

## **The R-INLA tutorial: SPDE models**

Warning: work in progress...

Suggestions are welcome to [`elias@r-inla.org`](mailto:elias@r-inla.org)

Elias T. Krainski and Finn Lindgren

October 8, 2013

## **Abstract**

In this tutorial we presents how we fit models to spatial point refered data, the called geostatistical model, using INLA and SPDE. Before a fast introduction on such models, we start with the mesh building (necessary to fit the models) and show the application on a toy example with details. In forward examples we presents some functionalities to fit more complex models, such as non-Gaussian, models with two likelihoods, such as a semi-continuous, missaligned and preferential sampling, where we model jointly the point pattern process. Also, we introduce the non-stationary SPDE models, space-time model, data-cloning.

2013 October 02:

- \* mesh improvements:

- . use the new Finn's functions: `inla.mesh.2d` and `inla.noncovexhull`
- . non-convex, `SpatialPolygons`

- \* toy-example improved (maybe more clear...)

- \* new chapters:

- . simulation and likelihood through SPDE
- . point process
- . preferential sampling
- . spatio temporal example
- . data cloning

2013 October 08:

- \* Finn's suggestions to make clear the stack with two likelihoods examples and on non-

to do:

- \*\* English improvement \*\*

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	The Gaussian random field . . . . .	4
1.2	Simulation of a data set . . . . .	6
1.3	Maximum likelihood estimation . . . . .	7
<b>2</b>	<b>The SPDE approach</b>	<b>10</b>
2.1	The [Lindgren et al., 2011] results . . . . .	10
2.2	The triangulation . . . . .	11
2.2.1	Getting started . . . . .	11
2.2.2	Non-convex hull meshes . . . . .	14
2.2.3	Meshes for the toy example . . . . .	15
2.2.4	Meshes for Paraná state . . . . .	17
2.2.5	Triangulation with a SpatialPolygonsDataFrame . . . . .	19
2.3	The projector matrix . . . . .	19
2.4	GRF sampling through SPDE . . . . .	21
2.5	Maximum likelihood inference . . . . .	23
<b>3</b>	<b>A toy example</b>	<b>26</b>
3.1	The stack functionality . . . . .	27
3.2	Model fitting and some results . . . . .	28
3.3	Prediction of the random field . . . . .	29
3.3.1	Jointly with the estimation process . . . . .	29
3.3.2	After the estimation process . . . . .	30
3.3.3	Projection on a grid . . . . .	31
3.4	Prediction of the response . . . . .	31
3.4.1	By the posterior distribution . . . . .	31
3.4.2	By sum of linear predictor components . . . . .	32
3.4.3	Response on a grid . . . . .	33
3.5	Results from different meshes . . . . .	35
<b>4</b>	<b>Non-Gaussian response: Precipitation on Paraná</b>	<b>38</b>
4.1	The data set . . . . .	38
4.2	The model and covariate selection . . . . .	39
4.3	Testing the significance of spatial effect . . . . .	43
4.4	Prediction of the random field . . . . .	43
4.5	Prediction of the response on a grid . . . . .	44
<b>5</b>	<b>Semicontinuous model to daily rainfall</b>	<b>47</b>
5.1	The data and the model . . . . .	47
5.2	Fitting the model and some results . . . . .	49
5.3	Results for the spatial random effect . . . . .	52

<b>6</b>	<b>Joint modeling a covariate with misalignment</b>	<b>55</b>
6.1	The model . . . . .	55
6.1.1	Simulation from the model . . . . .	55
6.2	Fitting the model . . . . .	56
6.3	The results . . . . .	57
<b>7</b>	<b>Non stationary model</b>	<b>61</b>
7.1	Introduction . . . . .	61
7.2	An example . . . . .	61
7.3	Simulation on the mesh vertices . . . . .	63
7.3.1	Simulation with linear constraint . . . . .	64
7.4	Estimation with data simulated on the mesh vertices . . . . .	65
7.5	Estimation with locations not on mesh vertices . . . . .	66
<b>8</b>	<b>Point process: inference on the log-Cox process</b>	<b>68</b>
8.1	Data simulation . . . . .	68
8.2	Inference . . . . .	69
<b>9</b>	<b>Geostatistical inference under preferential sampling</b>	<b>72</b>
9.1	Fitting the usual model . . . . .	72
9.2	Estimation under preferential sampling . . . . .	73
<b>10</b>	<b>A space time example</b>	<b>76</b>
10.1	Data simulation . . . . .	76
10.2	Data stack preparation . . . . .	78
10.3	Fitting the model and some results . . . . .	78
10.4	A look at the posterior random field . . . . .	79
<b>11</b>	<b>Data cloning example</b>	<b>82</b>

# Chapter 1

## Introduction

If we have data measured at some locations, so we have a point refereed data set. By locations we mean some coordinates reference system where are each data from, for example, the longitude and latitude coordiantes. The point refereed data is very common in many areas of science. These type of data appear on mining, climate modeling, ecology, agriculture and other areas. If we want do model that data incorporating the reference about where are each data from, we want do use a model to point refereed data.

It is possible that we build a regression model considering each coordinates as covariate. But in some cases it is necessary a very complicated function based on coordinates to get a adequate description of the mean. For example, any complex non-linear function or a non-parametric function. If these type of model are only to incorporate some trend on the mean based on the coordinates. Also, this type of model is a fixed effect model.

But, we want a model to measure the first geography law in a simple way. This law says: “Everything is related to everything else, but near things are more related than distant things”, [Tobler, 1970]. So, we need a model to incorporate the property that a observation is more correlated with a observation collected on a neighbour local, than another that is collected on more distant local. One option to modelled this dependence is the use of a random effect spatially structured. This type of model is a model to spatial dependency, different of the spatial trend. But, it is possible to include both terms in a model. Also, the models for spatial dependency are used within more general models as random effects, the spatially structured random effects.

In spatial statistics we have models to incorporate the spatial dependency when the locations are areas (states, cities, etc.) and when the locations are points. In the last case, it is also possible that such locations are fixed or are random. The models to point refereed data with a spatially structured random effect is commonly called a geostatistical models. There are a specific area of spatial statistics that study these models, the geostatistics. See [Cressie, 1993] for a good introduction on the spatial statistics.

### 1.1 The Gaussian random field

To introduce some notation, let  $s$  any location on the study area and  $X(s)$  is the random effect at this location. We have  $X(s)$  a stochastic process, with  $s \in \mathbf{D}$ , were  $\mathbf{D}$  is the domain area of the locations and  $\mathbf{D} \in \mathbb{R}^d$ . Suppose, for example, that we have  $\mathbf{D}$  any country and we have any data measured on geographical locations,  $d = 2$ , within this country.

Suppose that we assume that we have a realization of  $x(s_i)$ ,  $i = 1, 2, \dots, n$ , a realization of  $X(s)$  in  $n$  locations. It is common assumed that  $x(s)$  has a multivariate Gaussian distribution. Also, if we assume that  $X(s)$  continuous over space, we have a continuously indexed Gaussian field (GF). It is because we suppose that it is possible that we get data in any location within the study region. To complete the specification of the distribution of  $x(s)$ , is necessary to defines it mean and covariance.

A very simple option, is the definition of a correlation function based only on euclidean distance between locations. This assume that if we have two pairs of points separated same distance  $h$ , both pairs have same correlation. Also, is intuitive to choose any function decreasing with  $h$ . There is some work about the GF and correlation functions in [Abrahamsen, 1997].

A very popular correlation function is the Matérn correlation function, that depends on a scale parameter  $\kappa > 0$  and a smoothness parameter  $\nu > 0$ . Considering two locations  $s_i$  and  $s_j$ , the stationary and isotropic Matérn correlation function is:

$$Cor_M(X(s_i), X(s_j)) = \frac{2^{1-\nu}}{\Gamma(\nu)} (\kappa \|s_i - s_j\|)^\nu K_\nu(\kappa \|s_i - s_j\|) \quad (1.1)$$

where  $\| \cdot \|$  denotes the euclidean distance and  $K_\nu$  is the modified Bessel function of second kind order. Also, we defines the Matérn covariance function by  $\sigma_x Cor(X(s_i), X(s_j))$ , where  $\sigma_x$  is the marginal variance of the process.

If we have a realization  $x(s)$  on  $n$  locations, we write the joint correlation, or joint covariance, matrix  $\Sigma$  making each entry  $\Sigma_{i,j} = \sigma_x Cor_M(X(s_i), X(s_j))$ . It is common to assume that  $X(x)$  has a zero mean. So, we have completely defined a multivariate distribution to  $x(s)$ .

Now, suppose now that we have a data  $y_i$  observed at locations  $s_i$ ,  $i = 1, \dots, n$ . If we suppose that we have an underlie GF that generate these data, we are going to fit the parameters of this process, making the identity  $y(s_i) = x(s_i)$ , and  $y(s_i)$  is just a realization of the GF. In this case, we the likelihood function is just the multivariate distribution with mean  $\mu_x$  and covariance  $\Sigma$ . If we assume  $\mu_x = \beta_0$ , we have four parameters to estimate.

In many situations we assume that we have an underlie GF but we no observe it and observe a data with a measurement error, i. e.,  $y(s_i) = x(s_i) + e_i$ . Additionally, it is common to assume that  $e_i$  independent of  $e_j$  for all  $i \neq j$  and  $e_i \sim N(0, \sigma_e)$ . These additional parameter,  $\sigma_e$ , measures the noise effect, called nugget effect. In this case the covariance of marginal distribution of  $y(s)$  is  $\sigma_e^2 I + \Sigma$ . This model is a short extension of the basic GF model, and in this case, we have one additional parameter to estimate. To look more about this model see [Diggle and Ribeiro Jr, 2007].

It is possible to describe this model within a larger class of models, the hierarchical models. Suppose that we have observations  $y_i$  on locations  $s_i$ ,  $i = 1, \dots, n$ . We start with

$$\begin{aligned} y_i | \theta, \beta, x_i, F_i &\sim P(y_i | \mu_i, \phi) \\ \mathbf{x} &\sim GF(0, \Sigma) \end{aligned} \quad (1.2)$$

where  $\mu_i = h(F_i^T \beta + x_i)$ ,  $F$  is a matrix of covariates,  $x$  is the random effects,  $\theta$  are parameters of random effects,  $\beta$  are covariate coefficients,  $h(\cdot)$  is a function mapping the linear predictor  $F_i^T \beta + x_i$  to  $E(y_i) = \mu_i$  and  $\phi$  is a dispersion parameter of the distribution, in the exponential family, assumed to  $y_i$ . To write the GF with nugget effect on this class, we replace  $F_i^T \beta$  by  $\beta_0$ , consider the Gaussian distribution to  $y_i$ , with variance  $\sigma_e^2$  and  $x$  as GF.

We have many extensions of this basic hierarchical model. Now, we stop these extensions and return on some extensions in later sections. But, if we know the properties of the GF, we are able to study all the practical models that contain, or are based on, this random effect.

It is mentioned that the data, or the random effect, on a finite number of  $n$  points that we have observed data is considered a realization of a multivariate Gaussian distribution. But, to evaluate the likelihood function, or the random effect distribution, we need to make computations of the multivariate Gaussian density. So, we have, in the log scale, the expression

$$-\frac{1}{2} (n \log(2\pi) + \log(|\Sigma|) + [x(s) - \mu_x]^T \Sigma^{-1} [x(s) - \mu_x]) \quad (1.3)$$

where  $\Sigma$  is a dense  $n \times n$ . To compute this, we need a factorization of this matrix. Because this matrix is a dense, this is a operation of order  $O(n^3)$ , so is one 'big n problem'.

An alternative used in some software's to do geostatistical analysis, is the use of the empirical variogram to fit the parameters of the correlation function. This option don't use any likelihood for the data and the multivariate Gaussian distribution to the random effects. A good description of these techniques is made on [Cressie, 1993].

However, it is adequate to assume any likelihood for the data and a GF for the spatial dependence, the model based approach on geostatistics, [Diggle and Ribeiro Jr, 2007]. So, in some times we need the use the multivariate Gaussian distribution to the random effects. But, if the dimension of the GF is big, it is impractical to make model based inference.

In another area of the spatial statistics, the analysis of areal data, there is models specified by conditional distributions that implies a joint distribution with a sparse precision matrix. These models are called the Gaussian Markov random fields (GMRF), [Rue and Held, 2005]. So, the inference when we use GMRF is more easy to do than we use the GF, because to work with two dimensional GMRF models, we have cost of  $O(n^{3/2})$  on the computations with its precision matrix. So, it is more easy to make analysis with big 'n'.

## 1.2 Simulation of a data set

First we took a look at the model from parametrization used in a commonly software used to analyse the point refereed data in **R**, the **geoR** package, [Ribeiro Jr and Diggle, 2001]. In this section we show how we simulate a dataset with this parametrization.

We remember here that one realization of a GF is just one realization of a multivariate Gaussian distribution with an appropriate covariance matrix. To specify these matrix, we need a set of locations and the matrix of distance between each point with all others. Based on this matrix  $n \times n$  of distances, we compute the covariance matrix and do one simulation of a multivariate Gaussian distribution.

Suppose that we have a set of  $n = 100$  locations, on a square with area one with bottom left and top right limits: (0,0) and (1,1). We choose these locations with density in left bottom corner higher than top right corner. The **R** code to do it is:

```
n <- 200; set.seed(123)
pts <- cbind(s1=sample(1:n/n-0.5/n)^2, s2=sample(1:n/n-0.5/n)^2)
```

and for get the (lower triangle) matrix of distances we do

```
dmat <- dist(pts)
```

and for the Matérn covariance we need the parameters:  $\sigma_x^2$ ,  $\kappa$  and  $\nu$ . Additionally, we need the mean  $\beta_0$  and the nugget parameter  $\sigma_e^2$ . We declare an values to such parameters by

```
beta0 <- 10; sigma2e <- 0.3; sigma2x <- 5; kappa <- 7; nu <- 1
```

We, first, consider  $y$  as mean  $\beta_0$  and that covariance  $\sigma_e^2 I + \Sigma$ . And, we get the covariance by

```
mcor <- as.matrix(2^(1-nu)*(kappa*dmat)^nu*
                  besselK(dmat*kappa,nu)/gamma(nu))
diag(mcor) <- 1; mcov <- sigma2e*diag(n) + sigma2x*mcor
```

Now we going to simulate one realization of a geostatistical model. First, we do the simulation of one realization of the multivariate Gaussian distribution. Here we remember that if we want to get  $y \sim N(\mu, \Sigma)$  from  $z \sim N(0, I)$ , we do  $y = \mu + zL$ , where  $L$  is a matrix that  $L^T L = \Sigma$ . So, we use the Cholesky factorization of the covariance matrix, because if  $L$  is the Cholesky of the  $\Sigma$ , so  $L^T L = \Sigma$ . In **R** we use:



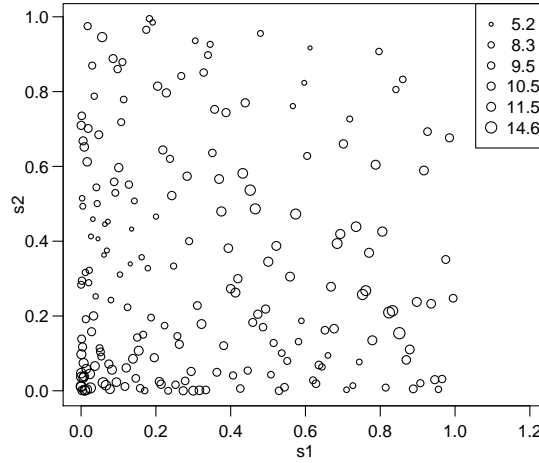


Figure 1.1: Visualization of the simulated data

```
L <- chol(mcov); set.seed(234); y1 <- beta0 + drop(rnorm(n)%*%L)
```

We show this simulated data in a graph of the locations with size of points proportional of the simulated values on Figure 1.1 with code below

```
par(mar=c(3,3,1,1), mgp=c(1.7, 0.7, 0), las=1)
plot(pts, asp=1, xlim=c(0,1.2), cex=y1/10)
q <- quantile(y1, 0:5/5)
legend('topright', format(q, dig=2), pch=1, pt.cex=q/10)
```

### 1.3 Maximum likelihood estimation

In this section we get the maximum likelihood estimation of the parameters of the model used to simulate the data on the previous section. We know that is adequate the use of the partial derivatives with respect to each parameter and the Fisher information matrix to use the Fisher scoring algorithm. Under Gaussian likelihood we have that the mean and variance are independently. So, taking the derivative with respect to  $\beta$  we get

$$(\mathbf{F}'\Sigma^{-1}\mathbf{F})\hat{\beta} = \mathbf{F}'\mathbf{F}$$

So, we want to write a likelihood function as a function of the unknown parameters and the data. For the others parameters we naively use one of the quasi-Newton methods, without using derivatives.

There are another details to a good performance on the maximum likelihood estimation. One is to write the variance parameters,  $\sigma_e^2$  and  $\sigma_x^2$ , in terms of the relative one. Such details are on the *Model Based Geostatistics* book, [Diggle and Ribeiro Jr, 2007], and on the **geoR** functions documentations.

Here we just write the negative of the multivariate normal density, considering the mean as a regression with identity link and the covariance with Matérn covariance plus the nugget effect. In the function implemented bellow, we have the also that the matrix of distances between the points because are one of the arguments, because it is not necessary to recalculate at each interaction of the minimization algorithm. We use the negative of the logarithm of the likelihood function because, by default, the algorithm find the minimum of a function. In the computation of the likelihood, we use the Cholesky factorization to compute the determinant and also the inverse of the covariance matrix. The implemented function is a function of the four length parameters vector ( $\sigma_e^2$ ,  $\sigma_x^2$ ,  $\kappa$  and  $\nu$ ), the covariate (design) matrix, the data and the matrix of distances:

```

nllf <- function(pars, ff, data, m) {
  m <- 2^(1-pars[4])*(pars[3]*m)^pars[4]*
    besselK(m*pars[3],pars[4])/gamma(pars[4])
  diag(m) <- 1; m <- pars[1]*diag(length(data)) + pars[2]*m
  m <- chol(m); ldet <- sum(log(diag(m))); m <- chol2inv(m)
  beta <- solve(crossprod(ff, m)%%ff, crossprod(ff, m)%%data)
  z <- data-ff)%%beta; ss <- drop(crossprod(z,m)%%z)
  return(ldet + (nrow(m)*log(2*pi) + ss)/2)
}

```

the first argument of this function is the vector of the five parameters of the model, such as  $\beta_0$ ,  $\sigma_e^2$ ,  $\sigma_s^2$ ,  $\kappa$  and  $\nu$ .

For test, we calculates the likelihood at true values of the parameters and at another set of the values

```

c(nllf(c(sigma2e, sigma2x, kappa, nu), matrix(1,n), y1, as.matrix(dmat)),
  nllf(c(0, sigma2x, kappa, nu), matrix(1,n), y1, as.matrix(dmat)))

```

```
[1] 282.1201 811.4252
```

We get the maximum likelihood estimates with 'L-BFGS-B' method implemented on `optim()` function. Also, using `hessian=TRUE`, we got the approximate Hessian matrix to use in asymptotic inference. We got the maximum likelihood estimates

```

(ores <- optim(c(sigma2e, sigma2x, kappa, nu), nllf, hessian=TRUE,
  ff=matrix(1,n), data=y1, m=as.matrix(dmat),
  method='L-BFGS-B', lower=rep(1e-5,4)))$par

```

```
[1] 0.2824531 3.2739337 9.4016876 1.0864144
```

and the squared root of the hessian diagonal

```
(se <- sqrt(diag(solve(ores$hessian))))
```

```
[1] 0.08000842 0.93752829 4.08361766 0.42725613
```

This naive solution by likelihood is to show how it works. But, in the **geoR** package, we have functions to get simulations and to do estimation. The difference is that the **geoR** package uses  $\phi = 1/\kappa$  for the scale parameter and it calls `kappa` for the smoothness parameter. The `grf()` function can be used to get samples of the geostatistical model from many correlation functions. To get exactly the same data, we use

```

require(geoR); set.seed(234)
grf1 <- grf(grid=pts, cov.pars=c(sigma2x, 1/kappa), mean=beta0,
  nugget=sigma2e, kappa=nu, messages=FALSE)

```

Also, we have the `likfit()` function to perform the maximum likelihood estimation. With this function we also got the maximum likelihood estimates for the smoothness parameter by

```

(glres <- likfit(grf1, ini=c(sigma2x, 1/kappa), messages=FALSE,
  nugget=sigma2e, kappa=nu, fix.kappa=FALSE))

```

likfit: estimated model parameters:

```

      beta      tausq      sigmasq      phi      kappa
"9.5457" "0.2824" "3.2741" "0.1064" "1.0862"

```

Practical Range with cor=0.05 for asymptotic range: 0.4403836

likfit: maximised log-likelihood = -281.1

In the **geoR** parametrization, the  $\sigma_e^2$  is called **tausq**,  $1/\kappa$  is called **phi** and the smoothness parameter  $\nu$  is called **kappa**.

But, for comparison with INLA results, we fix the smoothness parameter  $\nu$  on the true value used to simulate the data, and find the maximum likelihood estimates to another parameters by

```
(fit.l <- likfit(grf1, ini.cov.pars=c(sigma2x, 1/kappa),
               nugget=sigma2e, kappa=1, messages=FALSE))

likfit: estimated model parameters:
      beta      tausq  sigmasq      phi
"9.5349" "0.2709" "3.3234" "0.1156"
Practical Range with cor=0.05 for asymptotic range: 0.4620667

likfit: maximised log-likelihood = -281.1
```

Notice that we also estimate the mean parameter  $\beta$ , the parameter of fixed effect part of the model. Also, we note that the likelihood has similar value, just because the estimated value of  $\nu$  is very close to the fixed value.

## Chapter 2

# The SPDE approach

In the literature there is some ideas to fit GF by an approximation of the GF to any GMRF. The very good alternative found is such one that found a explicit link between the an stochastic partial differential equation (SPDE) and the Matérn Gaussian random fields, [Lindgren et al., 2011]. The solution of the that SPDE thought the Finite Element Method (FEM) provides a explicit link between GF to GMRF.

### 2.1 The [Lindgren et al., 2011] results

The SPDE approach is based on two main results. The first one extends the result obtained by [Besag, 1981]. This result is to approximate a GF with generalized covariance function, obtained when  $\nu \rightarrow 0$  in the Matérn correlation function. This approximation, considering a regular two-dimensional lattice with number of sites tending to infinite, is that the full conditional have

$$E(x_{ij}|x_{-ij}) = \frac{1}{a}(x_{i-1,j} + x_{i+1,j} + x_{i,j-1} + x_{i,j+1}) \quad (2.1)$$

and  $Var(x_{ij}|x_{-ij}) = 1/a$  for  $|a| > 4$ . In the representation using precision matrix, we have, for one single site, just the upper right quadrant and with  $a$  as the central element, that

$$\begin{matrix} -1 & & \\ & a & -1 \end{matrix} \quad (2.2)$$

Considering a GF  $x(\mathbf{u})$  with the Matérn covariance is a solution to the linear fractional SPDE

$$(\kappa^2 - \Delta)^{\alpha/2} x(\mathbf{u}) = \mathbf{W}(\mathbf{u}), \quad \mathbf{u} \in \mathbb{R}^d, \quad \alpha = \nu + d/2, \quad \kappa > 0, \quad \nu > 0, \quad (2.3)$$

[Lindgren et al., 2011] show that for  $\nu = 1$  and  $\nu = 2$  the GMRF representations are convolutions of (2.2). So, for  $\nu = 1$ , in that representation we have:

$$\begin{matrix} & 1 & & \\ -2a & & 2 & \\ 4 + a^2 & -2a & & 1 \end{matrix} \quad (2.4)$$

and, for  $\nu = 2$ :

$$\begin{matrix} & -1 & & & \\ & 3a & & -3 & \\ -3(a^2 + 3) & & 6a & & 3 \\ a(a^2 + 12) & -3(a^2 + 3) & 3a & & -1 \end{matrix} \quad (2.5)$$

This is an intuitive result, because if we have larger  $\nu$  on the Matérn correlation function of the GF, we need more non zero neighbours sites in the GMRF representation. Remember that it is the smoothness parameter, so if the process is more smooth, we need a more larger neighbourhood on the GMRF representation.

If the spatial locations are on a irregular grid, it is necessary the use of the second result on [Lindgren et al., 2011]. To extend the first result, a suggestion is the use of the finite element method (FEM) for a interpolation of the locations of observations to the nearest grid point. To do it, suppose that the  $\mathbb{R}^2$  is subdivided into a set of non-intersecting triangles, where any two triangles meet in at most a common edge or corner. The three corners of a triangle are named *vertices*. The suggestion is to start with the location of the observed points and add some triangles (heuristically) with restriction to maximize the allowed edge length and minimize the allowed angles. The the approximation is

$$x(\mathbf{u}) = \sum_{k=1}^n \psi_k(\mathbf{u}) w_k$$

for some chosen basis functions  $\psi_k$ , Gaussian distributed weights  $w_k$  and  $n$  the number of vertices on the triangulation. If the functions  $\psi_k$  are piecewise linear in each triangle,  $\psi_k$  is 1 at vertices  $k$  and 0 at all other vertices.

The second result is obtained using the  $n \times n$  matrices  $\mathbf{C}$ ,  $\mathbf{G}$  and  $\mathbf{K}$  with entries

$$C_{i,j} = \langle \psi_i, \psi_j \rangle, \quad G_{i,j} = \langle \nabla \psi_i, \nabla \psi_j \rangle, \quad (\mathbf{K}_{\kappa^2})_{i,j} = \kappa^2 C_{i,j} + G_{i,j} \quad (2.6)$$

to get the precision matrix  $\mathbf{Q}_{\alpha,\kappa}$  as a function of  $\kappa^2$  and  $\alpha$ :

$$\begin{aligned} \mathbf{Q}_{1,\kappa^2} &= \mathbf{K}_{\kappa^2}, \\ \mathbf{Q}_{2,\kappa^2} &= \mathbf{K}_{\kappa^2} \mathbf{C}^{-1} \mathbf{K}_{\kappa^2}, \\ \mathbf{Q}_{\alpha,\kappa^2} &= \mathbf{K}_{\kappa^2} \mathbf{C}^{-1} \mathbf{Q}_{\alpha-2,\kappa^2} \mathbf{C}^{-1} \mathbf{K}_{\kappa^2}, \quad \text{for } \alpha = 3, 4, \dots \end{aligned} \quad (2.7)$$

Here we have too the notion that if  $\nu$  increases, we need a more dense precision matrix.

The  $\mathbf{Q}$  precision matrix is generalized for a fractional values of  $\alpha$  (or  $\nu$ ) using a Taylor approximation, see the author's discussion response in [Lindgren et al., 2011]. From this approximation, we have the polynomial of order  $p = \lceil \alpha \rceil$  for the precision matrix

$$\mathbf{Q} = \sum_{i=0}^p b_i \mathbf{C} (\mathbf{C}^{-1} \mathbf{G})^i. \quad (2.8)$$

For  $\alpha = 1$  and  $\alpha = 2$  we have the (2.7). Because, for  $\alpha = 1$ , we have  $b_0 = \kappa^2$  and  $b_1 = 1$ , and for  $\alpha = 2$ , we have  $b_0 = \kappa^4$ ,  $b_1 = \alpha \kappa^4$  and  $b_2 = 1$ . For fractional  $\alpha = 1/2$   $b_0 = 3\kappa/4$  and  $b_1 = \kappa^{-1}3/8$ . And, for  $\alpha = 3/2$ , ( $\nu = 0.5$ , the exponential case),  $b_0 = 15\kappa^3/16$ ,  $b_1 = 15\kappa/8$ ,  $b_2 = 15\kappa^{-1}/128$ . Using these results combined with recursive construction, for  $\alpha > 2$ , we have GMRF approximations for all positive integers and half-integers.

## 2.2 The triangulation

The first step to fit the model is the construction of the 'mesh'. This step must be made VERY CAREFULLY. It is similar to choosing the integration points on a numeric ingration algorithm. We need regular ones? How many points is needed?

Additionally, here we need, yet, to add, ALSO CAREFULLY, additional points around the boundary, the outer extension. It is necessary to avoid the boundary effect were we have variance twice at border than that one within domain [Lindgren, 2012]. For more about it please see [Lindgren and Rue, 2013].

### 2.2.1 Getting started

For two dimentional mesh, we have a main function `inla.mesh.2d()` that is recommended to use for building a mesh. This function creates the Constrained Refined Delaunay Triangulation (CRDT) that we just call mesh. There are a several options on is function:

```
args(inla.mesh.2d)

function (loc = NULL, loc.domain = NULL, offset = NULL, n = NULL,
         boundary = NULL, interior = NULL, max.edge, min.angle = NULL,
         cutoff = 0, plot.delay = NULL)
NULL
```

We need some reference about the study region. It is possible to be provided by the location points, supplied on the `loc` argument, or by a location domain (a polygon for example). If we supply the points location, or the domain is supplied `loc.domain` argument, the algorithm find a convex hull mesh. A non convex hull mesh can be made when we provide a (list of) set of polygons on the `boundary` argument, each elements of this list of `inla.mesh.segment()` class.

The another mandatory argument is the `max.edge`. This argument specifies the maximum allowed triangle edge lengths in the inner domain and in the outer extension. So, it is a scalar or length two vector. This argument is a numeric on the same scale unit of the coordinates.

