U_{\eta}(\mathbf{s}_i), U_{\eta}(\mathbf{s}_j))^T$. In particular the bivariate pdf is given by (Alegria et al., 2017):

$$f_{\mathbf{U}_{\eta;ij}}(u_i, u_j; \boldsymbol{\theta}_{U_{\eta}}) = 2 \sum_{l=1}^2 \phi_2(\mathbf{u}_{ij} - \boldsymbol{\mu}_{ij}; \mathbf{A}_l) \Phi_n(\mathbf{c}_l; \mathbf{0}, \mathbf{B}_l) \quad (6)$$

where $\boldsymbol{\theta}_{U_{\eta}} = (\mu, \sigma^2, \eta, \alpha, \beta, \tau^2)^T$, $\boldsymbol{\mu}_{ij} = (\mu, \mu)^T$ and \mathbf{A}_l , \mathbf{B}_l , \mathbf{c}_l are specific quantities depending on the correlation and the parameters (see Alegria et al. (2017) for details).

The conditional pairwise likelihood function associated to Y is given by

$$pl(\boldsymbol{\theta}_Y) = \sum_{i=i}^N \sum_{j=1, j \neq i}^N [\log(f_{Y_{ij}}(y_i, y_j)) - \log(f_{Y_j}(y_j))] w_{ij} \quad (7)$$

and the conditional pairwise likelihood function associated to U_{η} is given by

$$pl(\boldsymbol{\theta}_{U_{\eta}}) = \sum_{i=i}^N \sum_{j=1, j \neq i}^N [\log(f_{\mathbf{U}_{\eta;ij}}(u_i, u_j)) - \log(f_{\mathbf{U}_{\eta;j}}(u_j))] w_{ij} \quad (8)$$

where w_{ij} are non-negative non-symmetric weights specified as:

$$w_{ij}(k) = \begin{cases} 1 & \mathbf{s}_i \in N_k(\mathbf{s}_j) \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

Here $N_k(\mathbf{s}_l)$ is the set of the neighbors of order k of the point \mathbf{s}_l . This kind of weights are computationally convenient since kd-tree type algorithms can be used to find exact or approximate nearest neighbors of order k for each location sites, that is $N_k(\mathbf{s}_l)$ for $l = 1, \dots, n$ with a computational complexity of order $N \log(N)$ and associated storage of order N . In particular the **GeoModels** package, exploits the functions implemented in the package **RANN**.

The conditional pairwise likelihood estimator $\hat{\boldsymbol{\theta}}_{pl}$ of the Gaussian and skew-Gaussian random fields is obtained maximizing (7) and (8) with respect to $\boldsymbol{\theta}_Y$ and $\boldsymbol{\theta}_{U_{\eta}}$ respectively.

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 **GeoFit2**. This function performs a preliminary estimation of the marginal parameters under the assumption of independence (the so-called independence composite likelihood) and then compute conditional pairwise likelihood estimation optimization using the estimates at first step as starting value.

In this example, we perform optimization of (7) and (8) using the function **nlminb** that allows for box-constrained optimization. However other type of optimization algorithms

can be used (BFGS or Nelder-Mead for instance). We use the following code to estimate the parameters θ_Y of the Gaussian random field (the option `neighbors` set the order of neighbors k in the weight function (9) and the `lower` and `upper` objects define the upper and lower bounds of the parametric space).

```
I=Inf
lower=list(mean=-I,sill=0,nugget=0,scale=0)
upper=list(mean=I,sill=I,nugget=1,scale=I)
start=list(mean=mean(z),sill=var(z),nugget=0.10,scale=140)
fixed=list(smooth=0.5)
## estimation
pcl1=GeoFit2(coordx=loc,corrmodel=corrmodel,data=z,
likelihood="Conditional",type="Pairwise",model="Gaussian",
optimizer="nllminb",lower=lower,upper=upper, neighbors=5,
start=start,fixed=fixed)
```

The object `pcl1` include information about the pairwise likelihood estimation:

```
pcl1
#####
Maximum Composite-Likelihood Fitting of Gaussian Random Fields
Setting: Conditional Composite-Likelihood
Model: Gaussian
Type of the likelihood objects: Pairwise
Covariance model: Matern
Optimizer: nllminb
Number of spatial coordinates: 7352
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: -33110.76
Estimated parameters:
      mean      nugget      scale      sill
    0.02616    0.07676   124.98471    0.79096
#####
```