The another arguments are used on as additional conditions. The `offset` is a numeric, or length two vector. If negative, interpreted as a factor relative to the approximate data diameter. If positive, is the extension distance on same scale unit to the coordinates provided.

The argument `n` is the initial number of points on the extended boundary. The `interior` is a list of segments to specify interior constraints, each one of `inla.mesh.segment` class. A good mesh need to have triangles more regular as possible, on size and shape. To help this requeriment, upward the `max.edge`, we have the `min.angle` argument, that can be scalar or length two vector, to specify the minimum internal angles of the triangles on the inner domain and on the outer extension. This values at most 21 guarantee the convergence of the algorithm.

Still on the shape of the triangles, we also have the `cutoff` argument that is the minimum allowed distance between points. It means that points at with less distance than the supplied value are replaced by a single vertex. So, it avoids small triangles and must be provided a positive number, specially when we have some very close points, on points location or on the domain.

To understand how this function works, we apply this function, varying some arguments, to first five locations of the toy dataset.

```
data(SPDEtoy)
coords <- as.matrix(SPDEtoy[,1:2]) ; p5 <- coords[1:5,]
```

We also build some meshes using the domain and not the points, we define the domain with

```
pl.dom <- cbind(c(0,1,1,0.7,0), c(0,0,0.7,1,1))
```

Creating some meshes for the first five points:

```
m1 <- inla.mesh.2d(p5, max.edge=c(0.5, 0.5))
m2 <- inla.mesh.2d(p5, max.edge=c(0.5, 0.5), cutoff=0.1)
m3 <- inla.mesh.2d(p5, max.edge=c(0.1, 0.5), cutoff=0.1)
m4 <- inla.mesh.2d(p5, max.edge=c(0.1, 0.5), offset=c(0,-0.65))
m5 <- inla.mesh.2d(, pl.dom, max.edge=c(0.3, 0.5), offset=c(0.03, 0.5))
m6 <- inla.mesh.2d(, pl.dom, max.edge=c(0.3, 0.5), offset=c(0.03, 0.5), cutoff=0.1)
m7 <- inla.mesh.2d(, pl.dom, max.edge=c(0.3, 0.5), n=5, offset=c(.05,.1))
m8 <- inla.mesh.2d(, pl.dom, max.edge=c(.3, 0.5), n=7, offset=c(.01,.3))
m9 <- inla.mesh.2d(, pl.dom, max.edge=c(.3, 0.5), n=4, offset=c(.05,.3))
```

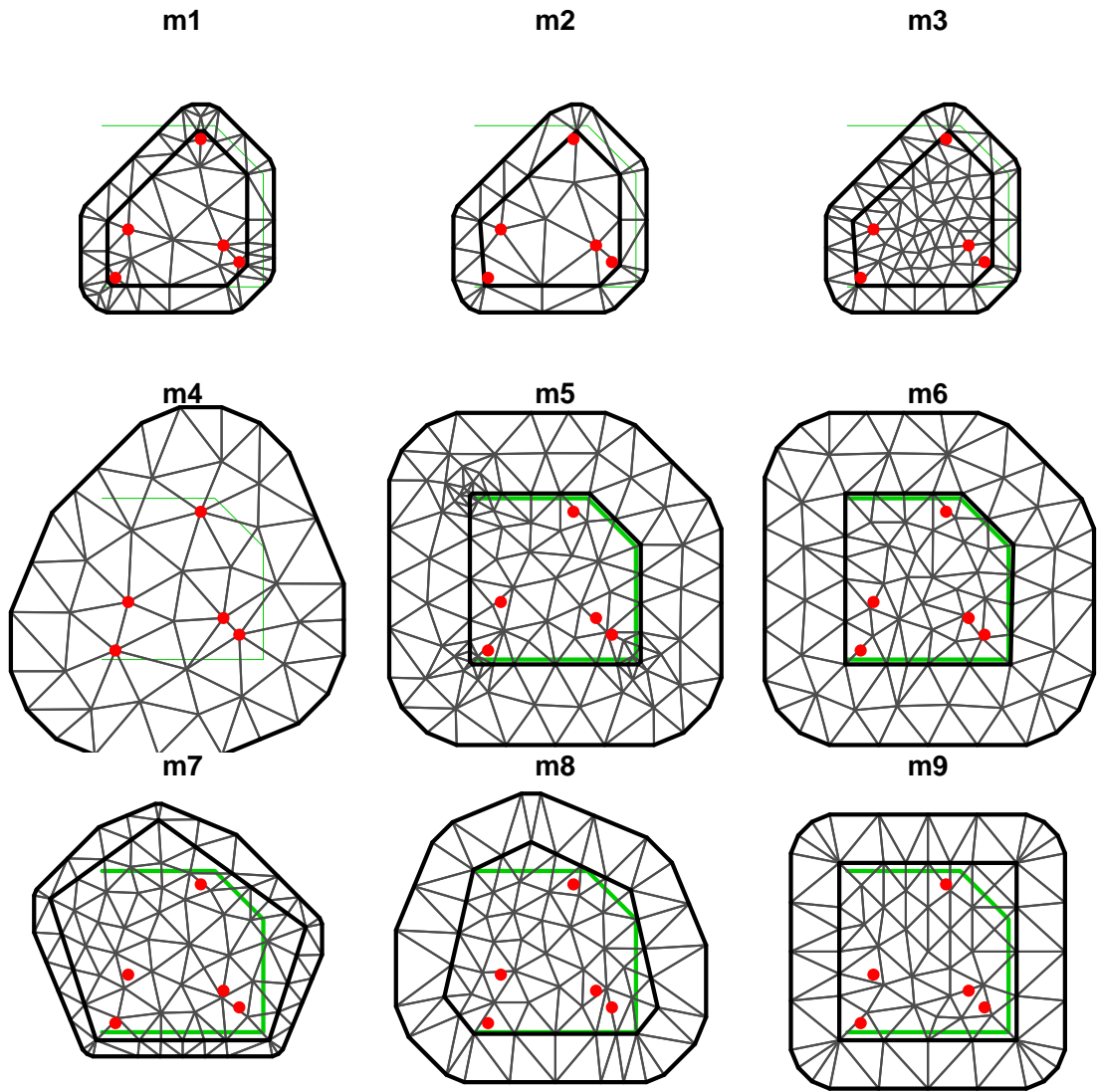


Figure 2.1: Triangulation with different restrictions.

We visualize these meshes on the Figure 2.1, with the code below

```
par(mfrow=c(3, 3), mar=c(0,0,1,0))
for (i in 1:9) {
  plot(pl.dom, type='l', col=3, lwd=2*(i>4), xlim=c(-0.57,1.57),
       main = paste('m',i,sep=''), asp=1, axes=FALSE)
  plot(get(paste('m', i, sep='')), add=TRUE)
  points(p5, pch=19, col=2)
}
```

In the **m1** mesh has two main problems: 1) some triangles with small inner angles, 2) some large triangles on inner domain. In the **m2** mesh, we relax the restriction on the locations, because points with distance less than the cutoff are considered single vertex. This avoids some of the triangles (at bottom right side) with small angles on the previous one. So the **cutoff** is a **VERY GOOD** idea! The **m3** mesh on the top right has its inner triangles has less than 0.1 and this mesh looks better than the two previous.

The **m4** was made with no first convex hull extension around the points. It has just the second one. In this case, the length of inner triangles does not works (first value on **max.edge** argument) and we have just triangles with length edges at most 0.5. The shapes of the triangles looks good, except for these ones with vertices including the two points at

bottom right side.

The `m5` mesh is made just using the domain polygon and it has shape similar to the domain area. In this mesh we have some small triangles at corners just due the fact that is build without `cutoff`. Also, we have a (relatively) small first extension and a (relatively) large second one. On the `m6` mesh we have add the a cutoff and get a better mesh than the previous one.

In the last tree meshes we change the initial number of extension points. Its can be usefull to change when we have some polygons domain to get convergence. Here we show the shape of the final mesh that we got with, for example, `n=5` on the `m7` mesh, This number produces a mesh that seems not adequate for this domain because we have non uniform extension behind the border. The `m9` mesh has very bad triangles shapes.

The object returned by `inla.mesh.2d()` function has class `inla.mesh` and contains a list of things:

```
class(m1)

[1] "inla.mesh"

names(m1)

[1] "meta"      "manifold" "n"          "loc"        "graph"      "segm"
[7] "idx"
```

The number of vertices on each mesh is

```
c(m1$n, m2$n, m3$n, m4$n, m5$n, m6$n, m7$n, m8$n, m9$n)

[1] 61 36 79 47 117 88 87 69 72
```

The 'graph' element represents the CRDT obtained. More, on 'graph' element we have the a matrix that represents the graph of neighborhood structure. For example, for `m1` we have 'A'

```
dim(m1$graph$vv)

[1] 61 61
```

The vertices that correspond the location points are identified on 'idx' element

```
m1$idx$loc

[1] 24 25 26 27 28
```

## 2.2.2 Non-convex hull meshes

All the meshes on Figure 2.1 are made to have a convex hull boundary. A convex hull is a polygon triangles out of the domain area, the extension made to avoid the boundary effect. A triangulation without additional border can be made by supplying the `boundary` argument instead the `location` or `loc.domain` argument. One way is to build a boundary for the points and supply it on `boundary` argument. But, in this case we need a boundary

We can build boundaries using the `inla.nonconvex.hull()` function

```
args(inla.nonconvex.hull)

function (points, convex = -0.15, concave = convex, resolution = 40,
  eps = NULL)
NULL
```



In this function we provide the points set and some constraint arguments. We are able to control the shape of the boundary by its convexity, concavity and resolution. We make some boundaries and build a mesh with each one to understand better it.

```
bound1 <- inla.nonconvex.hull(p5)
bound2 <- inla.nonconvex.hull(p5, convex=0.5, concave=-0.15)
bound3 <- inla.nonconvex.hull(p5, concave=0.5)
bound4 <- inla.nonconvex.hull(p5, concave=0.5, resolution=c(20, 20))
m10 <- inla.mesh.2d(boundary=bound1, cutoff=0.05, max.edge=c(.1,.2))
m11 <- inla.mesh.2d(boundary=bound2, cutoff=0.05, max.edge=c(.1,.2))
m12 <- inla.mesh.2d(boundary=bound3, cutoff=0.05, max.edge=c(.1,.2))
m13 <- inla.mesh.2d(boundary=bound4, cutoff=0.05, max.edge=c(.1,.2))
```

These meshes are visualized on Figure 2.5 by commands bellow

```
par(mfrow=c(2,2), mar=c(0,0,1,0))
for (i in 10:13) {
  plot(get(paste('m', i, sep='')), asp=1, main='')
  points(p5, pch=19, col=2); title(main=paste('m', i, sep=''))
}
```

The m10 mesh is builded with a boundary that we got with default arguments on `inla.nonconvex.hull()` function. The default `convex` and `concave` arguments are both equal 0.15 proportion of the points domain radius, that is computed by

```
max(diff(range(p5[,1])), diff(range(p5[,2])))*.15
```

```
[1] 0.12906
```

If we supply larger convex value, like the other used to generate the m11, we got a larger boundary. It's because all circles with centre on each point and radius less than the convex value are incide to the boundary. When we choose a larger concave value, the boundary used to the 'm12' and 'm13' meshes, we don't have circles with radius less than the concave value outside the boudary. If we choose a smaller resolution, we got a boundary with small resolution (in terms of number of points), for example, comparing the m12 and m13 meshes.

### 2.2.3 Meshes for the toy example

To analyze the toy data set, we use six triangulations options to make comparisons on section 3.5. The first mesh forces the location points as vertices of the mesh.

```
mesh1 <- inla.mesh.2d(coords, max.edge=c(0.035, 0.1))
mesh2 <- inla.mesh.2d(coords, max.edge=c(0.1, 0.2))
```

The second and third meshes are based on the points, but we use the cutoff greather than zero to avoid small triangles on regions where we have more density observations

```
mesh3 <- inla.mesh.2d(coords, max.edge=c(0.1, 0.2), cutoff=0.03)
```

We also build three other meshes based on the domain area. These are builded to have aproximately the same number of vertices than previous

```
mesh4 <- inla.mesh.2d(, pl.dom, max.e=c(0.04, 0.1), cutoff=0.015)
mesh5 <- inla.mesh.2d(, pl.dom, max.e=c(0.08, 0.2), cutoff=0.025)
mesh6 <- inla.mesh.2d(, pl.dom, max.e=c(0.10, 0.2), cutoff=0.045)
```

The number of nodes on each one of these meshes are

```
c(mesh1$n, mesh2$n, mesh3$n, mesh4$n, mesh5$n, mesh6$n)
```

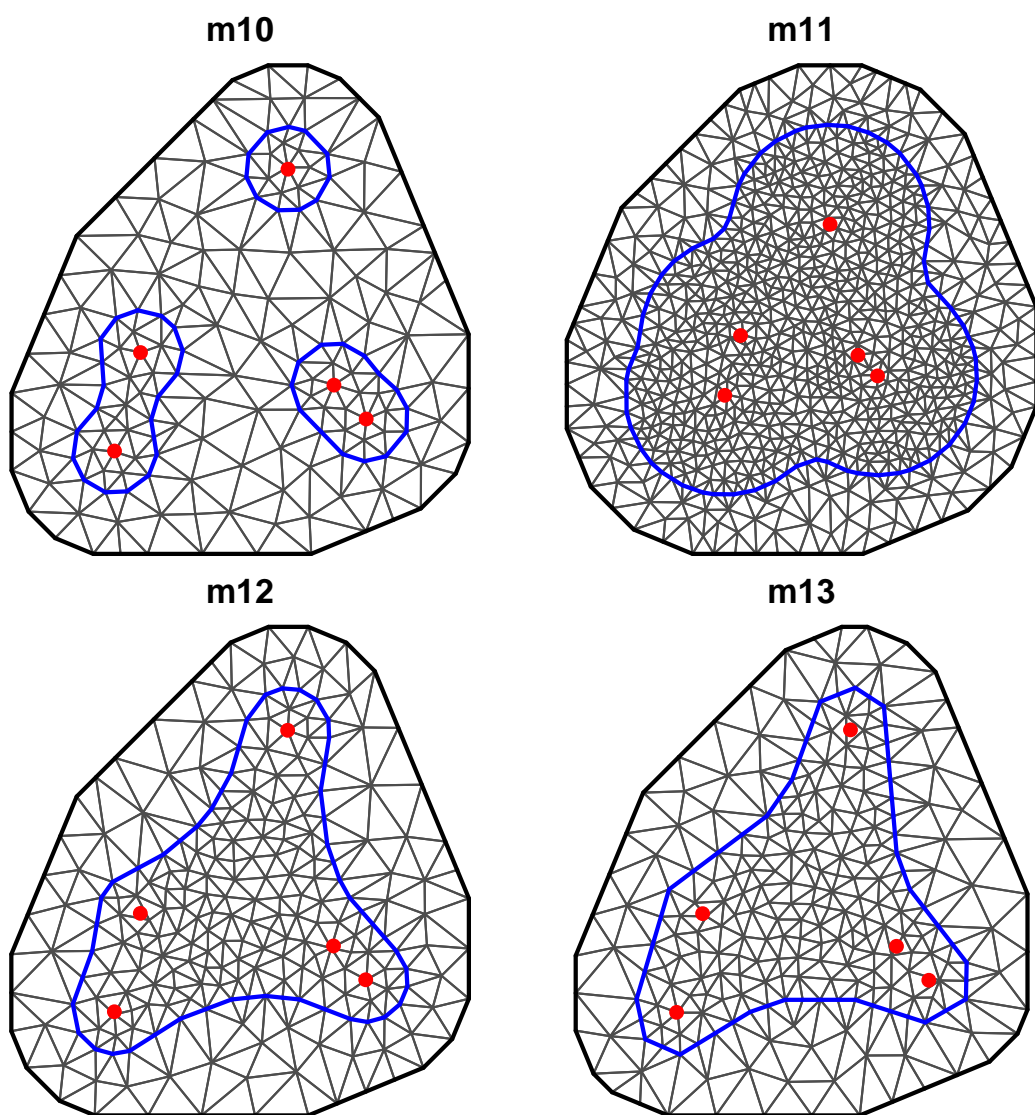


Figure 2.2: Non-convex meshes with different boundaries.

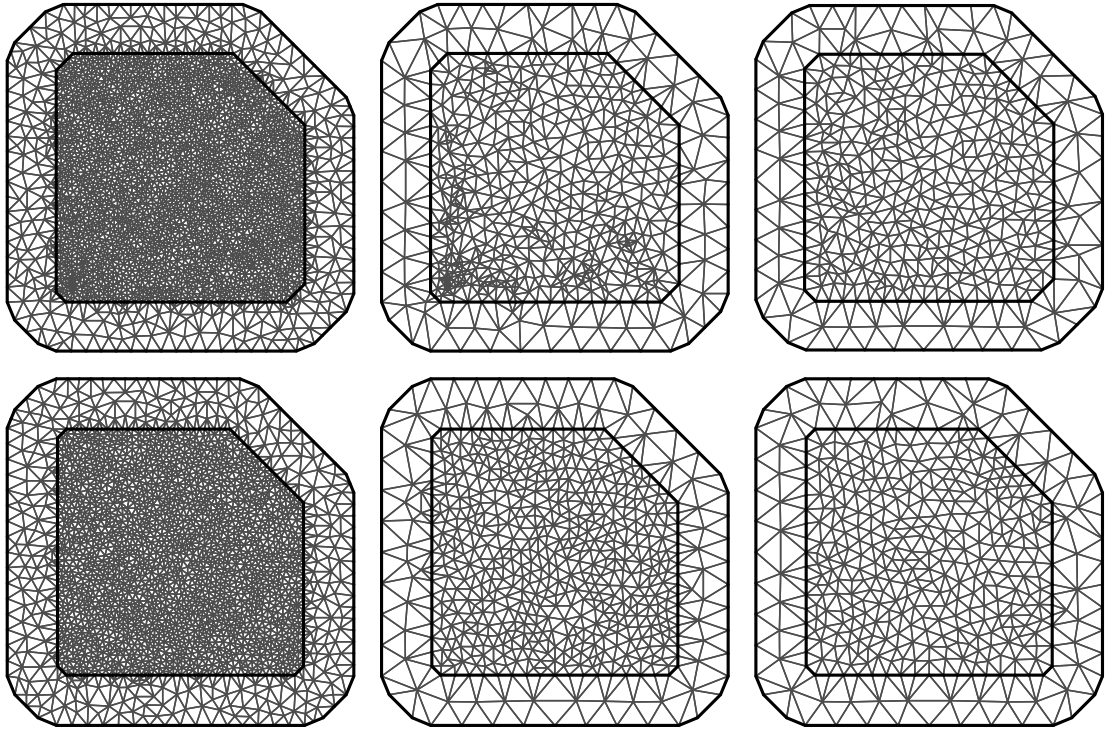


Figure 2.3: Six triangulation options for the toy example.

```
[1] 2900 593 435 2361 617 422
```

These six meshes are showed on Figure 2.3 with code below

```
par(mfrow=c(2,3), mar=c(0,0,0,0))
for (i in 1:6)
  plot(get(paste('mesh',i,sep='')), asp=1, main='')

```

## 2.2.4 Meshes for Paraná state

We have some examples using data collected on Paraná state, in Brasil. In this case we need to take into account two things: one is the shape of these domain area and other is the sistem of coordinates reference.

We have the daily rainfall data

```
data(PRprec)
dim(PRprec)
```

```
[1] 616 368
```

```
PRprec[1:3, 1:10]
```

	Longitude	Latitude	Altitude	d0101	d0102	d0103	d0104	d0105	d0106	d0107
1	-50.8744	-22.8511	365	0	0	0	0.0	0	0	2.5
3	-50.7711	-22.9597	344	0	1	0	0.0	0	0	6.0
4	-50.6497	-22.9500	904	0	0	0	3.3	0	0	5.1

that consists of the dayli rainfall data from 616 stations for each day of the 2012 year. The coordinates (two first columns) are on the latlong projection.

Also, we have the Paraná state polygon with

```

par(mfrow=c(1,2), mar=c(0,0,0,0))
plot(prmesh1, asp=1, main=""); lines(PRborder, col=3)
plot(prmesh2, asp=1, main=""); lines(PRborder, col=3)

```

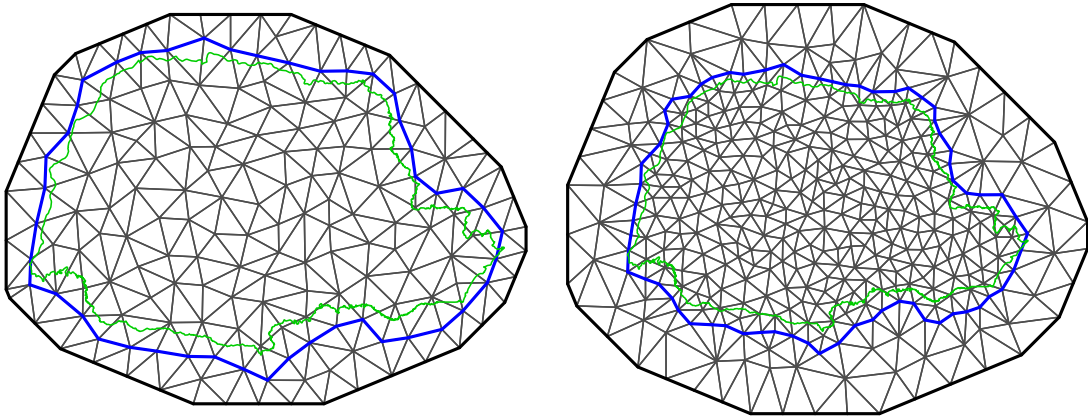


Figure 2.4: Mesh for Paraná state

```

data(PRborder)
dim(PRborder)

```

```
[1] 2055    2
```

that consists of a set of 2055 points on the latlong projection.

In this case is best to use a non-convex hull mesh. We start it with building a non-convex domain with

```

prdomain <- inla.nonconvex.hull(as.matrix(PRprec[,1:2]),
                                -0.03, -0.05, resolution=c(100,100))

```

with this defined domain we build two meshes with different resolution (max edge length) on the inner domain

```

(prmesh1 <- inla.mesh.2d(boundary=prdomain, max.edge=c(.7,.7),
                        cutoff=0.35, offset=c(-0.05, -0.05)))$n

```

```
[1] 187
```

```

(prmesh2 <- inla.mesh.2d(boundary=prdomain, max.edge=c(.45,1), cutoff=0.2))$n

```

```
[1] 382
```

We can visualize this both meshes on the Figure 2.4 with commands bellow

```

par(mfrow=c(1,2), mar=c(0,0,0,0))
plot(prmesh1, asp=1, main=""); lines(PRborder, col=3)
plot(prmesh2, asp=1, main=""); lines(PRborder, col=3)

```

### 2.2.5 Triangulation with a SpatialPolygonsDataFrame

Suppose that we have a map of the domain region. In **R** the representation of an spatial object is made by object classes on **sp** package, see [Pebesma and Bivand, 2005] and [Bivand et al., 2008]. To show an application in this case, we use the North Carolina map, on package **spdep**, [Bivand et al., 2012].

```
require(maptools)
nc.fl <- system.file("etc/shapes/sids.shp", package="spdep")[1]
nc.sids <- readShapePoly(nc.fl, ID="FIPSNO",
                        proj4string=CRS("+proj=longlat +ellps=clrk66"))
```

We simplify this map uniting all the areas together. To do it, we use the `unionSpatialPolygons()` **spdep** function that uses function of the **rgeos** package [Bivand and Rundel, 2013]. If we don't have the **rgeos** we can use the **gpclib** package, but we need to set a permission before the work with this library

```
gpclibPermit()

[1] TRUE

nc.border <- unionSpatialPolygons(nc.sids, rep(1, nrow(nc.sids)))
```

Now, we use the `inla.sp2segment()` to extract the boundary of the **SpatialPolygons** object that contains the border of the map

```
require(INLA)
nc.bdry <- inla.sp2segment(nc.border)
```

and creates the mesh

```
(nc.mesh <- inla.mesh.2d(boundary=nc.bdry, cutoff=0.15,
                        max.edge=c(0.3, 1)))$n

[1] 533
```

that is visualized on Figure 2.5 with the commands bellow.

```
par(mar=c(0,0,0,0))
plot(nc.mesh, asp=1, main="")
```

## 2.3 The projector matrix

Because the SPDE model is defined on the mesh, we need an appropriate specification of the linear predictor for the data response. For details on this please see [Lindgren, 2012].

The key is that we have a random field modeled at the mesh vertices, with dimension  $m$ , and a response at  $n$  locations. So, we need to define how the random field and other model components are linked to the response.

Let the matrix **A** that links the process on the mesh vertices triangles with the locations response. This matrix is called the projector matrix because it do the projection of the latent field modeled at mesh vertices to data locations.

The projector matrix can be builded with the `inla.spde.make.A` function. Considering that each mesh vertices has a weight, the value for one point within one triangle is the projection of the plane (formed by these three weights) at this point location. This projection is just a weighted average with the weights computed by the `inla.spde.make.A()` function. Using the toy data set and the mesh five we have

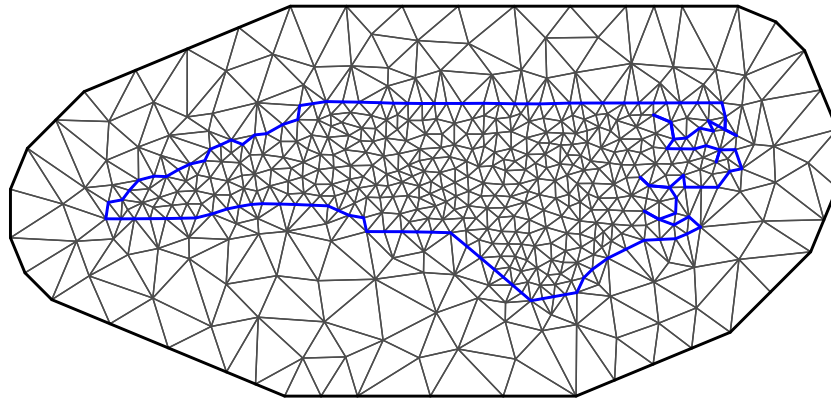


Figure 2.5: Mesh constructed using the North Carolina map

```
coords <- as.matrix(SPDEtoy[,1:2])
A5 <- inla.spde.make.A(mesh5, loc=coords)
```

This matrix has dimension equal to the number of data locations by the number of vertices on the mesh

```
dim(A5)
```

```
[1] 200 617
```

We have exactly three non-zero elements on each line

```
table(rowSums(A5>0))
```

```
3
200
```

because each location point is inside one of the triangles. These three non-zero elements on each line sum to one

```
table(rowSums(A5))
```

```
1
200
```

because multiplication with any vector of weights at mesh nodes by this matrix is the projection of these weights at the location points.

We have some columns on the projector matrix with all elements equal to zero.

```
table(colSums(A5)>0)
```

```
FALSE  TRUE
 332    285
```

This column corresponds to the triangles without any location point inside. This column can be dropped and the stack functionality (section 3.1) automatically handles this situation.

When we have a mesh where each location point is on the vertex set, the projector matrix has exactly one line with non-zero elements

```
A1 <- inla.spde.make.A(mesh1, loc=coords)
table(rowSums(A1>0))
```

```
1
200
```

and all these elements are equals one

```
table(rowSums(A1))
```

```
1
200
```

because in this case the the projection on the location points are just the weight at the corresponding node (at same location) of the mesh.

## 2.4 GRF sampling through SPDE

The main thing is that the SPDE make the model on the mesh vertices and we need to make simulation on a set o location points. We start by defining the spde model

```
spde5 <- inla.spde2.matern(mesh5, alpha=2)
```

and build the precision matrix associated.

It can be done by the `inla.spde2.precision()` function. On this function, we have to provide the SPDE model and the parameters  $\theta$  on the SPDE parametrization.  $\theta$  is a length two vector where the second element is the logarithim of the scale parameter,  $\log(\kappa)$ , and the second one is the 'local variance parameter'  $\tau$  such that the marginal variance  $\sigma_x^2$  is

$$\sigma_x^2 = \frac{1}{4\pi\tau^2\kappa^2}$$

and we have that

$$\log(\tau) = \theta_1 = -\log(4\kappa\sigma_x^2\pi)/2$$

.  
Considering  $\kappa = 5$  and  $\sigma_x^2 = 1$

```
kappa <- 5; s2x <- 1
```

we have

```
theta <- c(-log(2*kappa*s2x*sqrt(pi)), log(kappa))
```

and we have the precision matrix by

```
Q5 <- inla.spde2.precision(spde5, theta)
```

Now, we can use the `inla.qsample()` function to get a sample at mesh nodes by

```
x5 <- inla.qsample(Q=Q5)
```

and to get a sample at the locations, we just project it with the projector matrix by

```
x <- drop(A5%*%x5)
```

We can have a more general function to generate one or more samples, by defining a function that need just the coordinates and model parameters. A mesh can be provided and if not it is also generated within this function

```
rspde
```

```

function (coords, kappa, variance = 1, alpha = 2, n = 1, verbose = TRUE,
        mesh, return.attributes = TRUE)
{
  t0 <- Sys.time()
  theta <- c(NA, log(kappa))
  theta[1] <- -0.5 * log(4 * pi * variance * kappa^2)
  if (verbose)
    cat("theta =", theta, "\n")
  if (missing(mesh)) {
    mesh.pars <- 1/(kappa * c(3/2, 1, 7, 3, 1))
    if (verbose)
      cat("mesh.pars =", mesh.pars, "\n")
    attributes <- list(mesh = inla.mesh.2d(coords, max.edge = mesh.pars[1:2],
      cutoff = mesh.pars[3], offset = mesh.pars[4:5]))
    if (verbose)
      cat("n.mesh =", attributes$mesh$n, "\n")
  }
  else attributes <- list(mesh = mesh)
  attributes$spde <- inla.spde2.matern(attributes$mesh, alpha = alpha)
  attributes$Q <- inla.spde2.precision(attributes$spde, theta = theta)
  attributes$A <- inla.mesh.project(mesh = attributes$mesh,
    loc = coords)$A
  if (n == 1)
    result <- drop(attributes$A %*% inla.qsample(Q = attributes$Q))
  t1 <- Sys.time()
  result <- inla.qsample(n, attributes$Q)
  if (nrow(result) < nrow(attributes$A)) {
    result <- rbind(result, matrix(NA, nrow(attributes$A) -
      nrow(result), ncol(result)))
    dimnames(result)[[1]] <- paste("x", 1:nrow(result), sep = "")
    for (j in 1:ncol(result)) result[, j] <- drop(attributes$A %*%
      result[1:ncol(attributes$A), j])
  }
  else {
    for (j in 1:ncol(result)) result[1:nrow(attributes$A),
      j] <- drop(attributes$A %*% result[, j])
    result <- result[1:nrow(attributes$A), ]
  }
  t2 <- Sys.time()
  attributes$cpu <- c(prepare = t1 - t0, sample = t2 - t1, total = t2 -
    t0)
  if (return.attributes)
    attributes(result) <- c(attributes(result), attributes)
  return(result)
}

```

We can generate samples at a million of locations in some seconds by

```

pts1 <- matrix(runif(2e6), ncol=2)
x1 <- rspde(pts1, kappa=5)

theta = -2.87495 1.609438
mesh.pars = 0.1333333 0.2 0.02857143 0.06666667 0.2
n.mesh = 1018

```



```

summary(x1)

      sample1
Min.      :-1.87441
1st Qu.   :-0.36278
Median    : 0.11221
Mean      : 0.09358
3rd Qu.   : 0.61213
Max.      : 2.07145

table(x1>0)

FALSE     TRUE
441039 558961

and we have the attributes together

names(attributes(x1))

[1] "dim"      "dimnames" "mesh"     "spde"     "Q"        "A"
[7] "cpu"

attr(x1, 'cpu') ### time in seconds

      prep      sample      total
53.586810  4.238619 57.825429

attr(x1, 'mesh')$n ### number of mesh nodes

[1] 1018

dim(attr(x1, 'Q'))

[1] 1018 1018

dim(attr(x1, 'A'))

[1] 1000000 1018

```

This function is used on some of next Chapters.

## 2.5 Maximum likelihood inference

We have that the latent random field is distributed as

$$x \sim N(0, \mathbf{Q}) .$$

Under the toy model where we just add a noise, we have

$$y|x \sim N(\mathbf{A}x, \mathbb{1}\sigma_y^2) .$$

So, it is easy to integrate out  $x$  and get the marginal likelihood and this is

$$y \sim N(\mathbf{A}x, (\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}' + \mathbb{1}\sigma_y^2)) .$$

To eval the marginal likelihood efficiently we need to take care on the evaluation of the determinant of the covariance matrix  $(\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}' + \mathbb{1}\sigma_y^2)$  and also the evaluation of the precision matrix  $(\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}' + \mathbb{1}\sigma_y^2)^{-1}$ .

To evaluate the determinant, we use the matrix determinant lemma that, for any invertible matrix  $\mathbf{W}$  and also invertible matrix  $\mathbf{C}$  we have

$$\det(\mathbf{C} + \mathbf{U}\mathbf{W}\mathbf{V}') = \det(\mathbf{W}^{-1} + \mathbf{V}'\mathbf{C}^{-1}\mathbf{U})\det(\mathbf{W})\det(\mathbf{C}) .$$

Considering  $\mathbf{W} = \mathbf{Q}^{-1}$  and  $\mathbf{C} = \mathbb{1}\sigma_y^2$  we have the following function to implement it

```
prec.det.f <- function(U, Q, tau.y) {
  d1 <- determinant(crossprod(U*tau.y, U) + Q)$modulus
  (d1 - determinant(Q)$modulus) - nrow(U)*log(tau.y)
}
```

where `tau.y` is the likelihood precision.

The precision matrix can be evaluated using the Woodbury matrix identity. It says that for  $C$  and  $W$  invertible matrices we have

$$(C + UWV)^{-1} = C^{-1} - C^{-1}U(W^{-1} + VC^{-1}U)^{-1}VC^{-1}.$$

We implement it with the following function

```
precy.f <- function(U, Q, tau.y) {
  Bi <- Diagonal(nrow(U), rep(tau.y, nrow(U)))
  VBi <- crossprod(U, Bi)
  R <- inla.qinv(Q + VBi%*%U)
  forceSymmetric(Bi - crossprod(VBi, R)%*%VBi)
}
```

Considering  $\theta$  and  $\log(\sigma_y^2)$  the three parameter vector, we have the following negative of the log likelihood function

```
nllf <- function(pars, X, y, A, verbose=TRUE){
  if (verbose) cat(pars, "\n")
  m <- inla.spde2.precision(spde, theta=pars[1:2])
  ldet <- prec.det.f(A, m, exp(-pars[3]))
  m <- precy.f(A, m, exp(-pars[3]))
  Xm <- crossprod(X,m)
  beta <- drop(solve(Xm%*%X, sum(Xm*y)))
  y <- drop(y-X%*%beta)
  ss <- y%*%m%*%y
  r <- drop((ss + length(y)*log(2*pi) + ldet)/2)
  if (verbose) cat(beta, r, "\n")
  attr(r, 'beta') <- beta
  return(r)
}
```

We proceed the maximum likelihood estimation, for the toy example dataset and compare with the results from **geoR** package of previous Chapter. We remember also the true parameter values

```
beta0 <- 10; sigma2e <- 0.3; sigma2x <- 5; kappa <- 7; nu <- 1
```

We call the data

```
data(SPDEtoy)
```

build the mesh (using the true scale parameter to determine the edge lengths and `cutoff`), SPDE model (with  $\alpha = 2$ ) and the projector matrix

```
xy <- as.matrix(SPDEtoy[,1:2])
(mesh <- inla.mesh.2d(xy, max.edge=1/(c(1.7,0.7)*kappa),
  cutoff=1/(5*kappa)))$n
```

[1] 573

```
spde <- inla.spde2.matern(mesh)
A <- inla.mesh.project(mesh, loc=xy)$A
```

The maximum likelihood estimation is done using the `optim()` function. Using the default Nelder-Mead method

```
opt <- optim(c(-4,2,-1.5), nllf, hessian=TRUE,
            X=matrix(1,nrow(xy)), y=SPDEtoy[,3], A=A, verbose=FALSE)
```

and the BFGS method

```
opt.b <- optim(c(-4,2,-1.5), nllf, hessian=TRUE, method='BFGS',
            X=matrix(1,nrow(xy)), y=SPDEtoy[,3], A=A, verbose=FALSE)
```

we have the same results

```
rbind(opt$par, opt.b$par)

      [,1]      [,2]      [,3]
[1,] -3.814219  1.806219 -1.343733
[2,] -3.814413  1.806650 -1.343522
```

We also define a function to map the SPDE model parameters to the  $\sigma_x^2$  and  $\kappa$  user parameters. The function below receive  $\theta$ ,  $\log(\sigma_y^2)$  and  $\beta$  (regression parameters) and return  $\sigma_x^2$ ,  $\kappa$ ,  $\sigma_y^2$  and  $\beta$ .

```
par2user <- function(pars, A, X, y) {
  m <- inla.spde2.precision(spde, theta=pars[1:2])
  m <- precy.f(A, m, exp(-pars[3]))
  Xm <- crossprod(X,m)
  beta <- drop(solve(Xm%%X, sum(Xm*y)))
  k2s22 <- exp(sum(pars[1:2])*2)
  c(s2x=1/(4*pi*k2s22), kappa=exp(pars[2]), s2e=exp(pars[3]), beta=beta)
}
```

We use this function to compare the results with the results by `geoR` package. We have

```
rbind(spde=par2user(opt$par, A=A, X=matrix(1, nrow(A)), y=SPDEtoy[,3]),
      geoR=c(lkf$sigmasq, 1/lkf$phi, lkf$nugget, lkf$beta),
      true=c(sigma2x, kappa, sigma2e, beta0))

      s2x      kappa      s2e      beta
spde 4.414865 6.087385 0.2608701 9.77935
geoR 3.323371 8.653658 0.2708727 9.53488
true 5.000000 7.000000 0.3000000 10.00000
```

## Chapter 3

# A toy example

In this section we start to show the fitting process using the SPDE approach, [Lindgren et al., 2011]. This starter is by fitting a toy geostatistical model: Gaussian response without covariate. We use the Bayesian approach and found the posterior marginal distributions using the Integrated Nested Laplace Approximation - INLA, [Rue et al., 2009] implemented on the **INLA R** package. The ideas for application of the SPDE approach using the **INLA** package are well described on [Lindgren, 2012] and on [Lindgren and Rue, 2013].

The dataset used are a tree column `data.frame` simulated on the previous section and provided on **INLA** package. It can be called by

```
data(SPDEtoy)
```

this is a `data.frame` where the two first columns are the coordinates and the third is the response simulated at this locations

```
str(SPDEtoy)
```

```
'data.frame':      200 obs. of  3 variables:
 $ s1: num   0.0827 0.6123 0.162 0.7526 0.851 ...
 $ s2: num   0.0564 0.9168 0.357 0.2576 0.1541 ...
 $ y : num   11.52 5.28 6.9 13.18 14.6 ...
```

We consider the  $n$  observations  $y_i$  on locations the  $s_i$ ,  $i = 1, \dots, n$ , and we define the model

$$\begin{aligned} y_i | \beta_0, x_i, \sigma_e^2 &\sim N(\beta_0 + x_i, \sigma_e^2) \\ \mathbf{x} &\sim GF(0, \Sigma) \end{aligned} \quad (3.1)$$

We consider that  $x$  is a realization of a Gaussian Field, with Matérn correlation function parametrized by the smoothness parameter  $\nu$  and the scale  $\kappa$ , such the parametrization in [Lindgren et al., 2011].

With this toy example we show with details how we make a good triangulation, prepare the data, fit the model, extract the results from output and make predictions on locations where we don't have observed the response. In this section we use the default priors for all the parameters.

To the estimation mesh we need to define the mesh. We show it on Chapter 2.2 and here we use the fit mesh builded for the toy example on Chapter 2.2

With that mesh, we define the SPDE model using the function `inla.spde2.matern()`

```
args(inla.spde2.matern)
```

```
function (mesh, alpha = 2, param = NULL, constr = FALSE, extraconstr.int = NULL,
  extraconstr = NULL, fractional.method = c("parsimonious",
  "null"), B.tau = matrix(c(0, 1, 0), 1, 3), B.kappa = matrix(c(0,
  0, 1), 1, 3), prior.variance.nominal = 1, prior.range.nominal = NULL,
```

```
prior.tau = NULL, prior.kappa = NULL, theta.prior.mean = NULL,
theta.prior.prec = 0.1)
NULL
```

The principal arguments are the mesh object and the  $\alpha$  parameter, related to the smoothness parameter of the process.

The toy dataset was simulated with  $\alpha = 2$  and we use this value here (with the mesh five)

```
spde5 <- inla.spde2.matern(mesh5, alpha=2)
```

Also, from section 2.3 we define the projector matrix

```
coords <- as.matrix(SPDEtoy[,1:2])
A5 <- inla.spde.make.A(mesh5, loc=coords)
```

### 3.1 The stack functionality

The stack functionality is a very useful functionality to work with SPDE on **INLA** package. This allow to fit more general SPDE models, such as replicated or correlated ones, or more general models that includes more than one projector matrix. Using it we avoid errors in the index construction, [Lindgren, 2012]. Examples on more complex models, with the details, can be found on [Lindgren, 2012], [Cameletti et al., 2012] and [Lindgren and Rue, 2013].

In a section on the previous Chapter we define the projector matrix to project the latent field on the response locations. If we have have covariates measured at same locations (with same dimation of the response), we need a more general definition of the linear predictor. On the toy example, we have to include the effect of the random field and the intercept that is treated as a covariate. So, we define the new intercept  $\eta^*$  as

$$\eta^* = \mathbf{A}\mathbf{x} + \mathbf{1}\beta_0$$

that is a sum of two components each one can be represented as a product of a projector matrix and an effect.

The `inla.stack` function allow us to work 'easily' with more than one responses. The three main `inla.stack()` arguments are the `data` vectors list, a list of projector matrices (each one related to one block effect) and the effects.

We need two projector matrices, the projector matrix for the latend field and a matrix to map one-to-one the 'covariate' and the response. This last one can be just a constant instead a diagonal matrix. So, we have

```
stk5 <- inla.stack(data=list(resp=SPDEtoy$y), A=list(A5,1),
effects=list(i=1:spde5$n.spde,
m=rep(1,nrow(SPDEtoy))), tag='est')
```

The `inla.stack()` function automatically eliminates the elements when any column of each projector matrix has zero sum, generating a correspondent simplified projector matrix. The `inla.stack.A()` extracts a simplified predictor matrix to use on the `inla()` function and the `inla.stack.data()` function extract the correspondent organized data.

The simplified projector matrix from the stack is the binded by column the simplified projectors matrices from each effect block. So, in this case we have

```
dim(inla.stack.A(stk5))
```

```
[1] 200 286
```

one columns more than the number of columns with non null elements of the projector matrix.

## 3.2 Model fitting and some results

To fit the model need to remove the intercept from the formulae and add it as a covariate term, because we have projector matrix that allows it as a covariate effect. Of course, we must have to pass these predictor matrix on `control.predictor` argument of the `inla` function

```
res5 <- inla(resp ~ 0 + m + f(i, model=spde5),
             data=inla.stack.data(stk5),
             control.predictor=list(A=inla.stack.A(stk5)))
```

An object from `inla()` function has a set of several results. These includes summaries, marginal posterior densities of each parameter on the model: the regression parameters, each element of the latent field and all the hyperparameters.

The summary of  $\beta_0$  is obtained by

```
res5$summary.fix
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode	kld
m	9.541585	0.743176	7.985238	9.555832	11.014	9.578793	1.102504e-08

The summary of  $1/\sigma_e^2$  is obtained by

```
res5$summary.hy[1,]
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
	2.8801009	0.4565219	2.0786390	2.8480014	3.8684182	2.7876490

A marginal distribution on `inla()` output is just two vector, where one represents the parameter values and another the density. Any posterior marginal can be transformed. If we want the posterior marginal for  $\sigma_e$ , the square root of the  $\sigma_e^2$ , we use

```
post.se <- inla.tmarginal(function(x) sqrt(1/x),
                          res5$marginals.hy[[1]])
```

and now we are able to summarize this distribution

```
inla.emarginal(function(x) x, post.se)
```

```
[1] 0.5947447
```

```
inla.qmarginal(c(0.025, 0.5, 0.975), post.se)
```

```
[1] 0.5091092 0.5924102 0.6922622
```

```
inla.hpdmarginal(0.95, post.se)
```

	low	high
level:0.95	0.5052961	0.6877207

```
inla.pmargin(c(0.5, 0.7), post.se)
```

```
[1] 0.0137697 0.9822753
```

and, of course, we can visualize it.

The parameters of the latent field is parametrized as  $\log(\kappa)$  and  $\log(\tau)$ , where  $\tau$  is the local variance parameter. We have the posterior marginals for  $\kappa$ ,  $\sigma_x^2$  and for the nominal range (the distance that we have correlation equals 0.1). This can be done with the `inla.spde2.result` function

```
res5.field <- inla.spde2.result(res5, 'i', spde5, do.transf=TRUE)
```

and we get the posterior mean of each of these parameters by

```
inla.emarginal(function(x) x, res5.field$marginals.kappa[[1]])
```

```
[1] 7.777047
```

```
inla.emarginal(function(x) x, res5.field$marginals.variance.nominal[[1]])
```

```
[1] 3.952283
```

```
inla.emarginal(function(x) x, res5.field$marginals.range.nominal[[1]])
```

```
[1] 0.3818887
```

also we can get other summary statistics, HPD interval and visualize it.

### 3.3 Prediction of the random field

A very common objective when we have spatial data collected on some locations is the prediction on a fine grid to get high resolution maps. In this section we show two approaches to make prediction of the random field, one is after the estimation process and other is jointly on estimation process. To compare both approaches, we predict the random field on three target locations: (0.1,0.1), (0.5,0.55), (0.7,0.9).

```
pts3 <- rbind(c(.1,.1), c(.5,.55), c(.7,.9))
```

#### 3.3.1 Jointly with the estimation process

The prediction of the random field joint the parameter estimation process in Bayesian inference is the common approach. This approach is made by the computation of the marginal posterior distribution of the random field at target locations. If the target points are on the mesh, so we have automatically this distribution. If the target points are not on the mesh, we must define the projector matrix for the target points.

The predictor matrix for the target locations is

```
dim(A5pts3 <- inla.spde.make.A(mesh5, loc=pts3))
```

```
[1] 3 617
```

We can show the columns with non-zero elements of this matrix

```
(jj3 <- which(colSums(A5pts3)>0))
```

```
[1] 149 154 165 245 250 289 339 398 486
```

```
round(A5pts3[, jj3],3)
```

```
3 x 9 sparse Matrix of class "dgCMatrix"
```

```
[1,] .      0.094 .      0.513 .      0.393 .      .      .
[2,] .      .      0.106 .      0.616 .      .      0.278 .
[3,] 0.035 .      .      .      .      .      0.42 .      0.544
```

We have to define a data stack for the prediction and join it with the data stack of the observations. The prediction data stack contains the effect set, predictor matrices and assign NA to response

```
stk5p.rf <- inla.stack(data=list(resp=NA), A=list(A5pts3),
                      effects=list(i=1:spde5$n.spde), tag='prd5r')
```

Also, we join both stacks by

```
stk5.jp <- inla.stack(stk5, stk5p.rf)
```

and fit the model again with the full stack setting `compute=TRUE` on `control.predictor`

```
res5p <- inla(resp ~ 0 + m + f(i, model=spde5),
              data=inla.stack.data(stk5.jp),
              control.predictor=list(A=inla.stack.A(stk5.jp), compute=TRUE))
```

To access the posterior marginal distribution of the random field at the target locations, we extract the index from the full stack using the adequate `tag`.

```
(indd5p <- inla.stack.index(stk5.jp, tag='prd5r')$data)
```

```
[1] 201 202 203
```

The summary of the posterior distributions of the random field on the target locations is

```
round(res5p$summary.linear.pred[indd5p,], 4)
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode	kld
Apredictor.201	0.4562	0.8390	-1.1842	0.4433	2.1808	0.4232	0
Apredictor.202	2.9039	1.0673	0.8272	2.8927	5.0504	2.8731	0
Apredictor.203	-2.8149	1.2707	-5.3364	-2.8092	-0.3245	-2.7980	0

that includes the posterior mean, standard deviation, quantiles and mode.

Because it is a full bayesian analysis, we also we have the marginal distributions. We extract the marginals posterior distributions with

```
marg3 <- res5p$marginals.linear[indd5p]
```

and get the 95% HPD interval for the random field at the second target location by

```
inla.hpdmarginal(0.95, marg3[[2]])
```

	low	high
level:0.95	0.8009717	5.021002

and see that around the point (0.5,0.5) the random field has positive values, see Figure 3.1.

### 3.3.2 After the estimation process

If we need just the prediction we can do the prediction after the estimation process with a very small computational cost. It is just a matrix operation in way that we just project the posterior mean of the the random field on mesh nodes to target locations, using the correspondent projector matrix.

So, we 'project' the posterior mean of the latend random field to the target locations by

```
drop(A5pts3*%res5$summary.random$i$mean)
```

```
[1] 0.4545881 2.9023932 -2.8156423
```

or using the `inla.mesh.projector()` function



```
inla.mesh.project(inla.mesh.projector(mesh5, loc=pts3),
                  res5$summary.random$i$mean)
```

```
[1] 0.4545881 2.9023932 -2.8156423
```

and see that for the mean we have similar values than those on previous section.

Also, we can get the standard deviation

```
drop(A5pts3%%res5$summary.random$i$sd)
```

```
[1] 0.9305767 1.1912548 1.3592889
```

and we have a little difference.

```
sqrt(drop((A5pts3^2)%%(res5$summary.random$i$sd^2)))
```

```
[1] 0.6172511 0.8698557 0.9370294
```

### 3.3.3 Projection on a grid

The approach by the projection of the posterior mean random field is computationally cheap. So, it can be used to get the map of the random field on a fine grid. The `inla.mesh.projector()` function get the projector matrix automatically for a grid of points over a square that contains the mesh.

To get projection on a grid at the domain  $(0, 1) \times (0, 1)$  we just inform these limits

```
pgrid0 <- inla.mesh.projector(mesh5, xlim=0:1, ylim=0:1, dims=c(101,101))
```

and we project the posterior mean and the posterior standard deviation on the both grid with

```
prd0.m <- inla.mesh.project(pgrid0, res5$summary.ran$i$mean)
prd0.s <- inla.mesh.project(pgrid0, res5$summary.ran$i$s)
```

We visualize this values projected on the grid on Figure 3.1.

## 3.4 Prediction of the response

Another common result that we want on spatially continuous modelling is the prediction of the response on a target locations that we don't have data observed. In similar way that on past section, it is possible to find the marginal distribution or to make a projection of some functional of the response.

### 3.4.1 By the posterior distribution

In this case, we want to define a adequate predictor of the response and build the model again. This is similar to the stack to predict the random field, but here we add the intercept on the list of predictor matrix and on the list of effects

```
stk5.presp <- inla.stack(data=list(resp=NA), A=list(A5pts3,1),
                        effects=list(i=1:spde5$n.spde, m=rep(1,3)),
                        tag='prd5.resp')
```

and join with the data stack to build the model again

```
stk5.full <- inla.stack(stk5, stk5.presp)
r5presp <- inla(resp ~ 0 + m + f(i, model=spde5),
               data=inla.stack.data(stk5.full),
               control.predictor=list(A=inla.stack.A(stk5.full), compute=TRUE))
```

We find the index of the predictor that corresponds the predicted values of the response on the target locations. We extract the index from the full stack by

```
(indd3r <- inla.stack.index(stk5.full, 'prd5.resp')$data)

[1] 201 202 203
```

To get the summary of the posterior distributions of the response on target locations we do

```
round(r5presp$summary.fitted.values[indd3r,], 3)
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
fitted.Apredictor.201	9.996	0.385	9.235	9.998	10.750	10.001
fitted.Apredictor.202	12.444	0.770	10.933	12.443	13.961	12.442
fitted.Apredictor.203	6.726	1.050	4.668	6.723	8.802	6.718

Also, we extract the marginals posterior distributions with

```
marg3r <- r5presp$marginals.fitted.values[indd3r]
```

and get the 95% HPD interval for the response at second target location by

```
inla.hpdmarginal(0.95, marg3r[[2]])
```

	low	high
level:0.95	10.92848	13.95516

and see that around the point (0.5,0.5) we have the values of the response significantly larger than  $\beta_0$ , see Figure 3.1.

### 3.4.2 By sum of linear predictor components

A computational cheap approach is to (naively) sum the projected posterior mean to the regression term. In this toy example we just sum the posterior mean of the intercept to the posterior mean of the random field to get the posterior mean of the response.

If there are covariates, the prediction also can be made in similar way, see . That approach can be used here considering just the intercept

```
res5$summary.fix[1,1] + drop(A5pts3%*%res5$summary.random$i$mean)
```

```
[1] 9.996173 12.443978 6.725942
```

For the standard error, we need to take into account the error of the covariate values and regression coefficients.

```
summary(rvar <- res5$summary.random$i$sd^2)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.6119	1.0460	1.3780	2.3060	2.2310	11.2700

```
sqrt(1^2+res5$summary.fix[1,2]^2 + drop(A5pts3%*%rvar))
```

```
[1] 1.556573 1.734277 1.844263
```

### 3.4.3 Response on a grid

The computation of all marginal posterior distributions on a grid is computationally expensive. But, we usually not uses the marginal distributions. We usually uses just the mean and standard deviation. So, we don't need the storage of all the marginal distributions! Also, we don't need the quantiles of the marginal distributions.

On the code below, we build the model again but we disable the storage of the marginal posterior distributions to random effects and to posterior predictor values. Also, we disable the computation of the quantiles. Only the mean and standard defiation are stored.

We use the projector matrix on the projector object that we use to project the posterior mean on the grid

```
stkgrid <- inla.stack(data=list(resp=NA), A=list(pgrid0$proj$A,1),
                    effects=list(i=1:spde5$n.spde,
                                m=rep(1,101*101)), tag='prd.gr')
stk.all <- inla.stack(stk5, stkgrid)
res5g <- inla(resp ~ 0 + m + f(i, model=spde5),
              data=inla.stack.data(stk.all),
              control.predictor=list(A=inla.stack.A(stk.all),
                                    compute=TRUE), quantiles=NULL,
              control.results=list(return.marginals.random=FALSE,
                                  return.marginals.predictor=FALSE))
res5g$cpu
```

Pre-processing	Running inla	Post-processing	Total
0.7143722	11.8509769	0.1831388	12.7484879

We get the indexes

```
igr <- inla.stack.index(stk.all, 'prd.gr')$data
```

and use it to visualize, together the prediction of the random field on previous section, on Figure 3.1 with the commands bellow

```
require(gridExtra)
grid.arrange(levelplot(prd0.m, col.regions=topo.colors(99), main='latent field mean',
                      xlab='', ylab='', scales=list(draw=FALSE)),
              levelplot(matrix(res5g$summary.fitt[igr,1], 101),
                      xlab='', ylab='', main='response mean',
                      col.regions=topo.colors(99), scales=list(draw=FALSE)),
              levelplot(prd0.s, col.regions=topo.colors(99), main='latent field SD',
                      xlab='', ylab='', scales=list(draw=FALSE)),
              levelplot(matrix(res5g$summary.fitt[igr,2], 101),
                      xlab='', ylab='', main='response SD',
                      col.regions=topo.colors(99), scales=list(draw=FALSE)),
              nrow=2)
```

We see on Figure 3.1 that we have a variation from -4 to 4 on the spatial effect. Considering also that we have standard deviations around 0.8 to 1.6, the spatial dependence is significantly.

Another thing is that the standard deviation of both, random field and the response, are less near the corner (0, 0) and greater near the corner (1,1). This is just proportional to the locations density.

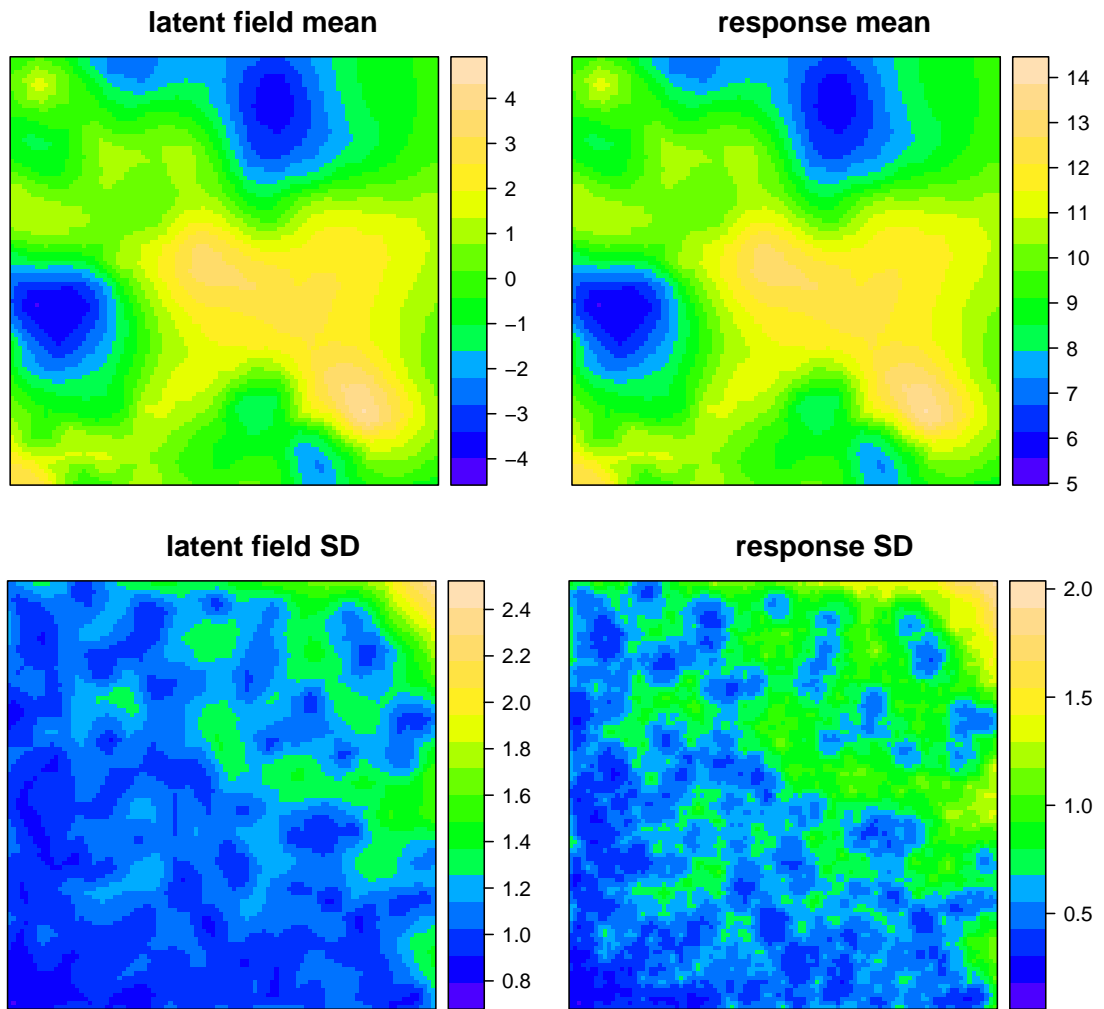


Figure 3.1: The mean and standard deviation of the random field (top left and bottom left) and the mean and standard variation of the response (top right and bottom right)

### 3.5 Results from different meshes

In this section we compare six results for the toy dataset based on the six different meshes builded on section 2.2. To do this comparison, we just plot the posterior marginal distributions of the model parameters. We evaluate the meshes by the addition of the true values used on the simulation of the toy dataset. Also, we add the maximum likelihood estimates from **geoR** package, [Ribeiro Jr and Diggle, 2001].

We fit the model, using each one of the six meshes, and put the results in a list with the code bellow

```
lrf <- lres <- l.dat <- l.spde <- l.a <- list()
for (k in 1:6) {
  l.a[[k]] <- inla.spde.make.A(get(paste('mesh', k, sep='')), loc=coords)
  l.spde[[k]] <- inla.spde2.matern(get(paste('mesh', k, sep='')), alpha=2)
  l.dat[[k]] <- list(y=SPDEtoy[,3], i=1:ncol(l.a[[k]]),
                    m=rep(1, ncol(l.a[[k]])))
  lres[[k]] <- inla(y ~ 0 + m + f(i, model=l.spde[[k]]),
                  data=l.dat[[k]], control.predictor=list(A=l.a[[k]]))
  lrf[[k]] <- inla.spde2.result(lres[[k]], 'i', l.spde[[k]], do.transf=TRUE)
}
```

The mesh size influences the computational time needed to fit the model. More nodes on the mesh need more computational time. The time running inla for these six meshes are

```
round(sapply(lres, function(x) x$cpu[2]), 2)

Running inla Running inla Running inla Running inla Running inla
      14.05      3.07      2.15      11.11      2.22
Running inla
      2.07
```

We compute the distribution for  $\sigma_e^2$  for each fitted model

```
s2.marg <- lapply(lres, function(m)
  inla.tmarginal(function(x) 1/x, m$marginals.hy[[1]]))
```

The true values are:  $\beta_0 = 10$ ,  $\sigma_e^2 = 0.3$ ,  $\sigma_x^2 = 5$ ,  $\kappa = 7$  and  $\nu = 1$ . The  $\nu$  parameter is fixed on the true value when we define  $\alpha = 2$  on definition of the SPDE model.