We now estimate the the skew-Gaussian random field.

```
start=list(nugget=0,scale=100,mean=mean(z),sill=var(z),skew=0.4)
```



```

fixed=list(smooth=.5)
lower=list(mean=-I,sill=0,skew=-I,scale=0,nugget=0)
upper=list(mean=I,sill=I,skew=I,scale=I,nugget=1)
pcl2=GeoFit2(coordx=loc,corrmodel=corrmodel,data=z,
  likelihood="Conditional",type="Pairwise",model="SkewGaussian",
  optimizer="nlminb",lower=lower,upper=upper,
  neighb=5,start=start,fixed=fixed)
pcl2

```

The object `pcl2` include informations about the conditional pairwise likelihood estimation:

```

pcl2
#####
#####
Maximum Composite-Likelihood Fitting of Skew Gaussian Random Fields
Setting: Conditional Composite-Likelihood
Model: SkewGaussian
Type of the likelihood objects: Pairwise
Covariance model: Matern
Optimizer: nlminb
Number of spatial coordinates: 7352
Number of dependent temporal realisations: 1
Type of the random field: univariate
Number of estimated parameters: 5
Type of convergence: Successful
Maximum log-Composite-Likelihood value: -33034.35
Estimated parameters:
      mean      nugget      scale      sill      skew
    -0.6514     0.1065    174.4599     0.5250     0.8533
#####

```

The estimation of the skew parameter shows the presence of negative asymmetry in the anomalies data. Additionally, it can be appreciated that the skew Gaussian case shows a higher maximum log-conditional pairwise likelihood value, as expected, since the Gaussian random field is a special case of the skew Gaussian random field.

Checking model assumptions

Given the estimation of the Gaussian and skew-Gaussian random fields, the estimated residuals are

$$\widehat{Y(\mathbf{s}_i)} = \frac{y(\mathbf{s}_i) - \hat{\mu}}{(\hat{\sigma}^2)^{\frac{1}{2}}} \quad i = 1, \dots, N \quad (10)$$

and

$$\widehat{U_\eta(\mathbf{s}_i)} = \frac{u(\mathbf{s}_i) - \hat{\mu}}{(\hat{\sigma}^2)^{\frac{1}{2}}} \quad i = 1, \dots, N \quad (11)$$

$\widehat{Y(\mathbf{s}_i)}$, for $i = 1, \dots, N$ can be viewed as a realization of a Gaussian random field with marginal distribution $N(0, 1)$ and with correlation function $\rho(\mathbf{h})$. Similarly $\widehat{U_\eta(\mathbf{s}_i)}$ for $i = 1, \dots, N$ can be viewed as a realization of a random field stationary of (3) with marginal distribution $SN(0, \omega, \delta)$ with $\delta = \eta/\sigma$, $\omega^2 = (\eta^2 + \sigma^2)/\sigma^2$ and with correlation function $\rho_{U_\eta}(\mathbf{h})$. The residuals can be computed using the `GeoResiduals` function:

```
resd1=GeoResiduals(pcl1); # residuals of Gaussian random field
resd2=GeoResiduals(pcl2); # residuals of skew-Gaussian random field
```

The marginal distribution assumption on the residuals can be graphically checked for instance with a qq-plot (see, Figure (4)) using the function `GeoQQ`:

```
### checking model residuals assumptions: marginal distribution
GeoQQ(resd1); #qq-plot residuals of Gaussian random field
GeoQQ(resd2); #qq-plot residuals of skew-Gaussian random field
```

It can be appreciated that the skew Gaussian case shows a better agreement between the theoretical and estimated quantiles with respect to the Gaussian case. Additionally, the covariance model assumption can be checked comparing the empirical and the estimated semi-variogram of the residuals using the `GeoVariogram` and `GeoCovariogram` functions (see Figure (4)).

```
### checking model residuals assumptions: covariance model

### semi-variogram residuals of Gaussian Random fields
vario1 <- GeoVariogram(data=resd2$data, coordx=loc,
                      maxdist=maxdist/4);
GeoCovariogram(resd1, show.vario=TRUE, vario1=evario, pch=20);

### semi-variogram residuals of skew-Gaussian Random fields
```

```
evariog2 <- GeoVariogram(data=resd2$data, coordx=loc,
                        maxdist=maxdist/4);
GeoCovariogram(resd2, show.vario=TRUE, vario=evariog2, pch=20);
```