```
beta0 <- 10; sigma2e <- 0.3; sigma2x <- 5; kappa <- 7; nu <- 1
```

and the maximum likelihood estimates are

```
lk.est

      beta      s2e      s2x      kappa
9.6779236 0.2889533 5.6276898 5.4072950
```

We want to visualize the posterior marginal distributions for  $\beta_0$ ,  $\sigma_e^2$ ,  $\sigma_x^2$ ,  $\kappa$ , nominal range and the local variance  $\tau$ . This can be done with the code bellow

```
rcols <- rainbow(6)##c(rgb(4:1/4,0:3/5,0), c(rgb(0,0:3/5,4:1/4)))
par(mfrow=c(2,3), mar=c(2.5,2.5,1,.5), mgp=c(1.5,.5,0), las=1)
xrange <- range(sapply(lres, function(x) range(x$marginals.fix[[1]][,1])))
yrange <- range(sapply(lres, function(x) range(x$marginals.fix[[1]][,2])))
plot(lres[[1]]$marginals.fix[[1]], type='l', xlim=xrange, ylim=yrange,
     xlab=expression(beta[0]), ylab='Density')
```

```

for (k in 1:6)
  lines(lres[[k]]$marginals.fix[[1]], col=rcols[k], lwd=2)
abline(v=beta0, lty=2, lwd=2, col=3)
abline(v=lk.est[1], lty=3, lwd=2, col=3)
xrange <- range(sapply(s2.marg, function(x) range(x[,1])))
yrange <- range(sapply(s2.marg, function(x) range(x[,2])))
plot.default(s2.marg[[1]], type='l', xlim=xrange, ylim=yrange,
             xlab=expression(sigma[e]^2), ylab='Density')
for (k in 1:6)
  lines(s2.marg[[k]], col=rcols[k], lwd=2)
abline(v=sigma2e, lty=2, lwd=2, col=3)
abline(v=lk.est[2], lty=3, lwd=2, col=3)
xrange <- range(sapply(lrf, function(r) range(r$marginals.variance.nominal[[1]][,1])))
yrange <- range(sapply(lrf, function(r) range(r$marginals.variance.nominal[[1]][,2])))
plot(lrf[[1]]$marginals.variance.nominal[[1]], type='l',
     xlim=xrange, ylim=yrange, xlab=expression(sigma[x]^2), ylab='Density')
for (k in 1:6)
  lines(lrf[[k]]$marginals.variance.nominal[[1]], col=rcols[k], lwd=2)
abline(v=sigma2x, lty=2, lwd=2, col=3)
abline(v=lk.est[3], lty=3, lwd=2, col=3)
xrange <- range(sapply(lrf, function(r) range(r$marginals.kappa[[1]][,1])))
yrange <- range(sapply(lrf, function(r) range(r$marginals.kappa[[1]][,2])))
plot(lrf[[1]]$marginals.kappa[[1]], type='l',
     xlim=xrange, ylim=yrange, xlab=expression(kappa), ylab='Density')
for (k in 1:6)
  lines(lrf[[k]]$marginals.kappa[[1]], col=rcols[k], lwd=2)
abline(v=kappa, lty=2, lwd=2, col=3)
abline(v=lk.est[4], lty=3, lwd=2, col=3)
xrange <- range(sapply(lrf, function(r) range(r$marginals.range.nominal[[1]][,1])))
yrange <- range(sapply(lrf, function(r) range(r$marginals.range.nominal[[1]][,2])))
plot(lrf[[1]]$marginals.range.nominal[[1]], type='l',
     xlim=xrange, ylim=yrange, xlab='nominal range', ylab='Density')
for (k in 1:6)
  lines(lrf[[k]]$marginals.range.nominal[[1]], col=rcols[k], lwd=2)
abline(v=sqrt(8)/kappa, lty=2, lwd=2, col=3)
abline(v=sqrt(8)/lk.est[4], lty=3, lwd=2, col=3)
xrange <- range(sapply(lrf, function(r) range(r$marginals.tau[[1]][,1])))
yrange <- range(sapply(lrf, function(r) range(r$marginals.tau[[1]][,2])))
plot(lrf[[1]]$marginals.tau[[1]], type='l',
     xlim=xrange, ylim=yrange, xlab=expression(tau), ylab='Density')
for (k in 1:6)
  lines(lrf[[k]]$marginals.tau[[1]], col=rcols[k], lwd=2)
legend('topright', c(paste('mesh', 1:6, sep=''), 'True', 'Likelihood'),
      lty=c(rep(1,6), 2, 3), lwd=rep(2, 6), col=c(rcols,3,3), bty='n')

```

At the Figure 3.2 we can see that the posterior marginal distribution for the intercept has mode on the likelihood estimate, considering the results from all six meshes.

1/kappa

[1] 0.1428571

The main differences are on the noise variance  $\sigma_e^2$  (the nugget effect). The result from the mesh based on the points and with small triangles mode less than the likelihood estimate, the second has mode near likelihood estimate and the third large. Considering

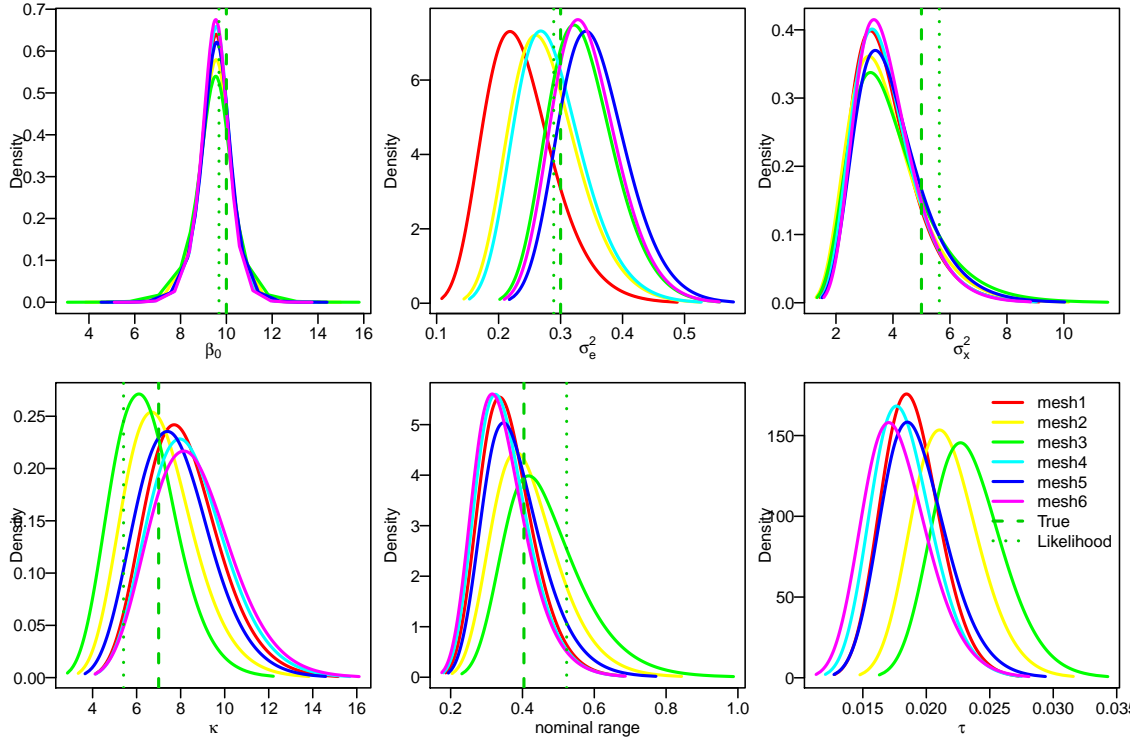


Figure 3.2: Marginal posterior distribution for  $\beta_0$  (top left),  $\sigma_e^2$  (top mid),  $\sigma_x^2$  (top right),  $\kappa$  (bottom left), nominal range (bottom mid) and  $\tau$  (bottom right).

the other meshes, the mesh four has mode around likelihood estimate and the other two little larger, similar to the third mesh, such is based on points but with some freedom (cutoff greater than zero).

For the marginal variance of the latent field,  $\sigma_x^2$ , the results with all meshes had mode near the likelihood estimate. For the scale parameter  $\kappa$  all meshes has mode less than the likelihood estimate. The posterior distribution from the meshes based on points are that ones with less mode and that the mode from third mesh are the less. For the practical range the opposite happens.

These results are not conclusive, but a general comment is that is good to have a mesh with some tune on the points locations, to access noise variance, but with some flexibility to avoid many variability on the triangles size and shape, to get good latent field parameters estimation.

## Chapter 4

# Non-Gaussian response: Precipitation on Paraná

A very common data in spatial statistics is the climate data. On this section we analyze a data set from Water National Agency in Brazil, in Portuguese it is *Agencia Nacional de Águas* - ANA. The ANA collect data on many locations over Brazil. All these data are freely available from the ANA website.

### 4.1 The data set

It have data on more one thousand of locations within Paraná state, a state at south of Brazil. We have daily rainfall data on each day of the year 2011 on 616 locations, including stations within the Paraná state and around its border.

We have these dataset on the **INLA** package and call it with

```
data(PRprec)
```

We have the coordinates at first two columns, altitude at third column and more 365 columns, one for each day with the daily accumulated precipitation.

We show bellow some data of four stations: the with missing altitude with less latitude, the stations with extremes longitudes and the station with greater altitude.

```
PRprec[ii <- c(which(is.na(PRprec$A))[which.min(PRprec$La[is.na(PRprec$A)])],  
               which(PRprec$Lo%in%range(PRprec$Lo)), which.max(PRprec$A)), 1:10]
```

	Longitude	Latitude	Altitude	d0101	d0102	d0103	d0104	d0105	d0106	d0107
1239	-48.9394	-26.1800	NA	20.5	7.9	8.0	0.8	0.1	15.6	31.0
658	-48.2167	-25.0831	9	18.1	8.1	2.3	11.3	23.6	0.0	22.6
1325	-54.4842	-25.6017	231	43.8	0.2	4.6	0.4	0.0	0.0	0.0
885	-51.5167	-25.7331	1446	0.0	14.7	0.0	0.0	28.1	2.5	26.6

We visualize this four stations in red points on the graph at right of Figure 4.1.

We have some problems on this data set. For example, we have a spatio-temporal data. Also, there are seven stations with missing altitude and missing data on daily rainfall. The red stations with missing altitude is the red points on graph at left of Figure 4.1. It's easy to get the altitude only from coordinates information. For example, using an digital elevation model, google earth or another on line information. For example, the station with missing altitude located at south and out of border of Paraná state is located on a coast city. So in this case it is reasonable to assign a small value. But, it is possible to build a stochastic model for altitude and predict these missing values.

However, the principal problem is that we have 47 stations with missing data and that 10 stations on 45 or more days. This problem is treated on next section. In this section,



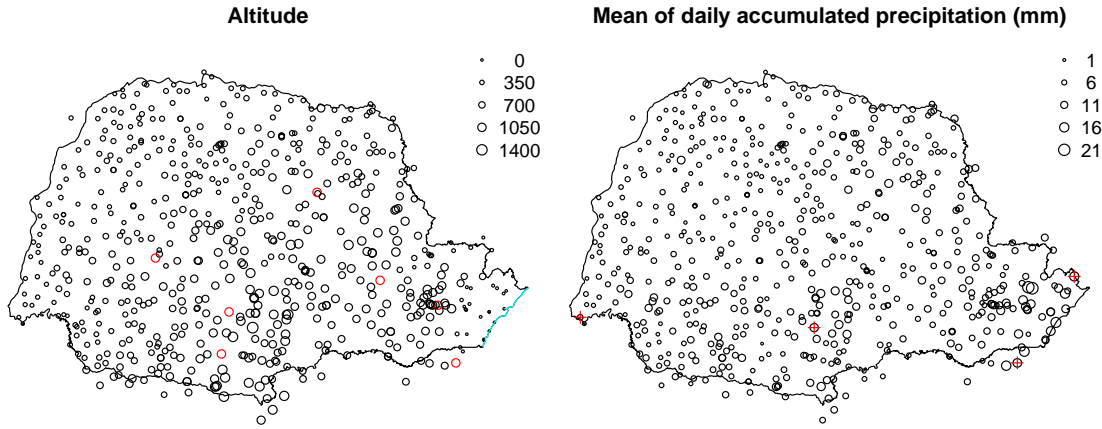


Figure 4.1: Locations of Paraná stations, altitude and average of daily accumulated precipitation (mm) on January.

we analyze the mean of daily accumulated precipitation on January of 2012. We compute this covariate with

```
PRprec$precMean <- rowMeans(PRprec[,3+1:31], na.rm=TRUE)
```

Also we have the Paraná state border

```
data(PRborder)
```

We visualize the locations on Figure 4.1, with the commands bellow

```
par(mfrow=c(1,2), mar=c(0,0,2,0))
plot(PRborder, type='l', asp=1, axes=FALSE, main='Altitude')
points(PRprec[1:2], col=is.na(PRprec$Alt)+1,
       cex=ifelse(is.na(PRprec$Alt), 1, .3+PRprec$Alt/1500))
legend('topright', format(0:4*350), bty='n', pch=1, pt.cex=.3+0:4*35/150)
lines(PRborder[1034:1078, ], col='cyan')
plot(PRborder, type='l', asp=1, axes=FALSE,
     main=paste('Mean of daily accumulated precipitation (mm)'))
points(PRprec[1:2], cex=0.3+PRprec$precMean/20)
legend('topright', format(seq(1,21,5)),
     bty='n', pch=1, pt.cex=0.3+seq(1,21,5)/20)
points(PRprec[1:2], pch=3, col=2)
```

The size of the points on left graph are proportional to altitude of the locations. The cyan line in this graph is the Paraná boundary with the Atlantic Ocean. So, we have height altitudes at mid and south, low altitudes near the sea and the altitude decrease at north and west of Paraná state. On the right graph, the points size are proportional to average of daily accumulated precipitation. We see that near the coast we have large values.

## 4.2 The model and covariate selection

In this section, we analyse the average of daily accumulated precipitation on that period. The average of daily precipitation must be positive. So, we consider in this section the model

$$\begin{aligned} y_i | F_i, \beta, x_i, \theta &\sim \text{Gamma}(a_i, b_i) \\ \log(\mu_i) &= F_i^T \beta + x_i \\ x_i &\sim GF(0, \Sigma) \end{aligned}$$

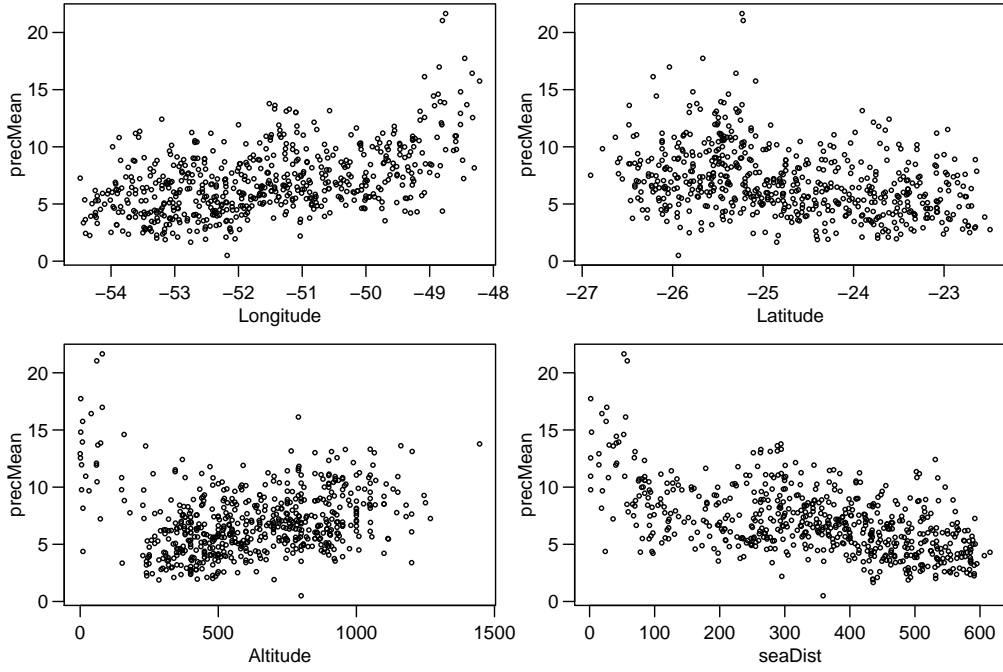


Figure 4.2: Dispersion plots of average of daily accumulated precipitation mean by Longitude (top left), Latitude (top right), Altitude (bottom left) and distance to sea (bottom right).

where,  $F_i$  is a vector of covariates (the location coordinates and altitude) as covariates, with the vector of coefficients  $\beta$ , and a latent Gaussian random field  $x$ , with its covariance matrix function of parameters  $\nu$ ,  $\kappa$  and  $\sigma_x^2$ . Also, with the Gamma likelihood, we consider that  $E(y_i) = \mu_i = a_i/b_i$  and  $V(y_i) = a_i/b_i^2 = \mu_i^2/\phi$ , where  $\phi$  is a precision parameter and we define a linear predictor to  $\log(\mu_i)$ .

To make an initial exploration of relationship between the precipitation and the covariates, we visualize some dispersion diagrams. After preliminary tests, we see that is more adequate the construction of a new covariate: the distance from each station to Atlantic Ocean. We found the coordinates of Paraná state border that share frontier with the sea (plotted as cyan line on left graph at Figure 4.1) and computes the distance from each station to the neighbor coordinate of this line.

We compute the distances in Km units using the `spDists()` function from `sp` package

```
coords <- as.matrix(PRprec[,1:2])
mat.dists <- spDists(coords, PRborder[1034:1078,], longlat=TRUE)
```

This function computes the distance between each location on the first points set to each one on the second points set. So, we need to take the minimum

```
PRprec$"seaDist" <- apply(mat.dists, 1, min)
```

We see the dispersion plots at Figure 4.2.

```
par(mfrow=c(2,2), mar=c(3,3,0.5,0.5), mgp=c(1.7,.7,0), las=1)
for (i in c(1:3, ncol(PRprec))) plot(PRprec[c(i,ncol(PRprec)-1)], cex=.5)
```

With the Figure 4.2, we conclude that have a not well defined non-linear relationship with Longitude, and there is a similar, but inverse, relation with sea distance. So, we build two models, one with longitude as covariate and another with distance to sea as covariate. And compute the DIC to decide what model is better.

To consider a the non-linear relationship, we consider the coefficients of these covariates with a semi-parametric trend on the relationship, such that the coefficients over the range

of the covariate values have a first order random walk prior. So we have the model adopted is of the form

$$\log(\mu_i) = \alpha + \sum_{j=1}^k w_k I(F_i = F0_k) + x_i$$

where  $F$  is the distance to sea or the longitude and  $F0_k$  is a knot (choose on range of the covariate values) and  $w_k$  are regression coefficients with first order random walk prior.

We use the mesh builded for the Paraná state on the Chapter 2.2.

The projector matrix can be done by

```
A <- inla.spde.make.A(prmesh2, loc=coords)
```

The SPDE model is defined by

```
spde <- inla.spde2.matern(prmesh2, alpha=2)
```

and the stack data is defined to include four effects: the GRF, intercept, west coordinate and distance to sea

```
stk.dat <- inla.stack(data=list(y=PRprec$precMean),
                      A=list(A,1), tag='dat',
                      effects=list(list(i=1:spde$n.spde),
                                   list(Intercept=rep(1, nrow(A)),
                                        gWest=inla.group(coords[,1]),
                                        gSeaDist=inla.group(PRprec$seaDist))))
```

We fit the two models using the same stack data. We just use different formula. For the model with west coordinate we have

```
f.west <- y ~ 0 + Intercept + f(gWest, model='rw1') + f(i, model=spde)
r.west <- inla(f.west, family='Gamma', control.compute=list(dic=TRUE),
              data=inla.stack.data(stk.dat),
              control.predictor=list(A=inla.stack.A(stk.dat), compute=TRUE))
```

For the model with distance to sea covariate we have

```
f.sead <- y ~ 0 + Intercept + f(gSeaDist, model='rw1') + f(i, model=spde)
r.sead <- inla(f.sead, family='Gamma', control.compute=list(dic=TRUE),
              data=inla.stack.data(stk.dat),
              control.predictor=list(A=inla.stack.A(stk.dat), compute=TRUE))
```

We have the DIC of each model with

```
c(r.west$dic$dic, r.sead$dic$dic)
```

```
[1] 2567.961 2565.174
```

and choose the model that with distance to sea as covariate.

We got the summary of posterior distribution of the intercept with

```
r.sead$summary.fixed
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
Intercept	1.944487	0.05137997	1.841249	1.944319	2.048408	1.943871

```

      kld
Intercept 1.632191e-09
```

To dispersion parameter of the gamma likelihood we have

```
r.sead$summary.hy[1,]
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
	12.602283	0.886516	10.931736	12.578540	14.415692	12.535246

To dispersion parameter of the coefficients of the sea distance covariate we have

```
r.sead$summary.hy[2,]
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
	8118.967	5144.412	2310.843	6822.871	21483.031	4923.374

And, to  $\log(\kappa)$  we have

```
r.sead$summary.hy[4,]
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
	1.3795478	0.3085711	0.7624957	1.3843872	1.9743398	1.3977385

The variance and range of the spatial process we need a post processing. We obtain the marginal distribution for both by

```
r.f <- inla.spde2.result(r.sead, 'i', spde, do.transf=TRUE)
```

The posterior mean for the marginal variance

```
inla.emarginal(function(x) x, r.f$marginals.variance.nominal[[1]])
```

```
[1] 0.06108992
```

And to the practical range of the process we have

```
inla.emarginal(function(x) x, r.f$marginals.range.nominal[[1]])
```

```
[1] 0.7460206
```

At the Figure 4.3 we look the posterior distribution to  $\beta_0$ ,  $\phi$ ,  $1/\kappa$ ,  $\sigma_x^2$ , practical range and to the mean and 95% credibility interval of the distance to sea effect at Figure 4.3. We choose  $1/\kappa$  instead  $kappa$  because  $1/\kappa$  is the range parameter and in this case is expressed in degrees units.

```
par(mfrow=c(2,3), mar=c(3,3.5,0,0), mgp=c(1.5, .5, 0), las=0)
plot(r.sead$marginals.fix[[1]], type='l', xlab='Intercept', ylab='Density')
plot(r.sead$summary.random[[1]][,1:2], type='l',
      xlab='Distance to sea (Km)', ylab='Coefficient'); abline(h=0, lty=3)
for (i in c(4,6)) lines(r.sead$summary.random[[1]][,c(1,i)], lty=2)
plot(r.sead$marginals.hy[[1]], type='l', ylab='Density', xlab=expression(phi))
plot.default(inla.tmarginal(function(x) 1/exp(x), r.sead$marginals.hy[[4]]),
              type='l', xlab=expression(kappa), ylab='Density')
plot.default(r.f$marginals.variance.nominal[[1]], type='l',
              xlab=expression(sigma[x]^2), ylab='Density')
plot.default(r.f$marginals.range.nominal[[1]], type='l',
              xlab='Practical range', ylab='Density')
```

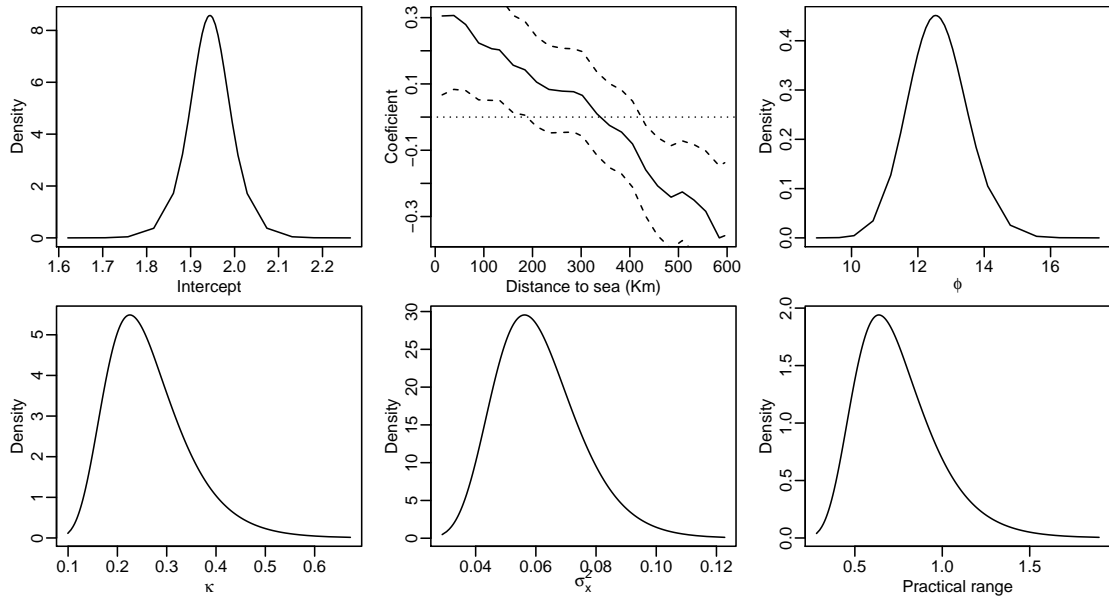


Figure 4.3: Posterior marginal distribution of  $\beta_0$  (top left), the posterior mean (continuous line) and 95% credibility interval (dashed lines) to effect of distance to sea (top mid), posterior marginal distribution to  $\phi$  (top right), posterior marginal distribution of  $\kappa$  (bottom left), posterior marginal distribution to nominal variance of the random field (bottom mid) and posterior marginal distribution of the practical range (bottom right).

### 4.3 Testing the significance of spatial effect

Now, we want to test the significance of the spatial random effect component on the model. To access the significance of this effect, can be use the DIC measure, comparing the model with this component with the DIC of the model without this component, fitted below

```
r0.sead <- inla(y ~ 0 + Intercept + f(gSeaDist, model='rw1'),
               family='Gamma', control.compute=list(dic=TRUE),
               data=inla.stack.data(stk.dat),
               control.predictor=list(A=inla.stack.A(stk.dat), compute=TRUE))
```

So, we have the DIC of both models with

```
c(r0.sead$dic$dic, r.sead$dic$dic)
```

```
[1] 2720.365 2565.174
```

and conclude that the spatial effect is significantly to explain the precipitation.

### 4.4 Prediction of the random field

To show the spatial effect, we make prediction of this effect on a fine grid. The `inla.mesh.projector()` get a projector matrix to a regular grid and with equal dimension in each direction by default. The Paraná state shape is more width than height. So, we choose diferent shape of grid to get an adequate cell size. We use

```
(stepsize <- 4*1/111)
```

```
[1] 0.03603604
```

```
(nxy <- round(c(diff(range(PRborder[,1])), diff(range(PRborder[,2])))/stepsize))
```

[1] 183 117

Here we divide the range of each coordinate by 4 kilometers on degree scale. We use (4/111) factor because each degree has approximately 111 kilometers. Its because we have the coordinates with units in kilometers and we want each cell on the grid with  $4 \times 4$  km.

Now, we get the projector with

```
projgrid <- inla.mesh.projector(prmesh2, xlim=range(PRborder[,1]),  
                               ylim=range(PRborder[,2]), dims=nxy)
```

and get the posterior mean and posterior standard deviation with

```
xmean <- inla.mesh.project(projgrid, r.sead$summary.random$i$mean)  
xsd <- inla.mesh.project(projgrid, r.sead$summary.random$i$sd)
```

To good visualization, we make NA the values corresponding of the points of the grid out of the border of the Paraná state. To do it, we use the function `inout()` on **splancs** package, with

```
require(splancs)  
table(xy.in <- inout(projgrid$lattice$loc,  
                    cbind(PRborder[,1], PRborder[,2])))
```

```
FALSE TRUE  
7865 13546
```

```
xmean[!xy.in] <- xsd[!xy.in] <- NA
```

We visualize the this on Figure 4.4. We see on top left graph at Figure 4.4 that the random field has variation from -0.6 to 0.4. It isn't despicable, because we have, on the top right graph, the standard errors around 0.2. The variation on the mean of this random effect express the remain variation after the consideration of the effect of the covariate on the model. The variation on its standard deviation is due to density of stations over the region. For example, on the top right graph we see, the first blue region from right to left is near Curitiba and around there many stations.

## 4.5 Prediction of the response on a grid

We want to predict the response on a grid. A naive approach is made by the projection of the sum of the linear predictor components and applying exponentiation to it, because we use the log link. A better alternative is the joint prediction with the estimation process. But it is computationally expensive when we have large grid.

Using the grid from previous section we have points discarded because they are not within the Paraná border. To make the predictions only on the points within the Paraná state we select the coordinates of the grid and the correspondent rows of the predictor matrix by

```
prdcoo <- projgrid$lattice$loc[which(xy.in),]  
Aprd <- projgrid$proj$A[which(xy.in), ]
```

Now we get the covariate values at each point on the grid. First we compute the distance to sea

```
seaDist0 <- apply(spDists(PRborder[1034:1078,], prdcoo, longlat=TRUE), 2, min)
```

and found the knot group, using the same used on the model

```

ug.seadist <- sort(unique(inla.group(PRprec$seaDist)))
ig0 <- apply(abs(outer(ug.seadist, seaDist0, '-')),2,which.min)
g.seadist <- ug.seadist[ig0]

```

and put it into on a stack to prediction and join with the stack data

```

stk.prd <- inla.stack(data=list(y=NA), A=list(Aprd,1,1),
                     effects=list(i=1:spde$n.spde,
                                   Intercept=rep(1,nrow(Aprd)),
                                   gSeaDist=g.seadist), tag='prd')
stk.all <- inla.stack(stk.dat, stk.prd)

```

Now, we fit the model joining the prediction data together the observed data stack. But, now we don't need the computation of DIC, marginal distributions and quantiles, just the mean and standard deviation of the response

```

r2.sead <- inla(f.sead, family='Gamma', data=inla.stack.data(stk.all),
               control.predictor=list(A=inla.stack.A(stk.all),
                                       compute=TRUE), quantiles=NULL,
               control.results=list(return.marginals.random=FALSE,
                                    return.marginals.predictor=FALSE))

```

Now, we extract the index on the predictor that corresponds to these wanted. Also, we put the mean and standard deviation on the matrix format, in the adequate positions.

```

id.prd <- inla.stack.index(stk.all, 'prd')$data
sd.prd <- m.prd <- matrix(NA, nxy[1], nxy[2])
m.prd[xy.in] <- r2.sead$summary.fitted.values$mean[id.prd]
sd.prd[xy.in] <- r2.sead$summary.fitted.values$sd[id.prd]

```

and we visualize on the Figure 4.4 the posterior mean and standard deviation to response with commands bellow on

```

grid.arrange(levelplot(xmean, col.regions=topo.colors(99),
                       xlab='', ylab='', scales=list(draw=FALSE)),
              levelplot(xsd, col.regions=topo.colors(99),
                       xlab='', ylab='', scales=list(draw=FALSE)),
              levelplot(m.prd, col.regions=topo.colors(99),
                       xlab='', ylab='', scales=list(draw=FALSE)),
              levelplot(sd.prd, col.regions=topo.colors(99),
                       xlab='', ylab='', scales=list(draw=FALSE)))

```

We see at this Figure a clear pattern on the average of the daily accumulated precipitation occurred on January of 2012. It's high near to sea and lower at north west side of the Paraná state.

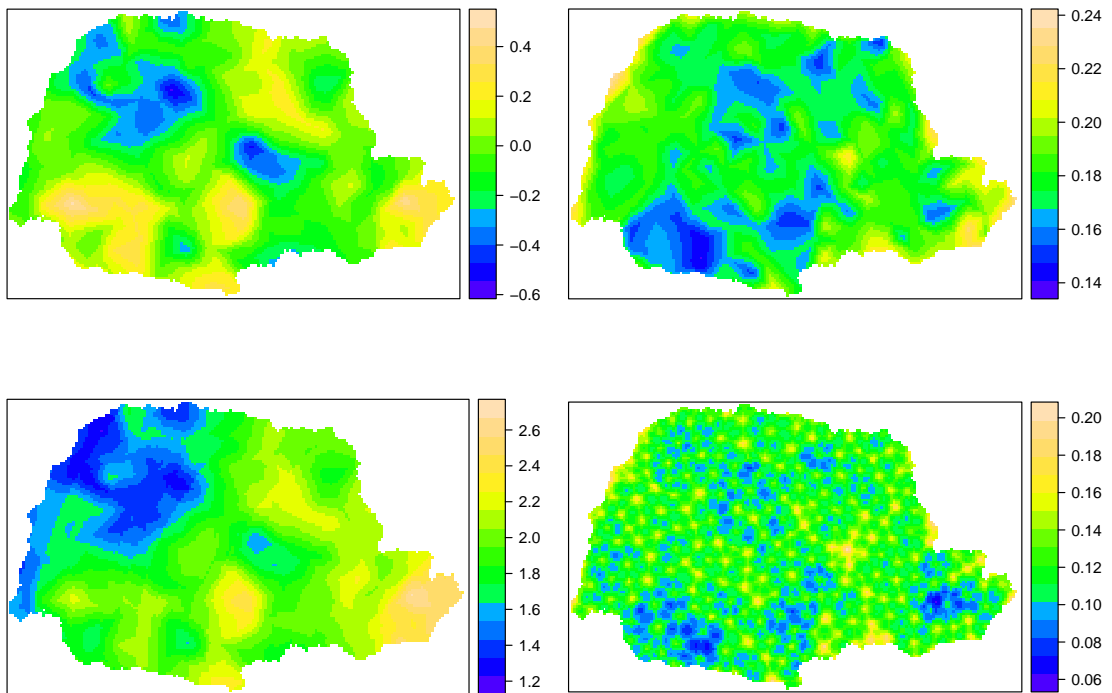


Figure 4.4: Posterior mean and standard deviation of the random field (top left and right respectively). Posterior mean and standard deviation of the response (top left and right respectively).



## Chapter 5

# Semicontinuous model to daily rainfall

The daily rainfall data is one type of data that is not easy to fit a simple model. Because there are not days when it rains and it rains in the days, the distribution of rainfall is very asymmetric. In this section we build a jointly model to rainfall taking into account this fact. We build a semicontinuous model that is a jointly model to occurrence and to positive values.

### 5.1 The data and the model

The data is called by

```
data(PRprec)
PRprec[1:3, 1:10]
```

	Longitude	Latitude	Altitude	d0101	d0102	d0103	d0104	d0105	d0106	d0107
1	-50.8744	-22.8511	365	0	0	0	0.0	0	0	2.5
3	-50.7711	-22.9597	344	0	1	0	0.0	0	0	6.0
4	-50.6497	-22.9500	904	0	0	0	3.3	0	0	5.1

and we have, for the first five days on January, the summary and the number of station with rain

```
sapply(PRprec[,4:9], summary)
```

	d0101	d0102	d0103	d0104	d0105	d0106
Min.	0.0000	0.00	0.000	0.000	0.000	0.000
1st Qu.	0.0000	0.00	0.000	0.000	0.000	0.000
Median	0.0000	0.00	0.300	0.000	0.000	0.000
Mean	0.8074	2.87	5.928	3.567	2.793	4.059
3rd Qu.	0.0000	2.10	7.000	3.475	0.600	3.850
Max.	43.8000	55.60	78.500	62.900	64.100	81.400
NA's	6.0000	7.00	7.000	6.000	6.000	8.000

```
apply(PRprec[,4:9]>0, 2, summary)
```

	d0101	d0102	d0103	d0104	d0105	d0106
Mode	"logical"	"logical"	"logical"	"logical"	"logical"	"logical"
FALSE	"538"	"399"	"294"	"373"	"429"	"368"
TRUE	"72"	"210"	"315"	"237"	"181"	"240"
NA's	"6"	"7"	"7"	"6"	"6"	"8"

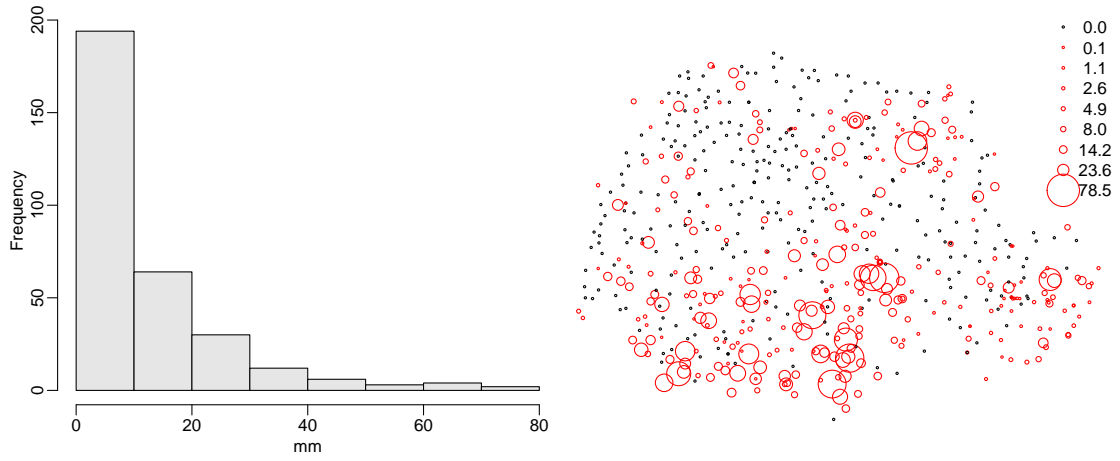


Figure 5.1: Histogram of positive precipitation and map of precipitation.

and we see that we have no rain in 538 stations on the first day, 399 on second and 294 on third. We also have missing values.

Let  $r_i$  the rain on the location  $s_i$ . We use the augmented data by definition of two new variables. We define the occurrence variable

$$z_i = \begin{cases} 1, & \text{if } r_i > 0 \\ 0, & \text{otherwise} \end{cases}$$

and the amount variable

$$y_i = \begin{cases} NA, & \text{if } r_i = 0 \\ r_i, & \text{otherwise} \end{cases}$$

Using the data on the third day we have

```
jday <- 6
z <- (PRprec[,jday]>0)+0
y <- ifelse(PRprec[,jday]>0, PRprec[,jday], NA)
```

The map of this data are showed on the left graph of the Figure 5.1. Also, the histogram of the positive precipitation are on the right graph of this Figure. Both produced by the code below

```
par(mfrow=c(1, 2), mar=c(3,3,.5,0), mgp=c(1.5,0.5,0))
hist(y, xlab='mm', col=gray(.9), main='')
par(mar=c(0,0,0,0))
plot(PRprec[,1:2], cex=0.3 + PRprec[,jday]/20, asp=1, col=z+1, axes=FALSE)
q.y <- c(0, quantile(y, 0:7/7, na.rm=T))
legend('topright', format(q.y, dig=2), pch=1, bty='n',
      pt.cex=0.3+q.y/20, col=c(1,rep(2,length(q.y)-1)))
```

We use a model with two likelihoods. We use the Bernoulli likelihood for  $z_i$  and the gamma for  $y_i$ .

$$z_i \sim \text{Bernoulli}(p_i) \quad y_i \sim \text{Gamma}(a_i, b_i)$$

We define the linear predictor to first component by

$$\text{logit}(p_i) = \alpha_z + x_i$$

where  $\alpha_z$  is an intercept and  $x_i$  is a random effect modeled by a Gaussian field through the SPDE approach.

To second component we consider that  $E(y_i) = \mu_i = a_i/b_i$  and  $V(y_i) = a_i/b_i^2 = \mu_i^2/\phi$ , where  $\phi$  is a precision parameter and we define a linear predictor to  $\log(\mu_i)$ . We have

$$\log(\mu_i) = \alpha_y + \beta x_i$$

where  $\alpha_y$  is an intercept and  $\beta$  is a scaling parameter to  $x_i$  that is a shared random effect with the first component of the model.

## 5.2 Fitting the model and some results

To define a SPDE model we need a mesh.

We use the mesh builded for the Paraná state on the Chapter 2.2. The SPDE model definition is made by

```
spde <- inla.spde2.matern(prmesh2, alpha=2)
```

We have to prepare the data in addequate form using the stack functionality. First we get the predictor matrix for the SPDE model effect with

```
A <- inla.spde.make.A(prmesh2, loc=as.matrix(PRprec[,1:2]))
```

Taking into account that we need to model two responses together, the response is a two column matrix. The data stack for the rain occurence is

```
stk.z <- inla.stack(tag='est.z',
  data=list(y=cbind(z, NA)), ### first column
  A=list(A, 1),
  effects=list(
    list(i.z=1:spde$n.spde),
    list(z.b0=rep(1,length(z))))))
```

The data stack for the precipitation amount is

```
stk.y <- inla.stack(tag='est.y',
  data=list(y=cbind(NA, y)), ### first column
  A=list(A, 1),
  effects=list(
    list(i.y=1:spde$n.spde),
    list(y.b0=rep(1,length(y))))))
```

The full data stack is

```
stk.zy <- inla.stack(stk.z, stk.y)
```

To fit the model, we make inference to  $\beta$  using the copy of the strategy with `fixed=FALSE`.

```
res <- inla(y ~ 0 + z.b0 + y.b0 +
  f(i.z, model=spde) + f(i.y, copy='i.z', fixed=FALSE),
  family=c('binomial', 'gamma'),
  data=inla.stack.data(stk.zy), control.compute=list(dic=TRUE),
  control.predictor=list(A=inla.stack.A(stk.zy), compute=TRUE))
```

Now, we look at the summary of the posterior distribution to  $\alpha_z$  and  $\alpha_y$

```
round(res$summary.fix, 4)
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode	kld
z.b0	0.0575	0.5175	-0.9887	0.0518	1.1352	0.0425	0
y.b0	2.1397	0.1924	1.7479	2.1426	2.5192	2.1503	0

We apply the inverse of the link function to look at response scale with

```
c(binomial(link='logit')$linkinv(res$summary.fix[1,1]),
  exp(res$summary.fix[2,1]))
```

```
[1] 0.5143776 8.4972557
```

and it looks similar to observed:

```
c(occ=mean(z, na.rm=TRUE), rain=mean(y, na.rm=TRUE))

      occ      rain
0.5172414 11.4606349
```

The summary of the posterior distribution of the precision parameter on second component is

```
res$summary.hy[1,]

      mean      sd 0.025quant    0.5quant 0.975quant      mode
0.86447598 0.06385172 0.74508410 0.86240797 0.99603502 0.85850475
```

And the summary of the posterior distributio of the scaling parameter  $\beta$  is

```
res$summary.hy[4,]

      mean      sd 0.025quant    0.5quant 0.975quant      mode
0.32702426 0.07468756 0.18127259 0.32673579 0.47437267 0.32602797
```

and we see on this summary that the random effect  $x_i$  is shared by  $z$  and  $y$ . In other words, the probability of rain occurrence is correlated with the amount of the rain, like the preferential sampling problem [Diggle et al., 2010].

We look the posterior marginal distributions of  $\alpha_z$ ,  $\alpha_y$ ,  $\phi$  and  $\beta$  on Figure 5.2 with commands below

```
par(mfrow=c(2,2), mar=c(3,3,.5,.5), mgp=c(1.5,.5,0), las=1)
plot(res$marginals.fix[[1]], type='l', ylab='Density',
      xlab=expression(alpha[z]))
plot(res$marginals.fix[[2]], type='l', ylab='Density',
      xlab=expression(alpha[y]))
plot(res$marginals.hy[[1]], type='l', ylab='Density',
      xlab=expression(phi))
plot(res$marginals.hy[[4]], type='l', ylab='Density',
      xlab=expression(beta))
```

To decide if the occurrence rain and the amount of the rain have spatial dependency, we want to test the significance of the random effect  $x$ .

One approach is looking at the 2,5% and 97,5% quantiles of each element of the random effect. If for some of them they have both quantiles with same signal, then they are significantly not null. We see on the Figure 5.3 the mean and these quantiles for each  $x$  in order of the posterior mean, with code below

```
ordx <- order(res$summary.random$i.z$mean)
par(mar=c(3,3,0.5,0.5), mgp=c(1.5,0.5,0), las=1)
plot(res$summary.random$i.z$mean[ordx], type='l', ylab='x',
      ylim=range(res$summary.random$i.z[, 4:6]))
for (i in c(4,6)) lines(res$summary.random$i.z[ordx,i], lty=2)
abline(h=0, lty=3)
```

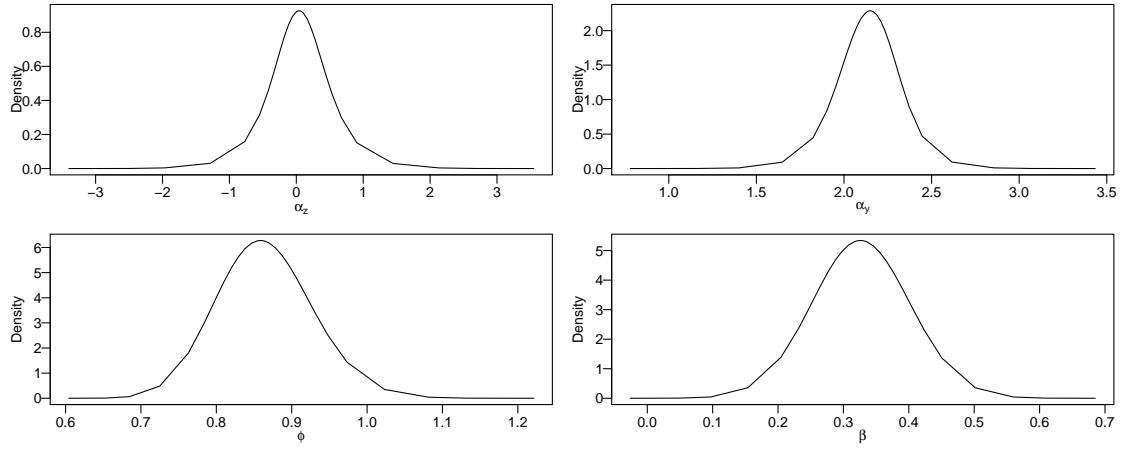


Figure 5.2: Posterior marginal distributions of  $\alpha_z$  (top left),  $\alpha_y$  (top right),  $\phi$  (bottom left) and  $\beta$  (bottom right).

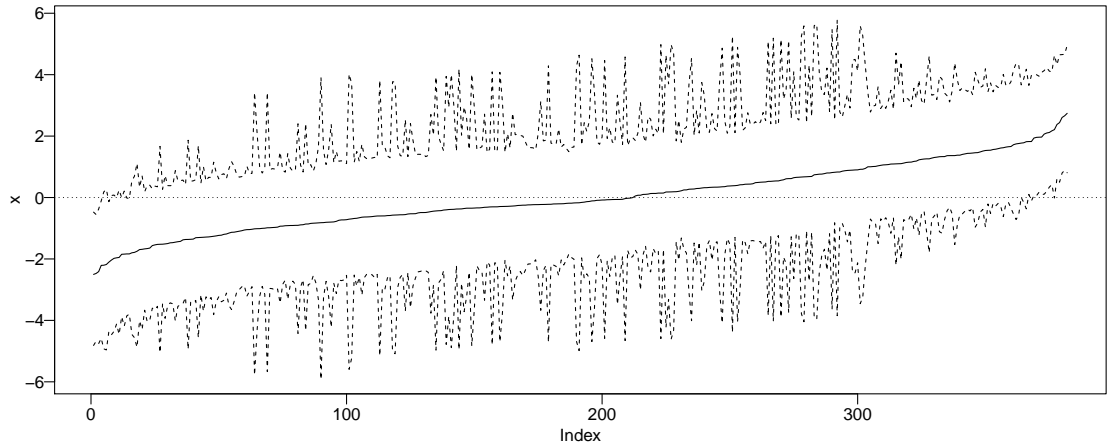


Figure 5.3: Posterior mean of each  $x_i$  (continuous line) and its 2.5% and 97.5% quantiles of the posterior marginal distributions (dashed lines).

Other approach is to compare the DIC of the model with  $x$  with the DIC of the model without it.

We fit the model without  $x$  with

```
res0 <- inla(y ~ 0 + z.b0 + y.b0,
            family=c('binomial', 'gamma'),
            data=inla.stack.data(stk.zy), control.compute=list(dic=TRUE),
            control.predictor=list(A=inla.stack.A(stk.zy), compute=TRUE))
```

and we have

```
c(res0$dic$dic, res$dic$dic)
```

```
[1] 3002.937 2840.135
```

and conclude that the spatial random effect is significative to occurence and to amount of the rain.

### 5.3 Results for the spatial random effect

The spatial random effect modeled by the SPDE approach has the  $\kappa$  scaling parameter, the marginal variance parameter  $\sigma_x^2$  and the practical range, [Lindgren et al., 2011].

The summary of the posterior distribution of  $\log(\kappa)$  is

```
res$summary.hy[3,]
```

```
      mean      sd 0.025quant    0.5quant 0.975quant      mode
0.74520125 0.35318467 0.04289549 0.74927226 1.43043055 0.76371444
```

To obtain results for the variance and for the practical range we need a post processing. We obtain the marginal distribution for both by

```
res.f <- inla.spde2.result(res, 'i.z', spde, do.transf=TRUE)
```

The posterior mean for the marginal variance

```
inla.emarginal(function(x) x, res.f$marginals.variance.nominal[[1]])
```

```
[1] 2.318503
```

And to the practical range of the process we have

```
inla.emarginal(function(x) x, res.f$marginals.range.nominal[[1]])
```

```
[1] 1.427326
```

At the Figure 5.4 we look the posterior distribution to  $\kappa$ ,  $\sigma_x^2$  and of the practical range. The  $\kappa$  and practical range are on degrees units.

```
par(mfrow=c(1,3), mar=c(3,3.5,0,0), mgp=c(1.5, .5, 0), las=0)
plot.default(inla.tmarginal(function(x) exp(x), res$marginals.hy[[3]]),
             type='l', xlab=expression(kappa), ylab='Density')
plot.default(res.f$marginals.variance.nominal[[1]], type='l',
             xlab=expression(sigma[x]^2), ylab='Density')
plot.default(res.f$marginals.range.nominal[[1]], type='l',
             xlab='Practical range', ylab='Density')
```

Another interesting result is the map of the random field. To get it we use get the projector with

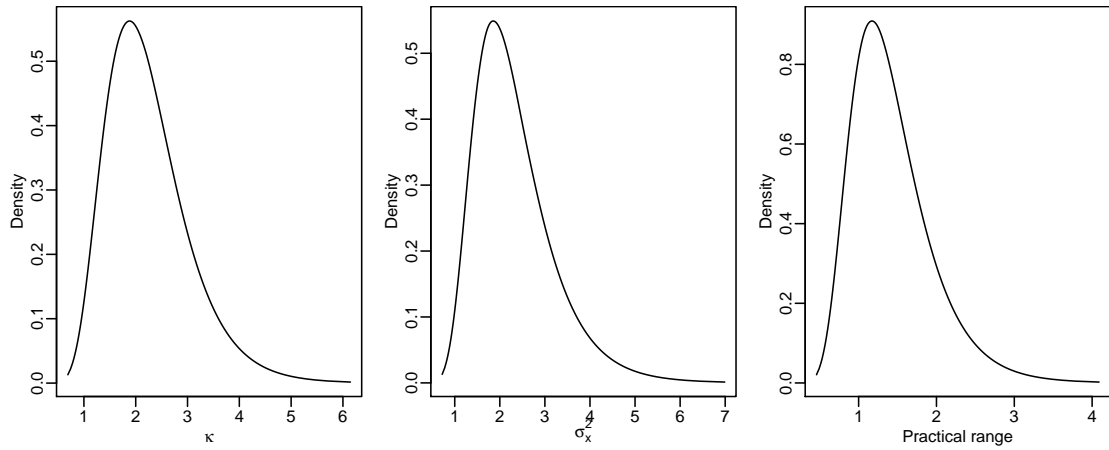


Figure 5.4: Posterior marginal distribution of  $\kappa$  (left),  $\sigma_x^2$  (mid) and of the practical range (right).

```
data(PRborder)
(nxy <- round(c(diff(range(PRborder[,1])), diff(range(PRborder[,2])))/.02))

[1] 330 210

projgrid <- inla.mesh.projector(prmesh2, xlim=range(PRborder[,1]),
                                ylim=range(PRborder[,2]), dims=nxy)
```

and get the posterior mean and posterior standard deviation with

```
xmean <- inla.mesh.project(projgrid, res$summary.random$i.z$mean)
xsd <- inla.mesh.project(projgrid, res$summary.random$i.z$sd)
```

To good visualization, we make NA the values corresponding of the points of the grid out of the border of the Paraná state. To do it, we use the function `inout()` on **splancs** package, with

```
table(xy.in <- inout(projgrid$lattice$loc, PRborder))

FALSE TRUE
25188 44112

xmean[!xy.in] <- xsd[!xy.in] <- NA
```

We visualize the this on Figure 5.5. with code below

```
grid.arrange(levelplot(xmean, col.regions=topo.colors(99),
                        xlab='', ylab='', scales=list(draw=FALSE)),
              levelplot(xsd, col.regions=topo.colors(99),
                        xlab='', ylab='', scales=list(draw=FALSE)), nrow=1)
```

Comparing the left graph on Figure 5.5 with the right graph on the Figure 5.1 we have a more clear pattern of the rain on third day of 2012 on Paraná state. We see that in fact we have rain near the sea (the right border of the map), at south and a small point on northeast. Also we have no rain in some parts of nothwest.

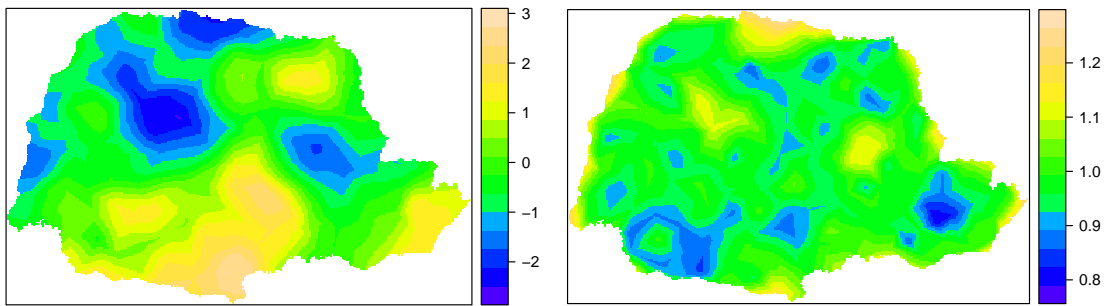


Figure 5.5: Posterior mean (left) and standard deviation (right) of the random field.



## Chapter 6

# Joint modeling a covariate with misalignment

Here we focus on the situation when we have a response  $y$  and a covariate  $c$ . But we have misalignment, i.e., we have  $y$  observed at  $n_y$  locations and  $x$  observed at  $n_c$  locations. In this example, we design a solution that not depends if we have or not some common observed locations for  $c$  and  $y$ .

A restriction is the assumption that  $c$  have spatial dependency. This restriction is made because we want a good model to predict  $c$  at locations of  $y$ . So, the goodness of prediction is proportional to the spatial dependency.