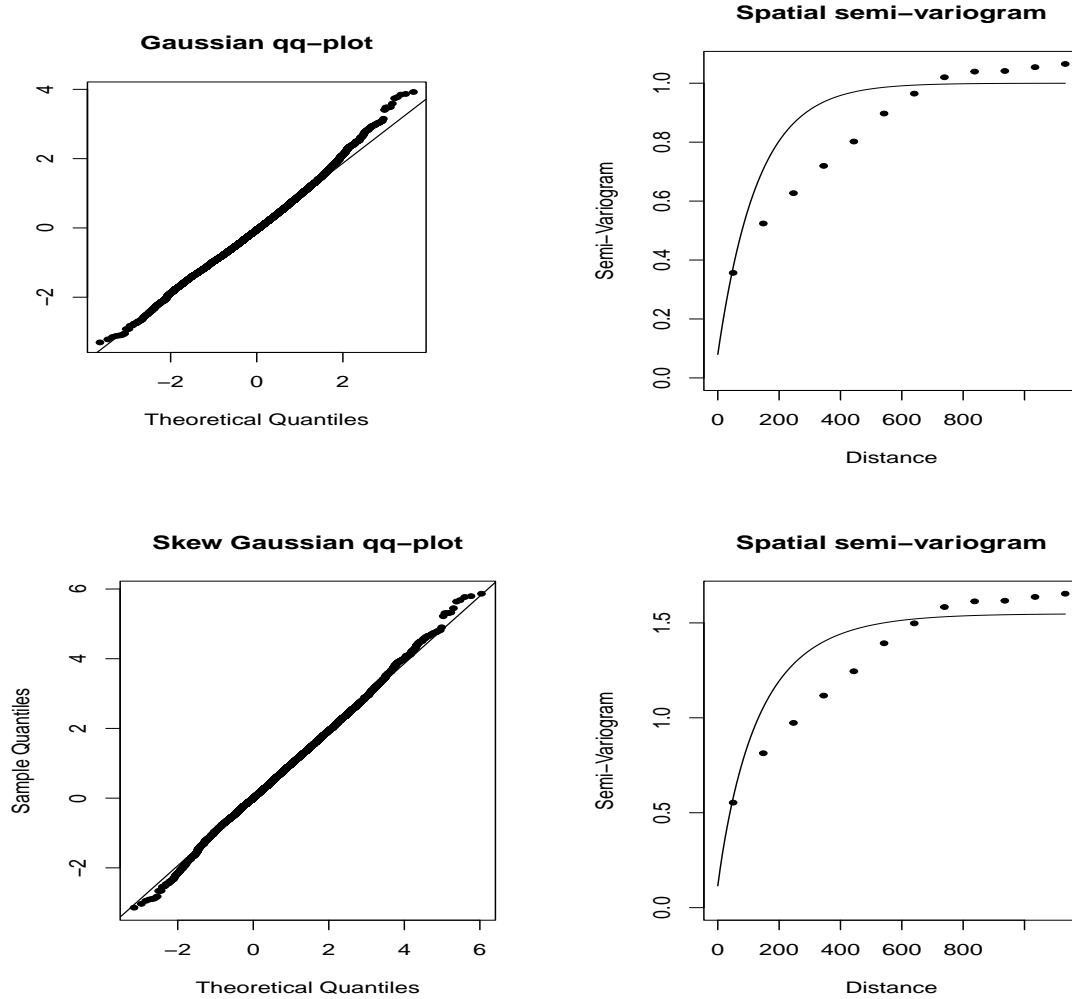


Figure 4: Upper part: qq-plot of the Gaussian residuals and empirical vs estimated semi-variogram of the residuals (from left to right). Bottom part: qq-plot of the skew Gaussian residuals and empirical vs estimated semi-variogram of the residuals (from left to right).

It can be seen that in the skew-Gaussian model the estimated semivariogram presents a better agreement with the empirical semivariogram.

Prediction

The package `GeoModels` allows to perform optimal linear prediction for the Gaussian and skew-Gaussian random fields. In the Gaussian case optimal linear prediction is equal to the optimal prediction (in the mean squared sense).