### 6.1 The model

Taking into account that  $c$  have spatial dependency, a simple approach is to define a model for  $c$ , predict it on locations that we have  $y$  and build the model for  $y$ . But, in this two stage model, we don't take into account the prediction error of  $c$  on second model. The measurement error models is an approach to solve similar problems, [Muff et al., 2013]. But, here we are able to consider the spatial dependency on  $c$ . So we build a spatial model for  $c$  and another spatial model for  $y$ , and do the estimation process jointly.

We consider the following likelihood for  $c$

$$c_i \sim N(m_i, \sigma_c^2)$$

where  $m_i$  is modeled as a random field and  $\sigma_c$  is a measurement error of  $c$ . So, when we predict  $c$  at location  $s$ , we have  $m(s)$ .

Let following model for  $y$

$$y_i \sim N(\alpha_y + \beta c_i + x_i, \sigma_y^2)$$

where  $\alpha_y$  is an intercept,  $\beta$  is the regression coefficient on  $c$ ,  $c_i$  is the covariate at location of  $y_i$ ,  $x_i$  is an zero mean random field and  $\sigma_y^2$  measures the error that remain unexplained on  $Y$ .

A particular case is when we don't have the  $x$  term in the model for  $y$ . Another case, is when  $\sigma_c^2 = 0$  and we don't have white noise in the covariate, i. e., the covariate is considered just a realization of a random field.

#### 6.1.1 Simulation from the model

To test the approach on next section, we do a simulation from this model. First, we set the model parameters.

```
n.y <- 123;      n.c <- 234
alpha.y <- -5;   beta <- -3
sigma2.y <- 0.5; sigma2.c <- 0.2
```

First, we do the simulation of the locations

```
set.seed(1)
loc.c <- cbind(runif(n.c), runif(n.c))
loc.y <- cbind(runif(n.y), runif(n.y))
```

Let the parameters of both random fields  $m$  and  $x$ :

```
kappa.m <- 3;          sigma2.m <- 5
kappa.x <- 7;          sigma2.x <- 3
```

and we use the `grf()` function of the **geoR**, [Ribeiro Jr and Diggle, 2001], package to do the simulation of both random fields. We need the simulation of  $m$  in both set of locations

```
library(geoR)
set.seed(2)
mm <- grf(grid=rbind(loc.c, loc.y), messages=FALSE,
           cov.pars=c(sigma2.m, 1/kappa.m), kappa=1)$data
xx <- grf(grid=loc.y, messages=FALSE,
           cov.pars=c(sigma2.x, 1/kappa.x), kappa=1)$data
```

and simulate the covariate and the response with

```
set.seed(3)
cc <- mm + rnorm(n.c+n.y, 0, sqrt(sigma2.c))
yy <- alpha.y + beta*mm[n.c+1:n.y] + xx +
     rnorm(n.y, 0, sqrt(sigma2.y))
```

## 6.2 Fitting the model

First we make the mesh

```
p101 <- matrix(c(0,1,1,0,0, 0,0,1,1,0), ncol=2)
(mesh <- inla.mesh.create.helper(, p101, cutoff=.03,
                                max.edge=c(.05, .1)))$n
```

```
[1] 1334
```

and use it for both the random fields. So, the both index set based on the mesh have the same values.

We do simulations of the covariate on the locations of the response just to simulate the response. But, in the problem that we want to solve in practice, we don't have the covariate on the response locations. The misalignment implies in different predictor matrix for response and covariate.

```
Ac <- inla.spde.make.A(mesh, loc=loc.c)
Ay <- inla.spde.make.A(mesh, loc=loc.y)
```

To predict the covariate jointly, we need to model it jointly and we need two likelihoods. So, the response is formed by two columns matrix. The data stack can be obtained by

```
stk.c <- inla.stack(data=list(y=cbind(cc[1:n.c], NA)),
                   A=list(Ac), tag='dat.c',
                   effects=list(m=1:mesh$n))
stk.y <- inla.stack(data=list(y=cbind(NA, yy)),
                   A=list(Ay, 1),
                   effects=list(list(c.y=1:mesh$n, x=1:mesh$n),
                                   list(a.y=rep(1,length(yy))))))
stk <- inla.stack(stk.c, stk.y)
```

The estimation of the regression coefficient in this approach is treated as a hyperparameter, such as copy parameter of an latent field. In this case, we need to do a good prior specification. For example, is possible to know, a priori, the signal. We set a  $N(-3, 25)$  prior to  $\beta$ . Also, we define the formulae for the model.

```
form <- y ~ 0 + f(m, model=spde) +
  a.y + f(x, model=spde) + f(c.y, copy='m', fixed=FALSE,
    hyper=list(theta=list(param=c(-3,25), initial=-3)))
```

define the spde model and fit the model with

```
spde <- inla.spde2.matern(mesh)
res <- inla(form, data=inla.stack.data(stk), family=rep('gaussian',2),
  control.predictor=list(compute=TRUE, A=inla.stack.A(stk)))
```

### 6.3 The results

The true values of the intercept and the summary of its posterior marginal

```
round(c(True=alpha.y, res$summary.fix), 4)
```

```
True
-5.0000 -4.5118  1.3817 -7.4194 -4.5176 -1.5663 -4.5254  0.0000
```

The true values of the precisions and the summary of the posterior marginals

```
round(cbind(True=1/c(Prec.c=sigma2.c, Prec.y=sigma2.y),
  res$summary.hy[1:2,]), 4)
```

	True	mean	sd	0.025quant	0.5quant	0.975quant	mode
Prec.c	5	5.7439	0.7719	4.3671	5.6965	7.3970	5.6051
Prec.y	2	1.4997	0.3844	0.8893	1.4508	2.3898	1.3581

The true value of the regression coefficient and the summary of the posterior marginal

```
round(c(True=beta, res$summary.hy[7,]), 4)
```

	True	mean	sd	0.025quant	0.5quant	0.975quant
	-3.0000	-2.7658	0.1585	-3.0754	-2.7662	-2.4533
mode						
	-2.7676					

The true values for the precision of the both random fields and the summary of the posterior marginals

```
m.rf <- inla.spde2.result(res, 'm', spde)
x.rf <- inla.spde2.result(res, 'x', spde)
round(cbind(True=c(s2m=sigma2.m, s2x=sigma2.x),
  mean=c(inla.emarginal(function(x) x,
    m.rf$marginals.variance.n[[1]]),
    inla.emarginal(function(x) x,
      x.rf$marginals.variance.n[[1]])),
  rbind(inla.hpdmarginal(.95, m.rf$marginals.variance.n[[1]]),
    inla.hpdmarginal(.95, x.rf$marginals.variance.n[[1]]))), 4)
```

	True	mean	low	high
s2m	5	4.4351	1.4046	8.4116
s2x	3	4.9245	1.2045	9.9349

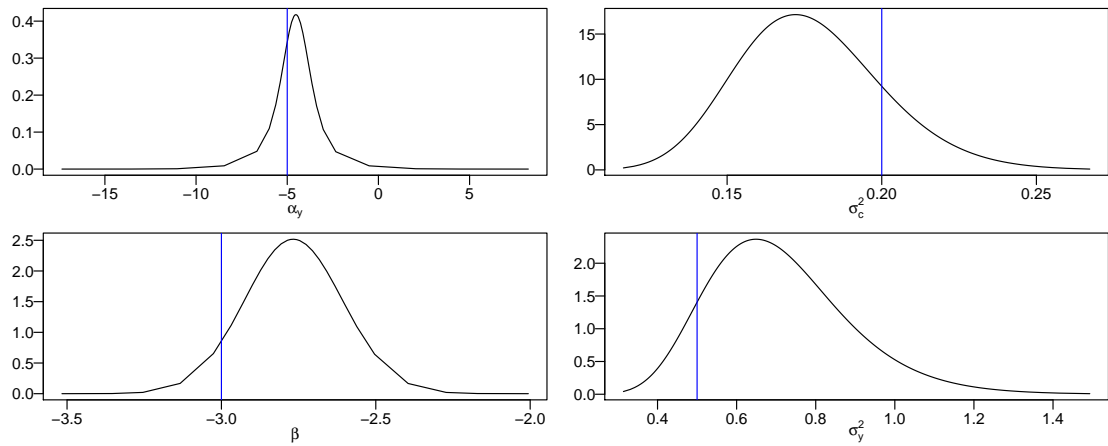


Figure 6.1: Posterior distribution of the likelihood parameters.

The true values for the scale parameter  $\kappa$  the mean of the posterior marginals and the 95% HPD credibility interval.

```
post.k <- lapply(res$marginals.hy[c(4,6)], function(y)
  inla.tmarginal(function(x) exp(x), y))
round(cbind(True=c(kappa.m=kappa.m, kappa.x=kappa.x),
  t(sapply(post.k, function(x)
    c(mean=inla.emarginal(function(y) y, x),
      CI=inla.hpdmarginal(.95, x))))), 4)
```

	True	mean	CI1	CI2
kappa.m	3	3.7800	2.0607	5.6362
kappa.x	7	5.9067	2.2330	10.3654

We see the posterior distribution of likelihood parameters on Figure 6.1 generated with comandns below

```
par(mfcol=c(2,2), mar=c(3,3,.1,.1), mgp=c(1.5,.5,0), las=1)
plot(res$marginals.fix[[1]], type='l',
  xlab=expression(alpha[y]), ylab='')
abline(v=alpha.y, col=4)
plot(res$marginals.hy[[7]], type='l',
  xlab=expression(beta), ylab='')
abline(v=beta, col=4)
plot.default(inla.tmarginal(function(x) 1/x, res$marginals.hy[[1]]),
  type='l', xlab=expression(sigma[c]^2), ylab='')
abline(v=sigma2.c, col=4)
plot.default(inla.tmarginal(function(x) 1/x, res$marginals.hy[[2]]),
  type='l', xlab=expression(sigma[y]^2), ylab='')
abline(v=sigma2.y, col=4)
```

We see on the Figure 6.1 that the posterior distribution covers the true values of all the parameters.

Now we want to see the posterior distribution of the both random fields. The practical range of the process is  $\sqrt{8\nu}/\kappa$ , where  $\nu = 1$  on this analisys. We see these parameters on Figure 6.2 generated with comandns below

```
par(mfcol=c(2,3), mar=c(3,3,.1,.3), mgp=c(1.5,.5,0), las=1)
plot.default(post.k[[1]], type='l', xlab=expression(kappa[m]), ylab='')
abline(v=kappa.m, col=4)
```

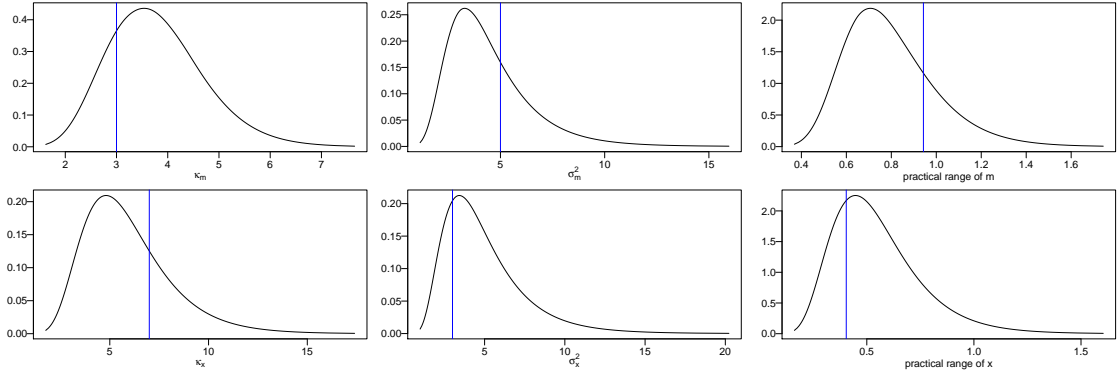


Figure 6.2: Posterior distribution of all the parameters of both random field.

```
plot.default(post.k[[2]], type='l', xlab=expression(kappa[x]), ylab='')
abline(v=kappa.x, col=4)
plot.default(m.rf$marginals.variance.n[[1]], type='l',
             xlab=expression(sigma[m]^2), ylab='')
abline(v=sigma2.m, col=4)
plot.default(x.rf$marginals.variance.n[[1]], type='l',
             xlab=expression(sigma[x]^2), ylab='')
abline(v=sigma2.x, col=4)
plot.default(m.rf$marginals.range.n[[1]], type='l',
             xlab='practical range of m', ylab='')
abline(v=sqrt(8)/kappa.m, col=4)
plot.default(x.rf$marginals.range.n[[1]], type='l',
             xlab='practical range of x', ylab='')
abline(v=sqrt(8)/kappa.x, col=4)
```

We see on Figure 6.2 that the posterior marginal distribution of the all parameters of both spatial process cover the true values well.

Another interesting result is the prediction of the covariate on the locations of the response. We have the simulated values of  $m$  on that locations. So, we are able to see if the predictions are good.

The predictor matrix used on the estimation proces maps the nodes from mesh vertices to the data locations. The first lines of the predictor matrix for the covariate can be used to access the predictions on the locations of the covariate. Also, we have the predictor matrix used to the response. The last lines of this matrix that maps the mesh vertices to the response locations. Because we have the covariate simulated in the both set of locations, we use the correspondent parts of both predictor matrix to project the posterior mean and the posterior variance on the locations.

We get this matrix by

```
mesh2locs <- rBind(Ac, Ay)
```

and the posterior mean and posterior standard deviations with

```
m.mprd <- drop(mesh2locs%%res$summary.ran$m$mean)
sd.mprd <- drop(mesh2locs%%res$summary.ran$m$sd)
```

With this aproach for this both posterior summary can be an aproximation to 95% credibility interval, with normally supposition. We see it this results with comand below

```
plot(m.mprd, mm, asp=1, type='n',
     xlab='Predicted', ylab='Simulated')
segments(m.mprd-2*sd.mprd, mm, m.mprd+2*sd.mprd, mm,
```

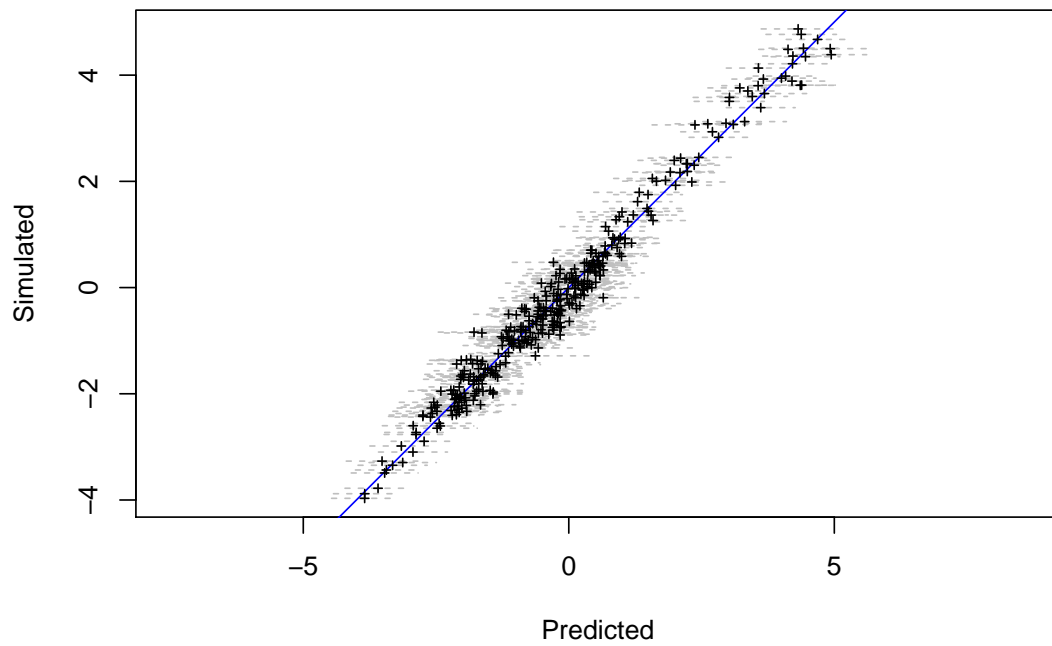


Figure 6.3: Simulated versus predicted values of  $m$  (+) and the approximated credibility intervals.

```
lty=2, col=gray(.75))
abline(c(0,1), col=4); points(m.mprd, mm, pch=3, cex=.5)
```

on the Figure 6.3. The blue line represents the situation where predicted is equal to simulated. We see that the prediction are very well.

# Chapter 7

## Non stationary model

### 7.1 Introduction

To introduce the non stationarity we need to remember the definition of the precision matrix of the GMRF that defines the RF. This matrix is defined on equations (2.7) and (2.8). The non stationarity is made by a redefinition of the precision matrix.

This new definition is

$$\mathbf{Q} = \mathbf{D}^{(0)}(\mathbf{D}^{(1)}\mathbf{M}^{(0)}\mathbf{D}^{(1)} + \mathbf{D}^{(2)}\mathbf{D}^{(1)}\mathbf{M}^{(1)} + (\mathbf{M}^{(1)})^T\mathbf{D}^{(1)}\mathbf{D}^{(2)} + \mathbf{M}^{(2)})\mathbf{D}^{(0)} \quad (7.1)$$

where  $\mathbf{M}^{(0)}$ ,  $\mathbf{M}^{(1)}$  and  $\mathbf{M}^{(2)}$ , are provided from the finite element method - FEM based on the mesh. For  $\alpha = 1$  ( $\nu = 0$ ), we have  $\mathbf{M}^{(0)} = \mathbf{C}$ ,  $(\mathbf{M}^{(1)})_{ij} = 0$  and  $\mathbf{M}^{(2)} = \mathbf{G}$ . For  $\alpha = 2$  ( $\nu = 1$ ), we have  $\mathbf{M}^{(0)} = \mathbf{C}$ ,  $\mathbf{M}^{(1)} = \mathbf{G}$  and  $\mathbf{M}^{(2)} = \mathbf{G}\mathbf{C}^{-1}\mathbf{G}$ .

All three  $\mathbf{D}^{(0)}$ ,  $\mathbf{D}^{(1)}$  and  $\mathbf{D}^{(2)}$  are diagonal with elements used to describe the non-stationarity. The definition of these matrixes are

$$\begin{aligned} \mathbf{D}^{(0)} &= \text{diag}\{\mathbf{D}_i^{(0)}\} = \text{diag}\{e^{\phi_i^{(0)}}\} \\ \mathbf{D}^{(1)} &= \text{diag}\{\mathbf{D}_i^{(1)}\} = \text{diag}\{e^{\phi_i^{(1)}}\} \\ \mathbf{D}^{(2)} &= \text{diag}\{\mathbf{D}_i^{(2)}\} = \text{diag}\{\phi_i^{(2)}\} \end{aligned}$$

where

$$\phi_i^{(k)} = \mathbf{B}_{i,0}^{(k)} + \sum_{j=1}^p \mathbf{B}_{i,j}^{(k)} \theta_j, \quad i = 1, \dots, n$$

with the  $\mathbf{B}^{(k)}$  :  $n$ -by- $(p+1)$  user defined matrix.

The default stationary SPDE model uses  $\mathbf{B} = [0 \ 1 \ 0]$  (one by three) matrix for the marginal variance and uses  $\mathbf{B} = [0 \ 0 \ 1]$  (one by three) matrix for the scaling parameter  $\kappa$ . When this matrix have just one line, the another lines are formed using the same element of this first line. In the next section, we add one of the location coordinates as a fourth column to build a non stationary model.

### 7.2 An example

In this section we define a model where the variance depends in one of the coordinates. Note that the any precision matrix defined on the equation (7.1) we need  $\mathbf{M}^{(0)}$ ,  $\mathbf{M}^{(1)}$  and  $\mathbf{M}^{(2)}$  that are based on any mesh.

First, we define a polygon to define a mesh. We define an unitary square

```
p101 <- cbind(c(0,1,1,0,0), c(0,0,1,1,0))
```

and build a mesh on the polygon with

```
(mesh <- inla.mesh.create.helper(, pl01, max.edge=c(0.07,.12)))$n
```

```
[1] 924
```

Now, we define the non stationary SPDE model. The non stationarity is defined on the two  $\mathbf{B}$  matrix, one for  $\tau$  and another for  $\kappa$ . We want to define a model where the variance depends on the first coordinate. We just choose to put it into a fourth column on the  $\mathbf{B}$  matrix for  $\tau$ . Also, we need a prior for the three dimensional  $\theta$ , simplified by definition of a vector of mean and precision (considering independence of the priors).

```
spde <- inla.spde2.matern(mesh, B.tau=cbind(0, 1, 0, sin(pi*mesh$loc[,1])),
                           B.kappa=cbind(0, 0, 1, 0), ##mesh$loc[,1]),
                           theta.prior.mean=rep(0, 3),
                           theta.prior.prec=rep(1, 3))
```

To finalize the precision matrix definition, we need to define the values of  $\theta$ . Just to test, we define two different precision matrix, both with non stationarity based on the same  $\mathbf{B}$  matrix. The  $\theta$  vector has length equal the number of columns of  $\mathbf{B}$  minus one. The two different  $\theta$  vectors are

```
theta1 <- c(-1, 2, -1)
theta2 <- c(-1, 2, 1)
```

and we have both the precision matrix with

```
Q1 <- inla.spde2.precision(spde, theta=theta1)
Q2 <- inla.spde2.precision(spde, theta=theta2)
```

The first and second values of these vectors are the same. The structure of the  $\mathbf{B}$  matrix indicate that we have

$$\tau_i = e^{\theta_1 + \theta_3 \sin(\pi \text{loc}[i,1])}$$

and we have that the second precision matrix with larger values of  $\tau_i$ , because its values are increased by a positive value.

To make more clear, we compute both covariance matrix implied. The covariance matrix of

$$x(s) = \sum_{k=1}^n A_k(s) w_k$$

is easy to obtain when we have  $s$  (the locations) equals the locations of the mesh vertices. It is just the inverse of the precision matrix defined.

```
cov1 <- inla.qinv(Q1);          cov2 <- inla.qinv(Q2)
```

and we see a summary of the variance (diagonal of the covariance matrix) for both covariance matrix

```
v1 <- diag(cov1);          v2 <- diag(cov2)
rbind(v1=summary(v1), v2=summary(v2))
```

```
      Min. 1st Qu.  Median    Mean 3rd Qu.  Max.
v1 0.004241 0.011520 0.031180 0.04007 0.06624 0.1786
v2 0.001495 0.002265 0.004924 0.01621 0.01432 0.1565
```

and we see that the first has larger variances. Its because the precision is less than the second.

We see the variance of both process on the Figure 7.1 by

```
par(mar=c(3,3,.5,.5), mgp=c(1.7, .5, 0), las=1)
plot(mesh$loc[,1], v1, ylim=range(v1,v2),
      xlab='x-coordinate', ylab='variance', las=1)
points(mesh$loc[,1], v2, col=2)
```



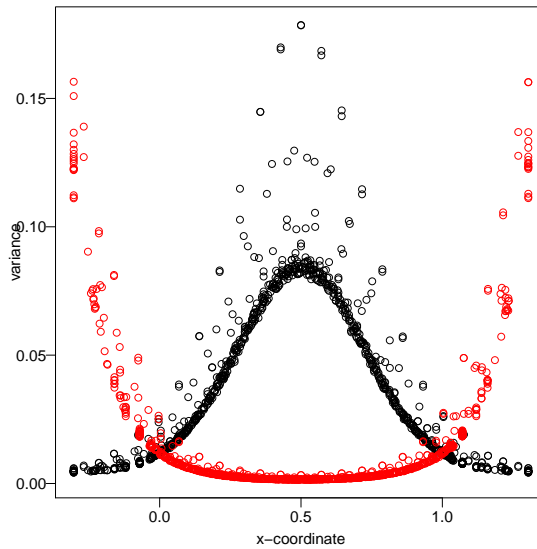


Figure 7.1: Variances implied by both non stationary process defined.

### 7.3 Simulation on the mesh vertices

Both precision matrix defined on previous the section consider that the locations are the triangles vertices of mesh. So, the simulation made with it is a realization of the random field on each point of the triangles vertices of the mesh. We use the same seed for each simulation, just to show it.

```
sample1 <- as.vector(inla.qsample(1, Q1, seed=1))
sample2 <- as.vector(inla.qsample(1, Q2, seed=1))
```

We compute the standard deviations for both the samples considering groups defined in accord to the first coordinate of the locations:

```
tapply(sample1, round(inla.group(mesh$loc[,1], 5),3), var)

-0.073      0.18      0.496      0.825      1.071
0.008349852 0.021988985 0.051450362 0.022831676 0.007851851

tapply(sample2, round(inla.group(mesh$loc[,1], 5),3), var)

-0.073      0.18      0.496      0.825      1.071
0.026432644 0.003051414 0.001069714 0.004342003 0.024451202
```

We observe that the variance of the sample from the first random field increase near 0.5 and decrease near 0 and near 1. For the sample of the second random field the opposite happens and we have larger values, because we have less precision.

We see the simulated values projected on a grid on Figure 7.2. We use a projector matrix to project the simulated values on the grid limited on the unit square with limits (0,0) and (1,1) with

```
proj <- inla.mesh.projector(mesh, xlim=0:1, ylim=0:1)
grid.arrange(levelplot(inla.mesh.project(proj, field=sample1),
                        xlab='', ylab='', scale=list(draw=FALSE),
                        col.regions=topo.colors(100)),
              levelplot(inla.mesh.project(proj, field=sample2),
                        xlab='', ylab='', scale=list(draw=FALSE),
                        col.regions=topo.colors(100)), nrow=1)
```

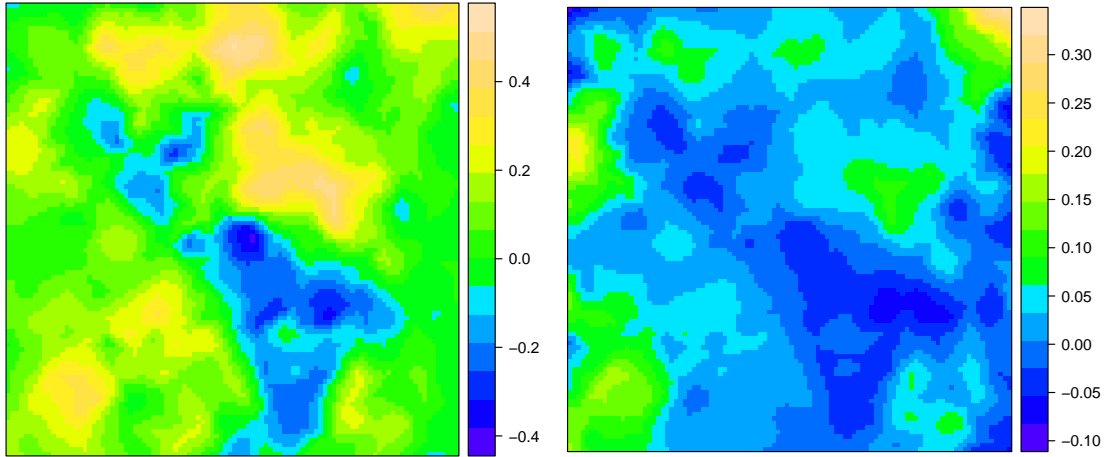


Figure 7.2: Two simulated random fields, using two different  $\theta$  on same basis functions.

### 7.3.1 Simulation with linear constraint

The linear constraint is common on models such as random walk, in one or two dimensions. Its because these models are improper. In this section we just want to show how we do linear constraint in the SPDE models.

Because the SPDE models are based on the Finite Element Method (FEM) approximation, the sum-to-zero restriction in this case is non trivial.

The issue is that

$$\sum_k x_k$$

doesn't mean anything for the mesh-based spde-models. Whereas

$$\int x(s)ds = \int \Psi(s)x_k ds$$

does mean something, and that integral is equal to

$$\sum_k C_k k x_k$$

so the constraint

$$\int x(s)ds = 0$$

is provided by

$$A * x = 0, \text{ where } A_{1k} = C_{kk}$$

Where  $C$  is the matrix used on the FEM.

Using  $A = (1, \dots, 1)$  instead of  $\text{diag}(C)$  leads to very bad behaviour for irregular meshes. So, if we want a linear constraint, we need to use  $C$ .

That matrix is obtained by the `inla.mesh.fem()` function and is available directly in outputs of `inla.spde2.matern()` function. So, we do the simulation with

```
s1r <- as.vector(inla.qsample(1, Q1, seed=1, constr=spde$f$extraconstr))
s2r <- as.vector(inla.qsample(1, Q2, seed=1, constr=spde$f$extraconstr))
```

and we have

```

rbind(s1r=summary(s1r), s2r=summary(s2r))

      Min.   1st Qu.   Median     Mean 3rd Qu.    Max.
s1r -0.4399 -0.03012 0.05297 0.06406 0.15710 0.7765
s2r -0.5973 -0.01840 0.02174 0.02401 0.06044 0.4819

c(cor1=cor(sample1, s1r), cor2=cor(sample2, s2r))

cor1 cor2
  1     1

```

where the mean of the process simulated on the mesh vertices have mean near zero.