For a given spatial location \mathbf{s}_0 , the optimal linear prediction of a Gaussian or skew-Gaussian random fields is given by:

$$\hat{L}(\mathbf{s}_0) = \mu + \mathbf{c}^T R^{-1}[\mathbf{l} - \mu], \quad (12)$$

with $\hat{L}(\mathbf{s}_0) = Y(\mathbf{s}_0)$, $\mathbf{l} = \mathbf{Y}$ or $\hat{L}(\mathbf{s}_0) = U_\eta(\mathbf{s}_0)$, $\mathbf{l} = \mathbf{U}_\eta$ for the Gaussian and skew-Gaussian cases respectively. In addition:

- $\mathbf{c} = (\text{cor}(L(\mathbf{s}_0), L(\mathbf{s}_1)), \dots, \text{cor}(L(\mathbf{s}_0), L(\mathbf{s}_N)))^T$.
- $R = [\text{cor}(L(\mathbf{s}_i), L(\mathbf{s}_j))]_{i,j=1}^N$.

both R and \mathbf{c} are computed by using $\rho(\mathbf{h})$ and $\rho_{U_\eta}(\mathbf{h})$ for the Gaussian and skew-Gaussian case respectively. Moreover the associated mean square error is given by:

$$MSE(\hat{L}(\mathbf{s}_0)) = \text{Var}(L(\mathbf{s})) (1 - \mathbf{c}^T R^{-1} \mathbf{c}). \quad (13)$$

where $\text{Var}(L(\mathbf{s}))$ is given by $\text{Var}(Y(\mathbf{s})) = \sigma^2$ and $\text{Var}(U_\eta(\mathbf{s})) = \sigma^2 + \eta^2(1 - 2/\pi)$ for the Gaussian and skew-Gaussian random field respectively. Both (12) and (13) can be computed replacing the parameters with the conditional pairwise likelihood estimates.

We evaluate the predictive performances of the Gaussian and skew Gaussian random fields using cross validation, with the function `GeoCV`.

```
## Gaussian case
a1=GeoCV(pcl1,K=100,estimation=TRUE, n.fold=0.25,seed=9,
        local=TRUE,neighb=100)
[1] 'Cross-validation kriging can be time consuming ...'
[1] 'Starting iteration from 1 to 100 ...'
## skew-Gaussian case
a2=GeoCV(pcl2,K=100,estimation=TRUE, n.fold=0.25,seed=9,
        local=TRUE,neighb=100)
[1] 'Cross-validation kriging can be time consuming ...'
[1] 'Starting iteration from 1 to 100 ...'
```

The function basically randomly choose 75% of the data for estimation and prediction and use the remaining 25% as validation dataset. To speed-up the computation, the (optimal linear) predictions are internally obtained using the `GeoKrigloc` function that perform local kriging using a fixed set of neighbors specified in the `neighb` option. Then some prediction scores as RMSE (root mean squared error) and MAE (mean absolute value) are constructed by comparing the predictions with the (known) values in the validation dataset. This is iterated 100 times (it can be computationally intensive for large datasets, as in this example). We can compare the two models from prediction performance viewpoint, using the empirical mean of the 100 RMSEs and MAEs

```
mean(a1$rmse);mean(a2$rmse);
[1] 0.4809941
[1] 0.4805738
mean(a1$mae);mean(a2$mae);
[1] 0.3642906
[1] 0.3640285
```

It can be appreciated that the estimated skew Gaussian random field perform slightly better from prediction viewpoint even if, in the skew-Gaussian case, the optimal local linear prediction is used.

A kriging map with associated MSE can be obtained using the `GeoKrig` or the `GeoKrigloc` function. For the given location sites, we first need to specify the border of the region and then to construct a fine grid inside the border. The following code perform this task:

```
Sr1 = Polygon(loc)
Srs1 = Polygons(list(Sr1), "s1")
SpP = SpatialPolygons(list(Srs1))
long1=min(loc[,1])-10;long2=max(loc[,1])+10
lat1=min(loc[,2])-10;lat2=max(loc[,2])+10
lat_seq=seq(lat1,lat2,24)
lon_seq=seq(long1,long2,24)
coords_tot=as.matrix(expand.grid(lon_seq,lat_seq))
gr.in <- locations.inside(coords_tot, SpP)
```

Then optimal (local) linear prediction (12) and associated MSE (13) can be computed (using the estimated parameters) for the Gaussian and skew Gaussian cases, with the following code:

```
pr1<-GeoKrigloc(loc=gr.in,coordx=loc,corrmodel=corrmodel,mse=TRUE,
  model="Gaussian",neighb=100,
  param=as.list(c(pcl1$param,pcl1$fixed)),data=z)
pr2<-GeoKrigloc(loc=gr.in,coordx=loc,corrmodel=corrmodel,mse=TRUE,
  model="SkewGaussian",neighb=100,
  param=as.list(c(pcl2$param,pcl2$fixed)),data=z)
```

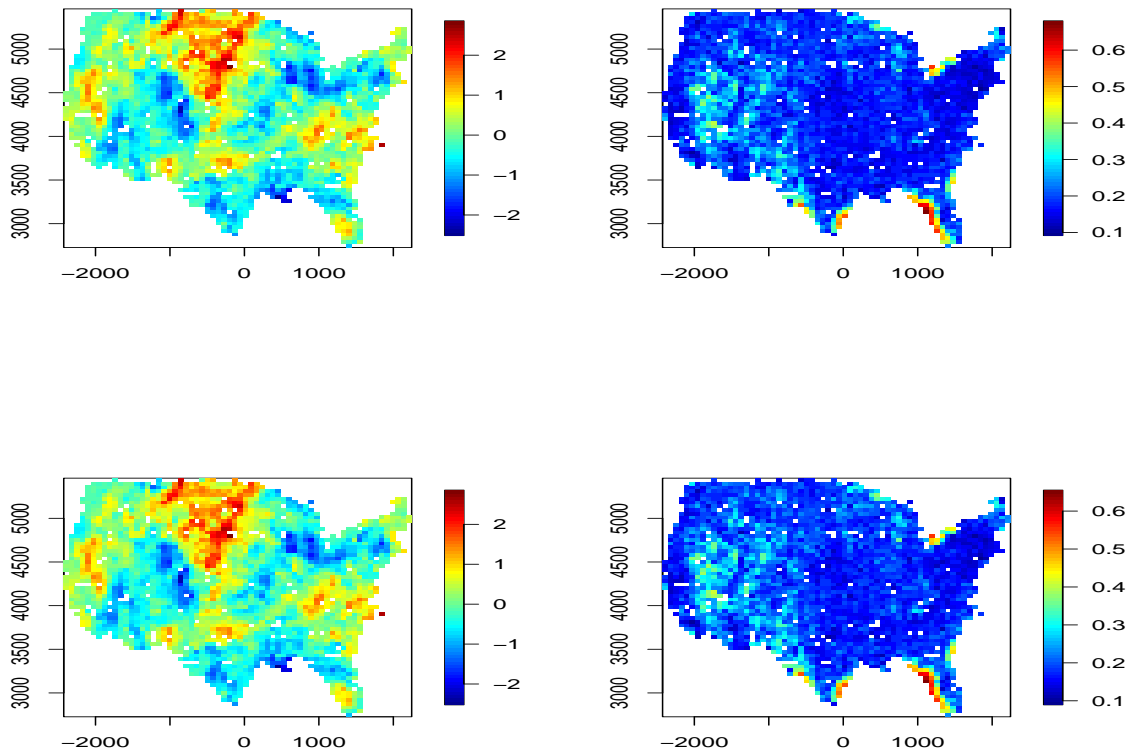


Figure 5: Kriging map and mean squared error map for the estimated. Gaussian (first row) and SkewGaussian (second row) random fields.

Finally a kriging map with associated mean square error (Figure 5) can be obtained with the following code:

```
quilt.plot(gr.in,pr1$pred)
quilt.plot(gr.in,pr1$mse)
quilt.plot(gr.in,pr2$pred)
quilt.plot(gr.in,pr2$mse)
```

References

- Alegria, A., S. Caro, M. Bevilacqua, E. Porcu, and J. Clarke (2017). Estimating covariance functions of multivariate skew-gaussian random fields on the sphere. *Spatial Statistics* 22, 388 – 402. Spatio-temporal Statistical Methods in Environmental and Biometrical Problems.
- Bevilacqua, M. and C. Gaetan (2015). Comparing composite likelihood methods based on pairs for spatial Gaussian random fields. *Statistics and Computing* 25, 877–892.
- 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(10), 1–25.
- Kaufman, C. G., M. J. Schervish, and D. W. Nychka (2008). Covariance tapering for likelihood-based estimation in large spatial data sets. *Journal of the American Statistical Association* 103, 1545–1555.
- Zhang, H. and A. El-Shaarawi (2010). On spatial skew-Gaussian processes and applications. *Environmetrics* 21(1), 33–47.