## 7.4 Estimation with data simulated on the mesh vertices

The model can be fitted easily with the data simulated on mesh vertices. Considering that we have data exactly on each vertice of the mesh, we don't need the use of any predictor matrix and the stack functionality. Because we have just realizations of the random field, we don't have noise and need to fix the precision of the Gaussian likelihood on high value, for example on the value  $e^{20}$

```

klik <- list(hyper=list(theta=list(initial=20, fixed=TRUE)))

```

Remember that we have a zero mean random field, so we also don't have fixed parameters to fit. We just do

```

formula <- y ~ 0 + f(i, model=spde)
fit1 <- inla(formula, control.family=klik,
             data=data.frame(y=sample1, i=1:mesh$n))
fit2 <- inla(formula, control.family=klik,
             data=data.frame(y=sample2, i=1:mesh$n))

```

We look at the summary of the posterior for  $\theta$  (joined with the true values). For the first sample

```

round(cbind(true=theta1, fit1$summary.hy), 4)

```

	true	mean	sd	0.025quant	0.5quant	0.975quant	mode
Theta1 for i	-1	-0.9554	0.0295	-1.0193	-0.9527	-0.9046	-0.9439
Theta2 for i	2	1.9045	0.1233	1.6893	1.8940	2.1704	1.8555
Theta3 for i	-1	-0.9676	0.0505	-1.0517	-0.9732	-0.8563	-0.9903

and for the second

```

round(cbind(true=theta2, fit2$summary.hy), 4)

```

	true	mean	sd	0.025quant	0.5quant	0.975quant	mode
Theta1 for i	-1	-1.0107	0.0333	-1.0753	-1.0111	-0.9447	-1.0120
Theta2 for i	2	2.0656	0.1059	1.8481	2.0698	2.2631	2.0809
Theta3 for i	1	1.0675	0.0443	0.9814	1.0672	1.1552	1.0663

We see that for both we have good results. In next section we see more results.

## 7.5 Estimation with locations not on mesh vertices

Suppose that we have the data on the locations simulated by the commands below

```
set.seed(2);      n <- 100
loc <- cbind(runif(n), runif(n))
```

Now, we project the data simulated on the mesh vertices to this locations. To do it, we need a projection matrix

```
projloc <- inla.mesh.projector(mesh, loc)
```

with

```
x1 <- inla.mesh.project(projloc, sample1)
x2 <- inla.mesh.project(projloc, sample2)
```

and we have the sample data on these locations.

Now, because the this locations aren't vertices of the mesh, we need to use the stack functionality. First, we need the predictor matrix. But this is the same used to 'sample' the data.

And we define the stack for each one of the samples

```
stk1 <- inla.stack(list(y=x1), A=list(projloc$proj$A), tag='d',
                  effects=list(data.frame(i=1:mesh$n)))
stk2 <- inla.stack(list(y=x2), A=list(projloc$proj$A), tag='d',
                  effects=list(data.frame(i=1:mesh$n)))
```

And we fit the model with

```
res1 <- inla(formula, data=inla.stack.data(stk1), control.family=clik,
             control.predictor=list(compute=TRUE, A=inla.stack.A(stk1)))
res2 <- inla(formula, data=inla.stack.data(stk2), control.family=clik,
             control.predictor=list(compute=TRUE, A=inla.stack.A(stk2)))
```

The true and summary of marginal posterior distribution for  $\theta$ :

```
round(cbind(True=theta1, res1$summary.hy), 4)
```

	True	mean	sd	0.025quant	0.5quant	0.975quant	mode
Theta1 for i	-1	-1.0496	0.1817	-1.4073	-1.0494	-0.6936	-1.0487
Theta2 for i	2	2.1947	0.1862	1.8175	2.1996	2.5481	2.2140
Theta3 for i	-1	-0.9677	0.2424	-1.4429	-0.9684	-0.4900	-0.9708

```
round(cbind(True=theta2, res2$summary.hy), 4)
```

	True	mean	sd	0.025quant	0.5quant	0.975quant	mode
Theta1 for i	-1	-0.9624	0.1771	-1.3091	-0.9630	-0.6124	-0.9654
Theta2 for i	2	2.1677	0.1900	1.7823	2.1733	2.5263	2.1894
Theta3 for i	1	0.9308	0.2413	0.4562	0.9311	1.4041	0.9323

To make the visualization more good, we take the logarithmum of the variance.

```
x1.mean <- inla.mesh.project(proj, field=res1$summary.ran$i$mean)
x1.var <- inla.mesh.project(proj, field=res1$summary.ran$i$sd^2)
x2.mean <- inla.mesh.project(proj, field=res2$summary.ran$i$mean)
x2.var <- inla.mesh.project(proj, field=res2$summary.ran$i$sd^2)
```

We visualize, for both random fields, the simulated, the predicted (posterior mean) and the posterior variance on Figure 7.3 with commands below

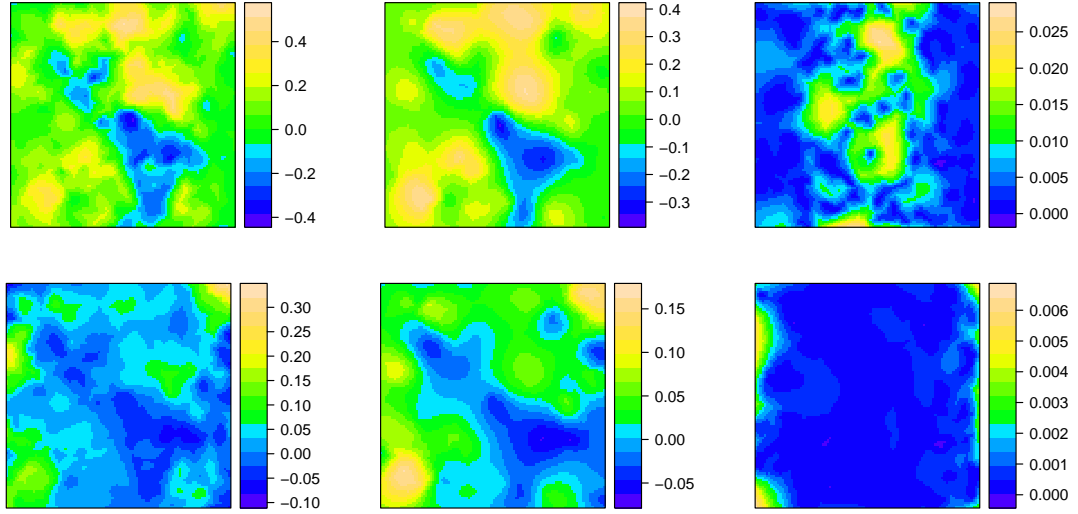


Figure 7.3: Simulated (top and bottom left), posterior mean (top and bottom mid) and the posterior variance (top and bottom right) for both random fields.

```
do.call(function(...) grid.arrange(..., nrow=2),
  lapply(list(inla.mesh.project(proj, sample1), x1.mean, x1.var,
    inla.mesh.project(proj, sample2), x2.mean, x2.var),
    levelplot, xlab='', ylab='',
    col.regions=topo.colors(100), scale=list(draw=FALSE)))
```

We see on the Figure 7.3 that the predicted values are similar to the simulated ones. Also, we see that the posterior variance of the first model increase near 0.5 on the first coordinate. And we see the oposite for the second random field. Also, we see that the variance of the first is greater than the second.

## Chapter 8

# Point process: inference on the log-Cox process

Under the log-Cox model assumption, there is a latent Gaussian Random Field (LGRF) and the inference can be done using **INLA** [Illian et al., 2012].

A common approach to fit the log-Cox process is to divide the study region into cells, that forms a lattice, and count the number of points into each one. This counts are modeled using the Poisson likelihood.

A good approach for inference of log-Cox model is to use the SPDE approach instead occurrence counts on cells, [Simpson et al., 2011]. There are the main advantages on this approach the consideration of the loc-Cox likelihood directly.

### 8.1 Data simulation

The data simulated on this section is used also on the Chapter 9.

In this section we use the `rLGCP()` function from **spatstat** package to do the simulation. By default that function do simulation on window over the  $(0, 1) \times (0, 1)$  square. We choose to do simulation over the  $(0, 3) \times (0, 3)$  square.

```
require(spatstat)
win <- owin(c(0,3), c(0,3))
```

This function uses the `GaussRF()` function from the **RandomFields** package. The `rLGCP` uses the `GaussRF()` function to do simulation of the LGRF over a grid on the provided window and use it to do the point process simulation.

There is an internal parameter to control the resolution of the grid. We change it to

```
spatstat.options(npixel=100)
```

First we define the model parameter for the model is the mean of the LGRF. This is directly related to expected number of points of the spatial pattern. The expected number of points is its exponential times the area fo the window. We use

```
beta0 <- 3
```

So, the expected number of points is

```
exp(beta0) * diff(range(win$x)) * diff(range(win$y))
```

```
[1] 180.7698
```

Is also possible to use a functional.

On the estimation process we use the Matern covariance function with  $\nu = 1$ . So, here we just fix it on this value. The other parameters are the variance and scale

```
sigma2x <- 1;      kappa <- 1
```

Doing the simulation

```
set.seed(1)
lg.s <- rLGCP('matern', beta0,
              c(0, variance=sigma2x, nugget=0, scale=1/kappa, nu=1), win=win)
```

Both, the LGRF and the point pattern, are returned. The point pattern locations are

```
(n <- nrow(xy <- cbind(lg.s$x, lg.s$y)))
```

```
[1] 137
```

The exponential of simulated values of the LGRF are returned as the `Lambda` attribute of the object. We extract the  $\Lambda$  and see a summary of the  $\log(\Lambda)$  below

```
Lam <- attr(lg.s, 'Lambda')
summary(as.vector(rf.s <- log(Lam$v)))
```

```
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.2869   2.0440   2.5770   2.5580   3.0650   4.4530
```

The values of the LGRF is collected (and a summary) at closest grid centers with

```
summary(z <- log(Lam$v[Reduce(
  'cbind', nearest.pixel(xy[,1], xy[,2], Lam))]))
```

```
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.8845   2.5060   3.0310   3.0120   3.4750   4.3890
```

at point pattern locations we have larger average than the entire simulated random field.

We see the exponential of the simulated LGRF over the grid, left graph of Figure 8.1, and the point pattern with size of points proportional to the LGRF at closest grid center.

```
par(mfrow=c(1,2), mar=c(0,0,1,.5), mgp=c(1,0.5,0))
plot(Lam, main='Lambda')
par(mar=c(2,2,2.5,0.5))
plot(xy, cex=0.3+(z-min(z))/diff(range(z)), asp=1, xlab='', ylab='',
      main='Simulated point pattern\npoint size proportional of values at')
```

## 8.2 Inference

Following[Simpson et al., 2011] we can estimate the parameters of the log-Cox point process model using few command lines.

Creating the mesh

```
(nv <- (mesh <- inla.mesh.2d(loc.dom=as.matrix(expand.grid(win$x, win$y)),
                           cutoff=.15, max.edge=c(.3, .5)))$n)
```

```
[1] 439
```

Defining the SPDE model

```
spde <- inla.spde2.matern(mesh=mesh, alpha=2)
```

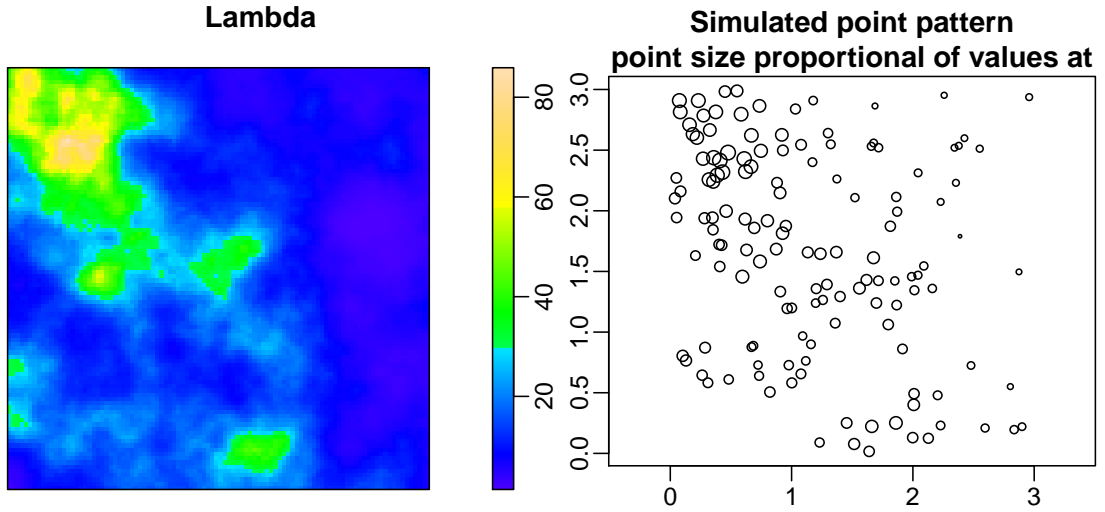


Figure 8.1: Simulated intensity of the point process (left), simulated point pattern (right).

The trick to the estimation of the log-Cox model is to use the 'relative' area of the each mesh vertices as the expected number of counts on this locations (the integration points). But, at data locations, it needs to set the expected counts as zero (it is just a point). So, we define a vector of data augmented to have length equals the number of edges on the mesh plus the number of events.

```
y.pp <- rep(0:1, c(nv, n))
e.pp <- c(diag(spde$param.inla$M0), rep(0,n))
```

We must have to define the projector matrix to do inference using the SPDE approach, [Lindgren, 2012]. For the observed points locations we have

```
lmat <- inla.spde.make.A(mesh, xy)
```

We need also a projector matrix for the integration points and this is just a diagonal matrix because this locations are just the mesh vertices.

```
imat <- Diagonal(nv, rep(1, nv))
```

So, the entire projector matrix is

```
A.pp <- rBind(imat, lmat)
```

The data stack can be made by

```
stk.pp <- inla.stack(data=list(y=y.pp, e=e.pp),
                    A=list(A.pp), tag='pp',
                    effects=list(i=1:nv))
```

The posterior marginals for the parameters of the log-Cox model (on the SPDE approach scale) and for the latent random field at integration and location points are obtained by

```
pp.res <- inla(y ~ 0 + f(i, model=spde),
              family='poisson', data=inla.stack.data(stk.pp),
              control.predictor=list(A=inla.stack.A(stk.pp)),
              E=inla.stack.data(stk.pp)$e)
```

To get the model parameters on the user-scale parameters such as the scale  $\kappa$ , nominal variance  $\sigma_x^2$  and nominal range we use the `inla.spde2.result()` function

The posterior distribution of the log-Cox model parameters are visualized on the Figure 8.2.



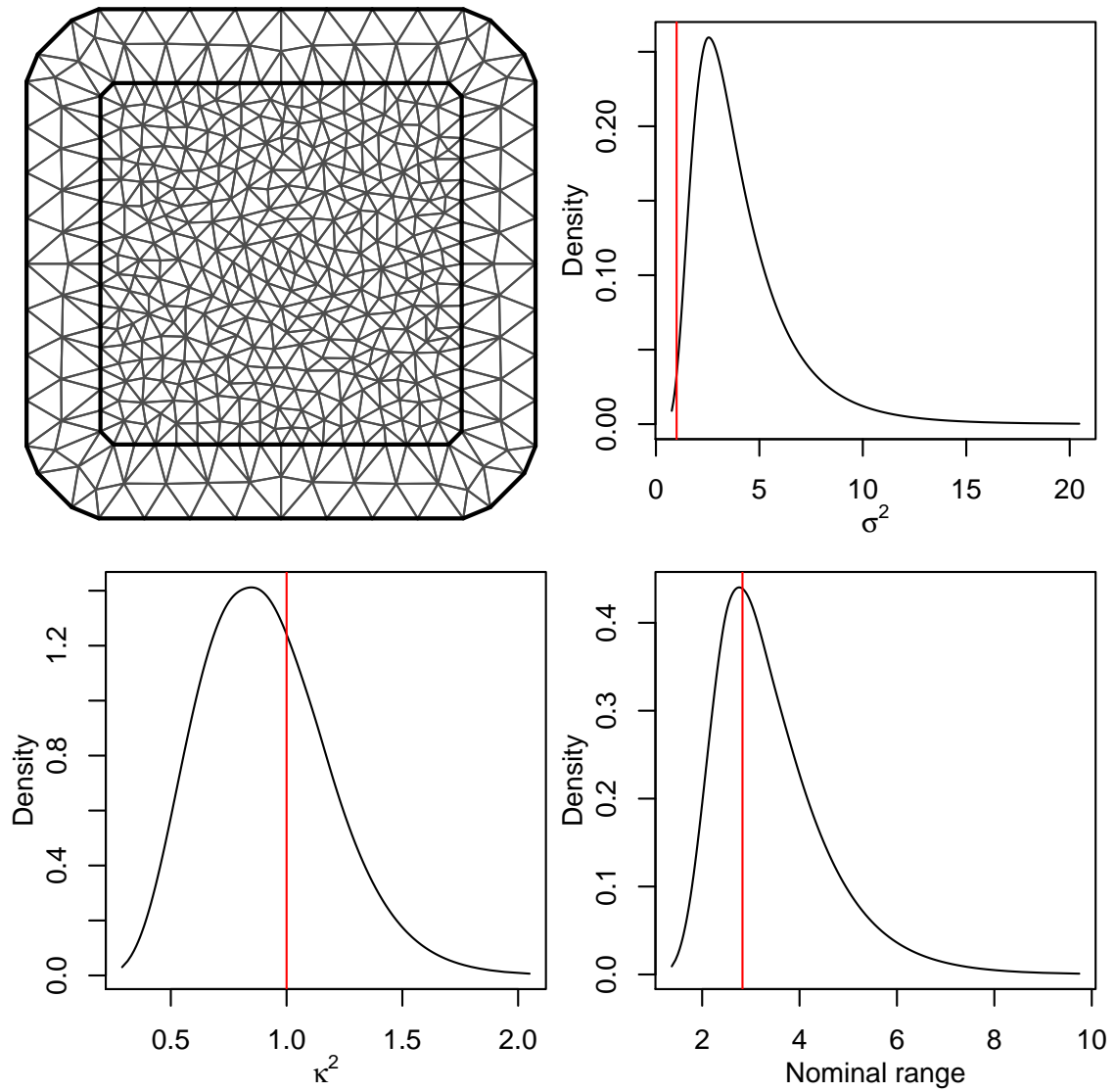


Figure 8.2: Posterior distribution for the parameters of the log-Cox model  $\sigma^2$  (left),  $\kappa$  (mid) and the nominal range (right)

```

par(mfrow=c(2,2), mar=c(0,0,0,0))
plot(mesh, asp=1, main='')
par(mar=c(3,3,1,0.3), mgp=c(2,1,0))
plot(pp.rf$marginals.variance.nominal[[1]], type='l',
      xlab=expression(sigma^2), ylab='Density')
abline(v=sigma2x, col=2)
plot(pp.rf$marginals.kappa[[1]], type='l',
      xlab=expression(kappa^2), ylab='Density')
abline(v=kappa, col=2)
plot(pp.rf$marginals.range.nominal[[1]], type='l',
      xlab='Nominal range', ylab='Density')
abline(v=sqrt(8*1)/kappa, col=2)

```

## Chapter 9

# Geostatistical inference under preferential sampling

In some cases the effort on sampling depends on the response. For example, is more common to have stations collecting data about pollution on industrial area than on rural ones. To make inference in this case, we can test if we have a preferential sampling problem in our data. One approach is to build a joint model considering a log-Cox model for the point pattern (the locations) and the response, [Diggle et al., 2010]. So, we need also to make inference for a point process model jointly.

An illustration of the use **INLA** for the preferential sampling problem is on the case studies section of the **INLA** web page, precisely on <http://www.r-inla.org/examples/case-studies/diggle09>. This example uses the two dimensional random walk model for the latent random field. Here, we show geoestatistical inference under preferential sampling using SPDE.

We use the values of the LGRF simulated on the Chapter 8 to define the response. We just take the values of closest grid centers to each location of the point pattern, change the signal and divide by two. So, the response values now are inversly proportional to the points density

```
set.seed(2)
summary(resp <- -z + rnorm(length(z), 0, 0.2))
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-4.4850	-3.5470	-2.9270	-3.0040	-2.5650	-0.8266

### 9.1 Fitting the usual model

Here, we just fit the geoestatistical model using the usual approach. In this approach we just use the locations as fixed. We use the mesh of the previous Chapter.

```
stk.u <- inla.stack(data=list(y=resp), A=list(lmat, 1),
                   effects=list(i=1:nv, beta=rep(1,length(resp))))
u.res <- inla(y ~ 0 + beta + f(i, model=spde),
             data=inla.stack.data(stk.u),
             control.predictor=list(A=inla.stack.A(stk.u)))
```

Summaries

```
rbind(beta0=u.res$summary.fix[, 1:6], noise.prec=u.res$summary.hy[1,])
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
beta0	-2.523138	0.5161457	-3.592109	-2.532296	-1.41126	-2.546876
noise.prec	14.781615	2.8274426	9.956685	14.536643	21.02973	14.073799

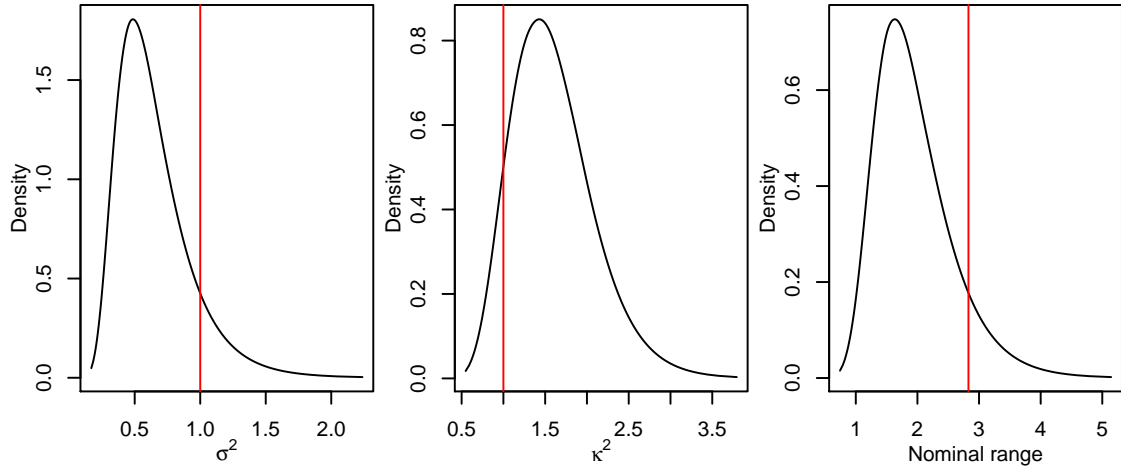


Figure 9.1: Posterior distribution for  $\sigma^2$ ,  $\kappa$  and the nominal range just using the response.

We have to build also the posterior marginals on the user-scale parameters: scale, nominal variance and nominal range

```
u.rf <- inla.spde2.result(u.res, 'i', spde)
```

In the Figure 9.1 we see the marginals posterior distributions for the model parameters with code bellow

```
par(mfrow=c(1,3), mar=c(3, 3, 0.3, 0.3), mgp=c(2,1,0))
plot(u.rf$marginals.variance.nominal[[1]], type='l',
     xlab=expression(sigma^2), ylab='Density')
abline(v=sigma2x, col=2)
plot(u.rf$marginals.kappa[[1]], type='l',
     xlab=expression(kappa^2), ylab='Density')
abline(v=kappa, col=2)
plot(u.rf$marginals.range.nominal[[1]], type='l',
     xlab='Nominal range', ylab='Density')
abline(v=sqrt(8*1)/kappa, col=2)
```

## 9.2 Estimation under preferential sampling

In this situation we fit the model where a LGRF is considered to model both point pattern and the response. Using **INLA** it can be done using two likelihoods, one for the point pattern and another for the response. To do it we need a matrix response and a new index set to specify the model for the LGRF. It is more easy by using the `inla.stack()` following previous examples for two likelihood models.

We consider the point pattern 'observation' on the first column and the response values on the second column. So, we just redefine the stack for the response and also for the point process. We put the response on the first column and the Poisson data for the point process as the second column. Also, to avoid the expected number of cases as NA for the Poisson likelihood, we set it as zero on the response data stack. For the SPDE effect on the point process part we have to model it as a copy of the SPDE effect at response part. We do it by defining a index set with different name and use it on the copy feature later.

```
stk2.y <- inla.stack(data=list(y=cbind(resp,NA), e=rep(0,n)),
                    A=list(lmat, 1), tag='resp2',
                    effects=list(i=1:nv, beta=rep(1,n)))
stk2.pp <- inla.stack(data=list(y=cbind(NA,y.pp), e=e.pp),
```

```
A=list(A.pp), tag='pp2',
effects=list(j=1:nv))
```

and join both together

```
j.stk <- inla.stack(stk2.y, stk2.pp)
```

Now, we fit the geostatistical model under preferential sampling. To put the LGRF on both likelihood, we have to use the copy strategy.

```
j.res <- inla(y ~ 0 + f(i, model=spde) +
              f(j, copy='i', fixed=FALSE),
              data=inla.stack.data(j.stk),
              family=c('gaussian', 'poisson'),
              control.predictor=list(A=inla.stack.A(j.stk)),
              E=inla.stack.data(j.stk)$e)
```

Computing the marginals posterior distributions on the user-scale random field parameters

```
j.rf <- inla.spde2.result(j.res, 'i', spde)
```

We can visualize the posterior distribution on the Figure9.2. The negative signal to the copy parameter  $\beta$  is due to the fact that we define the response with opposite signal of the latent field used to do the simulation of the point process.

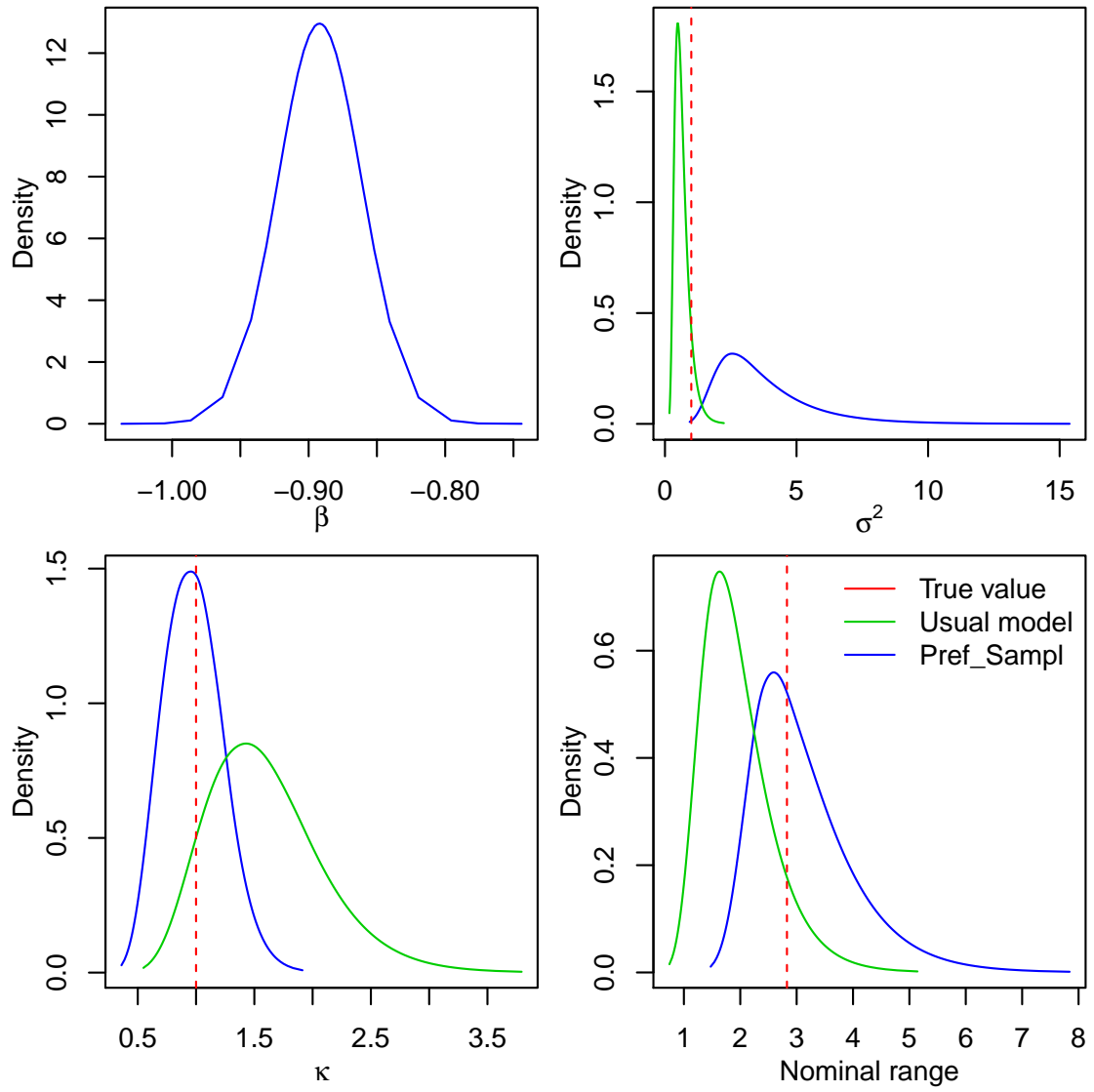


Figure 9.2: Posterior distribution for  $\beta_0$ ,  $\sigma^2$ ,  $\kappa$  and the nominal range under preferential sampling.

## Chapter 10

# A space time example

In this chapter we show an example on fitting a spatio-temporal model. This model is a separable one described on [Cameletti et al., 2012]. Basically the model is defined as a SPDE model for the spatial domain and another one for the temporal domain. The spatio-temporal separable model is defined by the kronecker product between the precision of these two models.

### 10.1 Data simulation

We use the Paraná state border, available on **INLA** package, as the domain.

```
data(PRborder)
```

We start by defining the spatial model. Because we need that the example run faster, we use the low resolution mesh for Paraná state border created on the section 2.2.

There is two options to simulate from Cameletti's model. One is based on the marginal distribution of the latent field and another is on the conditional distribution at each time. The second one is easy for us by use the function defined on the section 2.4 to simulate an independent random field realization each time.

First we define a set of points as the same on the **PRprec** data

```
data(PRprec)
coords <- as.matrix(PRprec[,1:2])
```

and set  $k = 10$ , the time dimention

```
k <- 12
```

The  $k$  independent realizations can be done by

```
params <- c(variance=1, kappa=1/2)
set.seed(1)
x.k <- rspde(coords, kappa=1/2, n=k, mesh=prmesh1)
```

```
theta = -0.5723649 -0.6931472
```

```
dim(x.k)
```

```
[1] 616 12
```

Now, we define the autoregressive parameter  $\rho$

```
rho <- 0.7
```

and get the sample

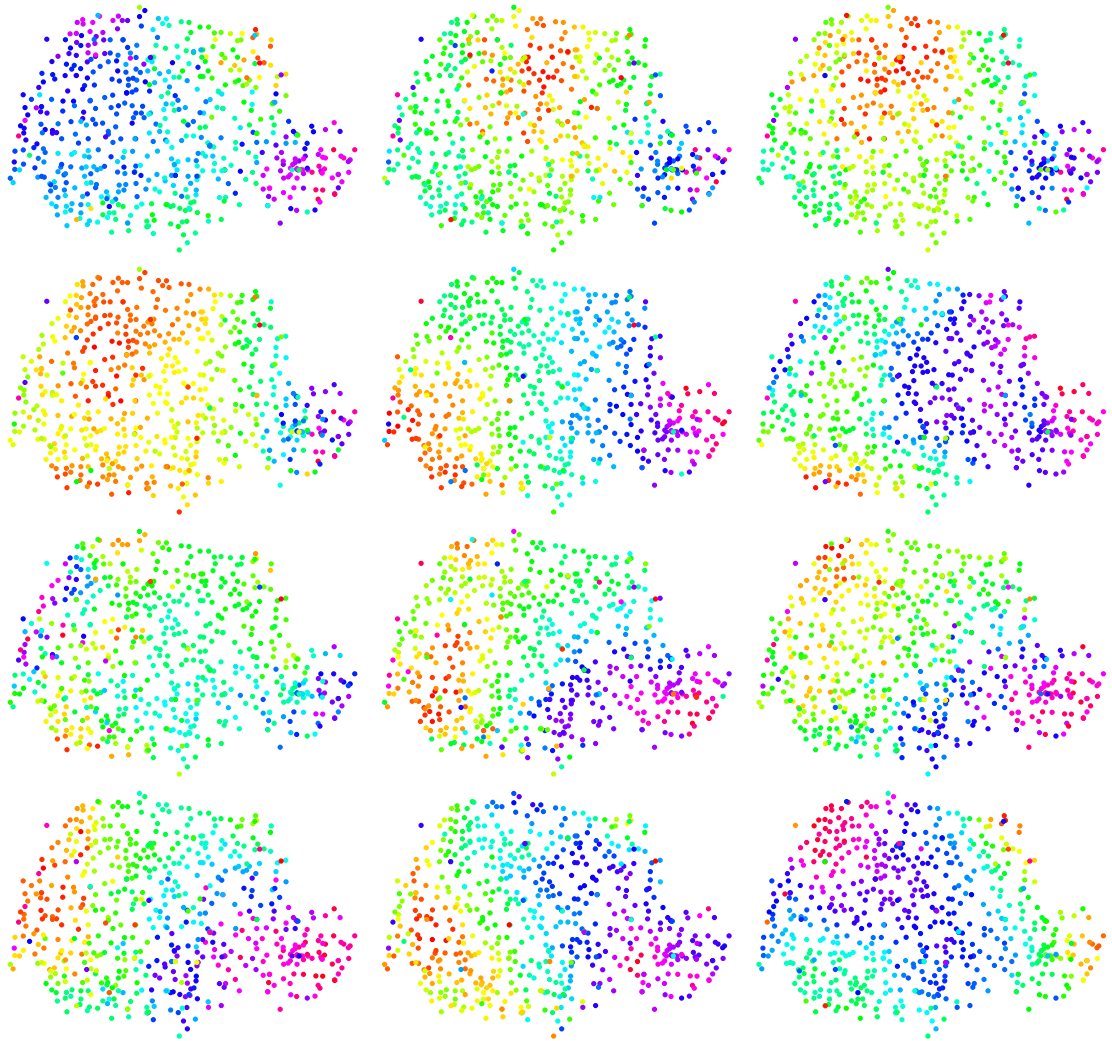


Figure 10.1: Realization of the spatio temporal model.

```
x <- x.k
for (j in 2:k)
  x[,j] <- rho*x[,j-1] + x.k[,j]
```

We can visualize the realization at the figure 10.1 with commands bellow

```
c100 <- rainbow(101)
par(mfrow=c(4,3), mar=c(0,0,0,0))
for (j in 1:k)
  plot(coords, col=c100[round(100*(x[,j]-min(x[,j]))/diff(range(x[,j])))],
        axes=FALSE, asp=1, pch=19, cex=0.5)
```

In this example we need to show the use of a categorical covariate. First we do the simulation of the covariate as

```
n <- nrow(coords)
set.seed(2)
table(ccov <- factor(sample(LETTERS[1:3], n*k, replace=TRUE)) )
```

```

  A    B    C
2458 2438 2496
```

and the regression parameters as

```
beta <- -1:1
```

The response is

```
y <- beta[unclass(ccov)] + x
tapply(y, ccov, mean)
```

```
      A      B      C
-2.128641 -1.149361 -0.138450
```

## 10.2 Data stack preparation

Now, we need the data preparation to build the spatio-temporal model. The index set is made taking into account the number of weights on the SPDE model and the number of groups

```
iset <- inla.spde.make.index('i', n.spde=attr(x.k, 'spde')$n.spde, n.group=k)
```

The projector matrix must be defined considering the temporal dimension. This can be defined using the group argument on the `inla.spde.make.A()` function

```
A <- inla.spde.make.A(prmesh1, cbind(rep(coords[,1], k),
                                     rep(coords[,2], k)),
                    group=rep(1:k, each=n))
```

The effects on the stack are a list with two elements, one is the index set and another the categorical covariate

```
str(list(iset, ccov=ccov))
```

List of 2

```
$      :List of 3
..$ i      : int [1:2244] 1 2 3 4 5 6 7 8 9 10 ...
..$ i.group: int [1:2244] 1 1 1 1 1 1 1 1 1 1 ...
..$ i.repl : int [1:2244] 1 1 1 1 1 1 1 1 1 1 ...
$ ccov: Factor w/ 3 levels "A","B","C": 1 3 2 1 3 3 1 3 2 2 ...
```

The stack data is defined as

```
sdat <- inla.stack(tag='stdata', data=list(y=as.vector(y)),
                  A=list(A,1), effects=list(iset,ccov=ccov))
```

## 10.3 Fitting the model and some results

We set the prior on the temporal autoregressive parameter with zero mean and variance equals 5, and define the initial value equals the true with

```
h.spec <- list(theta=list(initial=0.7, param=c(0, 5)))
```

The likelihood hyperparameter is fixed on a high precision, just because we haven't noise. We choose the `strategy='gaussian'` to it run faster. To deal with the categorical covariate we need to set `expand.factor.strategy='inla'` on the `control.fixed` argument list.



```

res <- inla(y ~ 0 + ccov +
  f(i, model=attr(x.k, 'spde'), group=i.group,
    control.group=list(model='ar1', hyper=h.spec)),
  data=inla.stack.data(sdat),
  control.predictor=list(compute=TRUE, A=inla.stack.A(sdat)),
  control.family=list(initial=20, fixed=TRUE),
  control.inla=list(strategy='gaussian'),
  control.fixed=list(expand.factor.strategy='inla'))

```

The summary of the covariate coefficients (together the observed mean for each covariate level)

```
round(cbind(observed=tapply(y, ccov, mean), res$summary.fixed), 4)
```

	observed	mean	sd	0.025quant	0.5quant	0.975quant	mode	kld
A	-2.1286	-2.094	0.9778	-4.0366	-2.0942	-0.1532	-2.0945	0
B	-1.1494	-1.094	0.9778	-3.0366	-1.0942	0.8468	-1.0945	0
C	-0.1384	-0.094	0.9778	-2.0366	-0.0942	1.8468	-0.0945	0

The summary for the posterior marginal distribution of the temporal correlation

```
round(res$summary.hy[3,], 5)
```

	mean	sd	0.025quant	0.5quant	0.975quant	mode
	0.66390	0.01694	0.63020	0.66400	0.69679	0.66404

Bellow, we get the marginals distributions for random field parameters on the user scale

```
rf <- inla.spde2.result(res, 'i', attr(x.k, 'spde'), do.transf=TRUE)
```

and we see these distributions, also the marginal distribution for the temporal correlation, on the Figure 10.2 with the commands bellow

```

par(mfrow=c(2,2), mar=c(3,3,1,0.1), mgp=2:0)
plot(res$marginals.hyper[[3]], type='l',
  xlab=expression(beta), ylab='Density')
abline(v=rho, col=2)
plot(rf$marginals.variance.nominal[[1]], type='l',
  xlab=expression(1/sigma[x]), ylab='Density')
abline(v=1/params[1], col=2)
plot(rf$marginals.kappa[[1]], type='l',
  xlab=expression(kappa), ylab='Density')
abline(v=params[2], col=2)
plot(rf$marginals.range.nominal[[1]], type='l',
  xlab='range nominal', ylab='Density')
abline(v=sqrt(8)/params[2], col=2)

```

## 10.4 A look at the posterior random field

The first look at the random field posterior distribution is to compare the realized random field with the posterior mean, median or/and mode and any quantile. We just compute the correlation between the simulated data response and the posterior mean of the predicted values.

First, we found the index for the random field at data locations

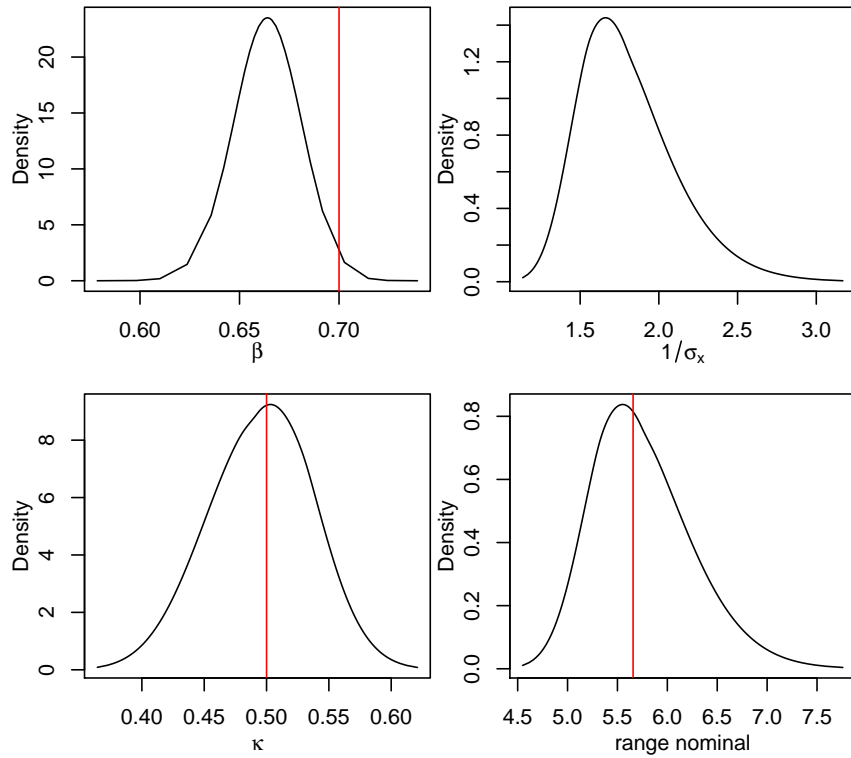


Figure 10.2: Marginal posterior distribution for  $\beta$  (top left),  $\sigma_x$  (top right),  $\kappa$  (bottom left) and the nominal range (bottom right). The red vertical lines are placed at true value.

```
str(idat <- inla.stack.index(sdat, 'stdata')$data)
```

```
int [1:7392] 1 2 3 4 5 6 7 8 9 10 ...
```

and compute the correlation between the the posterior mean and the response by

```
cor(as.vector(y), res$summary.linear.predictor$mean[idat])
```

```
[1] 1
```

We also can do prediction for each time and visualize it. First, we define the projection grid in the same way that on the example on Chapter 4.

```
stepsize <- 4*1/111
nxy <- round(c(diff(range(coords[,1])), diff(range(coords[,2])))/stepsize)
projgrid <- inla.mesh.projector(prmesh1, xlim=range(coords[,1]),
                                ylim=range(coords[,2]), dims=nxy)
```

The prediction for each time can be done by

```
xmean <- list()
for (j in 1:k)
  xmean[[j]] <- inla.mesh.project(
    projgrid, res$summary.random$i$mean[iset$i.group==j])
```

We found what points of the grid are inside the Paraná state border.

```
require(splancs)
xy.in <- inout(projgrid$lattice$loc, cbind(PRborder[,1], PRborder[,2]))
```

And, to get better visualization, we set NA to the points of the grid out of the Paraná border.

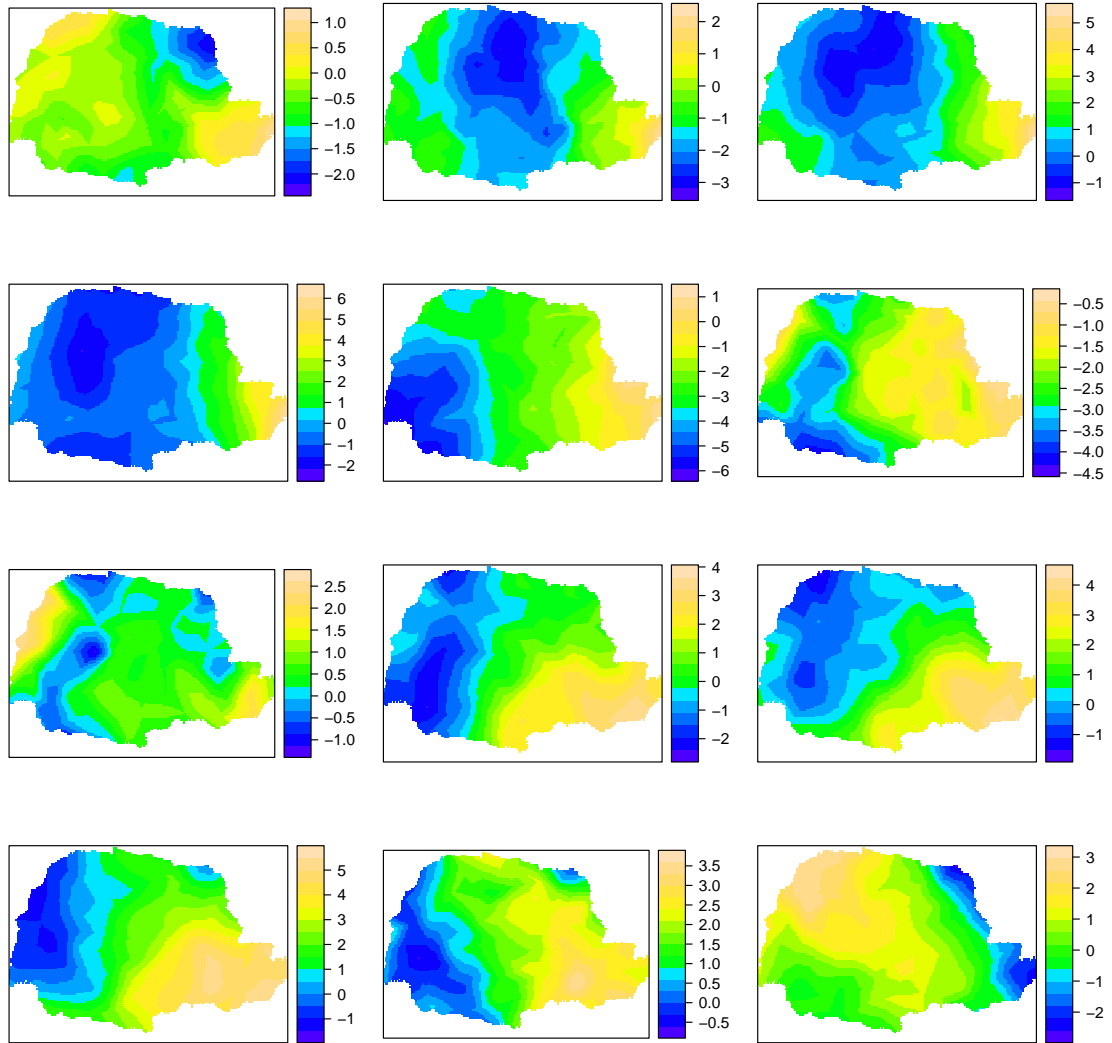


Figure 10.3: Visualization of the spatio temporal random field.

```
for (j in 1:k)  xmean[[j]][!xy.in] <- NA
```

The visualization at Figure 10.3 can be made by the comand bellow

```
require(gridExtra)
do.call(function(...) grid.arrange(..., nrow=4),
        lapply(xmean, levelplot, xlab="", ylab="",
               col.regions=topo.colors(100), scale=list(draw=FALSE)))
```

## Chapter 11

# Data cloning example

The data clone algorithm for spatial models is described in <http://www.sciencedirect.com/science/article/pii/S0167947310004329>. In the **INLA** package we have an example with the hybrid data clone when the random effect are 'iid' on <http://www.math.ntnu.no/inla/r-inla.org/papers/FS8-2010-R.pdf>

Here we use the hybrid data cloning approach to see what happens with the parameters of the toy spatial continuous model considering the SPDE model.

We use the toy data set

```
data(SPDEtoy)
coords <- as.matrix(SPDEtoy[,1:2])
```

with the mesh

```
(mesh <- inla.mesh.create.helper(coords,
                                max.e=c(0.15, 0.2), cutoff=0.02))$n
```

[1] 371

to build the model

```
spde <- inla.spde2.matern(mesh)
```

We consider a set of values for the number of clones. By cloning we use the `replica` option on the latent field definition function `f()`. The `inla.spde.make.index` function allow us to make the replica index.

In the following code we do a loop on the number o clones an fit the model with replications

```
n <- nrow(coords); k <- c(1,2,3,5,10); names(k) <- k
resk <- list()
for (i in 1:length(k)) {
  kk <- k[i]
  A <- inla.spde.make.A(mesh, loc=coords,
                        index=rep(1:n, kk), repl=rep(1:kk, each=n))
  ind <- inla.spde.make.index(name='i', n.spde=spde$n.spde, n.repl=kk)
  st.dat <- inla.stack(data=list(resp=rep(SPDEtoy[,3], kk)),
                      A=list(A,1),
                      effects=list(ind, list(m=rep(1,n*kk))), tag='est')
  resk[[i]] <- inla(resp ~ 0 + m + f(i, model=spde, replicate=i.repl),
                    data=inla.stack.data(st.dat),
                    control.predictor=list(A=inla.stack.A(st.dat),
                                           compute=TRUE))
}
```

Also we get the posterior marginal distribution of the marginal variance and practical range with

```
res.f <- lapply(1:length(resk), function(i)
  inla.spde2.result(resk[[i]], 'i', spde, do.transf=TRUE))
```

We collect the results of the posterior mean and posterior variance of each parameter with

```
r <- list()
r$intercept <- sapply(resk, function(x)
  c(x$summary.fix[1,1], x$summary.fix[1,2]^2))
r$variance.lik <- sapply(resk, function(x) {
  d <- inla.tmarginal(function(x) 1/x, x$marginals.hyperpar[[1]])
  r <- inla.emarginal(function(x) x, d)
  c(r, inla.emarginal(function(x) (x-r)^2, d))
})
r$marginal.variance <- sapply(res.f, function(x) {
  e <- inla.emarginal(function(x) x, x$marginals.variance.nom[[1]])
  c(e, inla.emarginal(function(x) (x - e)^2,
    x$marginals.variance.nominal[[1]]))
})
r$kappa <- sapply(resk, function(x) {
  e <- inla.emarginal(function(x) exp(x), x$marginals.hy[[3]])
  c(e, inla.emarginal(function(x) (exp(x) - e)^2, x$marginals.hy[[3]]))
})
```

The main think in data clonning is to visualize the posterior mean and the posterior variance of the parameters. When we don't have identificability problem we see that the mean converges to the likelihood estimate and the posterior variance converges to zero.

We visualize the posterior mean and posterior variance of each parameter ( $\beta_0$ ,  $\sigma_y^2$ ,  $\sigma_x^2$  and  $\kappa$ ) relative to the result without clone ( $k = 1$ ). Also, to compare, we visualize the likelihood estimate. This graphs are on the Figure 11.1 with

```
par(mfrow=c(2, 2), mar=c(3,3,2,1), mgp=c(1.5,.5,0))
for (i in 1:length(r)) {
  plot(k, r[[i]][2, ]/r[[i]][2,1], type='o', ylim=c(0, 1.2),
    ylab='relative to k=1', main=names(r)[i])
  lines(k, r[[i]][1, ]/r[[i]][1,1], type='o', lty=2, col=2)
  abline(h=1, col=4, lty=3)
  abline(h=1k.est[i]/r[[i]][1,1], col=3, lty=3)
}
legend('bottomleft', c('Posterior mean', 'Posterior variance'),
  lty=1, col=2:1, bty='n')
```

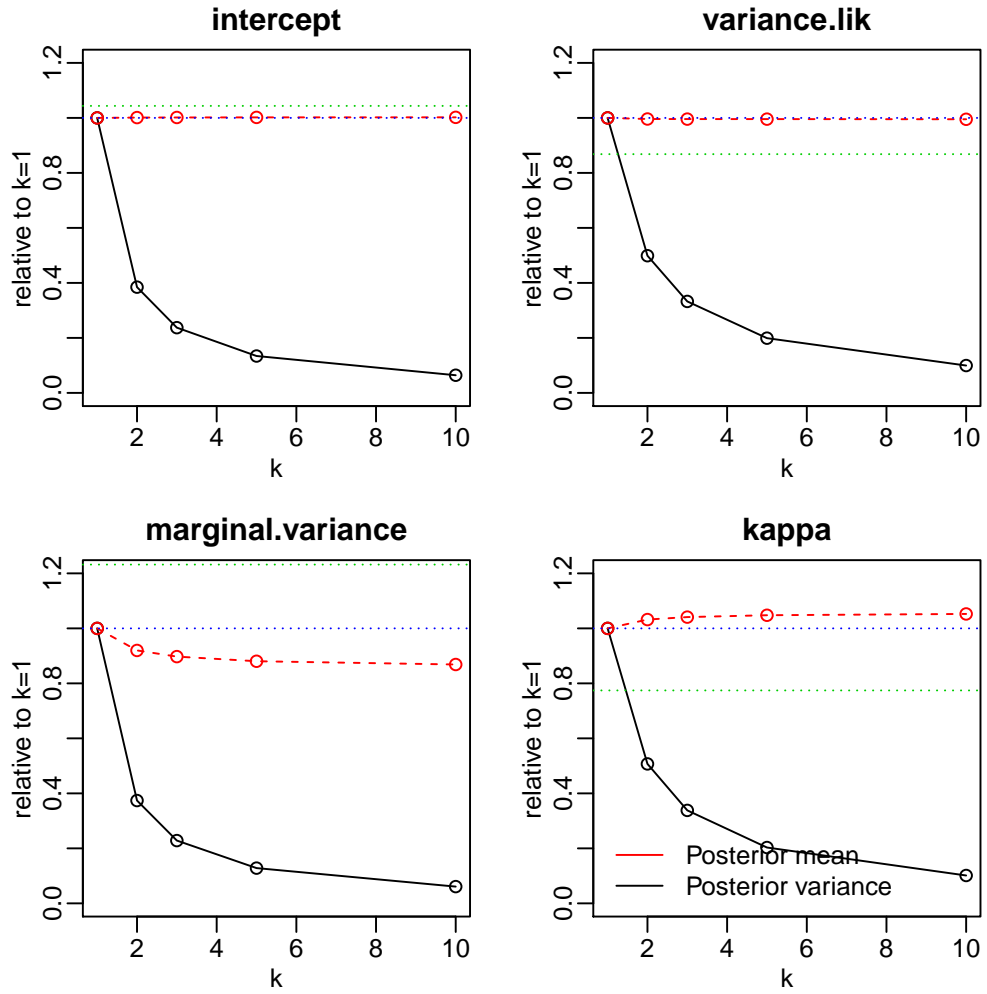


Figure 11.1: Relative posterior mean (blue dashed line), likelihood estimate (green dashed line) and posterior variance (black line) for each parameter of the toy model.

# Bibliography

- [Abrahamsen, 1997] Abrahamsen, P. (1997). A review of gaussian random fields and correlation functions. Norwegian Computing Center report No. 917.
- [Besag, 1981] Besag, J. (1981). On a system of two-dimensional recurrence equations. *J. R. Statist. Soc. B*, 43(3):302–309.
- [Bivand and Rundel, 2013] Bivand, R. and Rundel, C. (2013). *rgeos: Interface to Geometry Engine - Open Source (GEOS)*. R package version 0.2-13.
- [Bivand et al., 2012] Bivand, R., with contributions by Micah Altman, Anselin, L., ao, R. A., Berke, O., Bernat, A., Blanchet, G., Blankmeyer, E., Carvalho, M., Christensen, B., Chun, Y., Dormann, C., Dray, S., Halbersma, R., Krainski, E., Legendre, P., Lewin-Koh, N., Li, H., Ma, J., Millo, G., Mueller, W., Ono, H., Peres-Neto, P., Piras, G., Reder, M., Tiefelsdorf, M., and Yu, D. (2012). *spdep: Spatial dependence: weighting schemes, statistics and models*. R package version 0.5-55.
- [Bivand et al., 2008] Bivand, R. S., Pebesma, E. J., and Gomez-Rubio, V. (2008). *Applied spatial data analysis with R*. Springer, NY.
- [Cameletti et al., 2012] Cameletti, M., Lindgren, F., Simpson, D., and Rue, H. (2012). Spatio-temporal modeling of particulate matter concentration through the spde approach. *Advances in Statistical Analysis*.
- [Cressie, 1993] Cressie, N. (1993). *Statistics for Spatial Data*. Wiley, N. Y. 990p.
- [Diggle et al., 2010] Diggle, P. J., Menezes, R., and Su, T.-l. (2010). Geostatistical inference under preferential sampling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 59(2):191–232.
- [Diggle and Ribeiro Jr, 2007] Diggle, P. J. and Ribeiro Jr, P. J. (2007). *Model-Based Geostatistics*. Springer Series in Statistics. Hardcover. 230p.
- [Illian et al., 2012] Illian, J. B., Srbye, S. H., and Rue, H. (2012). A toolbox for fitting complex spatial point process models using integrated nested laplace approximation (inla). *Annals of Applied Statistics*, 6(4):1499–1530.
- [Lindgren, 2012] Lindgren, F. (2012). Continuous domain spatial models in r-inla. *The ISBA Bulletin*, 19(4). URL: <http://www.r-inla.org/examples/tutorials/spde-from-the-isba-bulletin>.
- [Lindgren and Rue, 2013] Lindgren, F. and Rue, H. (2013). Bayesian spatial and spatio-temporal modelling with r-inla. *Journal of Statistical Software*.
- [Lindgren et al., 2011] Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between gaussian fields and gaussian markov random fields: the stochastic partial differential equation approach. *J. R. Statist. Soc. B*, 73(4):423–498.
- [Muff et al., 2013] Muff, S., Riebler, A., Rue, H., Saner, P., and Held, L. (2013). Measurement error in glmms with inla. *submitted*.

- [Pebesma and Bivand, 2005] Pebesma, E. J. and Bivand, R. S. (2005). Classes and methods for spatial data in R. *R News*, 5(2):9–13.
- [Ribeiro Jr and Diggle, 2001] Ribeiro Jr, P. J. and Diggle, P. J. (2001). geoR: a package for geostatistical analysis. *R-NEWS*, 1(2):14–18. ISSN 1609-3631.
- [Rue and Held, 2005] Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*. Monographs on Statistics & Applied Probability. Boca Raton: Chapman and Hall.
- [Rue et al., 2009] Rue, H., Martino, S., and Chopin, N. (2009). Approximate bayesian inference for latent gaussian models using integrated nested laplace approximations (with discussion). *Journal of the Royal Statistical Society, Series B*, 71(2):319–392.
- [Simpson et al., 2011] Simpson, D. P., Illian, J. B., Lindren, F., Srbye, S. H., and Rue, H. (2011). Going off grid: Computationally efficient inference for log-gaussian cox processes. *submitted*.
- [Tobler, 1970] Tobler, W. R. (1970). A computer movie simulating urban growth in the detroid region. *Economic Geography*, 2(46):234–240.