

# Spatial Statistics (NMST543)

Zbyněk Pawlas

September 30, 2022

## Contents

<b>1</b>	<b>Point processes</b>	<b>3</b>
1.1	Estimation of summary characteristics . . . . .	3
1.1.1	Edge effects . . . . .	3
1.1.2	Estimation of intensity function . . . . .	4
1.1.3	Estimation of $K$ -function . . . . .	5
1.1.4	Estimation of inhomogeneous $K$ -function . . . . .	7
1.1.5	Estimation of pair correlation function . . . . .	7
1.1.6	Estimation of nearest neighbour distance distribution function . . . . .	7
1.1.7	Estimation of contact distribution function . . . . .	9
1.1.8	Estimation of $J$ -function . . . . .	10
1.1.9	Estimation of aggregation index . . . . .	10
1.2	Hypothesis testing . . . . .	11
1.3	Estimation of model parameters . . . . .	12
1.3.1	Method of moments . . . . .	12
1.3.2	Minimum contrast method . . . . .	13
1.3.3	Maximum likelihood method . . . . .	14
1.3.4	Maximum pseudolikelihood method . . . . .	14
1.3.5	Second-order composite likelihood . . . . .	16
1.3.6	Palm likelihood . . . . .	16
1.3.7	Takacs–Fiksel method . . . . .	17
1.4	Model diagnostics . . . . .	19
<b>2</b>	<b>Marked point processes</b>	<b>21</b>
2.1	Estimation of summary characteristics . . . . .	21
2.2	Tests of independence . . . . .	23
2.2.1	Testing independent marking . . . . .	23
2.2.2	Independence of marks and locations . . . . .	24
<b>3</b>	<b>Geostatistics</b>	<b>25</b>
3.1	Variogram estimation . . . . .	25
3.1.1	Non-parametric estimators . . . . .	25
3.1.2	Parametric methods . . . . .	27
3.1.3	Model validation . . . . .	29
3.2	Kriging . . . . .	29
3.2.1	Simple kriging . . . . .	29
3.2.2	Ordinary kriging . . . . .	30
3.2.3	Universal kriging . . . . .	31
3.2.4	Other possibilities . . . . .	32
3.3	Influence of covariance parameters estimation . . . . .	33
3.4	Bayesian approach . . . . .	34

<b>4</b>	<b>Lattice data</b>	<b>37</b>
4.1	Modelling and estimation for areal data . . . . .	37
4.2	Testing of spatial autocorrelation . . . . .	38
<b>5</b>	<b>Appendix</b>	<b>39</b>
5.1	Random censoring . . . . .	39

# 1 Point processes

In this chapter, we will deal with the statistical analysis of simple point processes on  $\mathbb{R}^d$ . We will start with the estimation of summary characteristics. Afterwards, the task of testing various statistical hypotheses will be discussed. Finally, we will consider parametric models and discuss the problem of model fitting and diagnostics.

Recall that a spatial point process is defined as a random locally finite counting measure on  $\mathbb{R}^d$ . A simple point process can also be viewed as a random locally finite set. We will use both of these approaches. Thus, we write  $\Phi(B)$  for the number of points (atoms) of the process  $\Phi$  in the set  $B$ . By  $X \in \Phi$  we mean that  $X$  is a point (atom) of  $\Phi$ .

## 1.1 Estimation of summary characteristics

Assume that we have a single realization of  $\Phi$  in a bounded Borel set  $W$ , so-called *observation window*. The window is usually a  $d$ -dimensional rectangle, but its shape may be more complicated. Our aim is to estimate various summary characteristics of the point process  $\Phi$  based on the given realization in  $W$ . A list of basic estimates follows. Most of them are implemented in the R package `spatstat` [1]. Therefore, we always mention also the corresponding R function.

### 1.1.1 Edge effects

When estimating numerical and functional summary characteristics, *edge effects* play an important role. They are caused by the fact that a point process is observed in a bounded window  $W$ . For example, we can base the estimate of  $K$ -function  $K(r)$  on the number of points in balls of radius  $r$  centred at the points of the process. However, we are unable to determine this number when the distance of the point to the boundary of  $W$  is smaller than  $r$ . The situation is illustrated in Figure 1 (left) – true number of points in  $b(X, r)$  is 5 but in  $W$  we only observe 3 points. For another example, we can consider the estimation of  $G$ -function. We are looking for the nearest neighbour of the point  $X \in \Phi$ . Based on the point pattern inside the observation window  $W$ , we would determine  $Y$  to be the nearest neighbour of  $X$ , see Figure 1 (right). Actually, the point  $Z$  that lies outside  $W$  is the true nearest neighbour of  $X$ . It is clear that our conclusion about the characteristics of the point process could be distorted by ignoring edge effects.

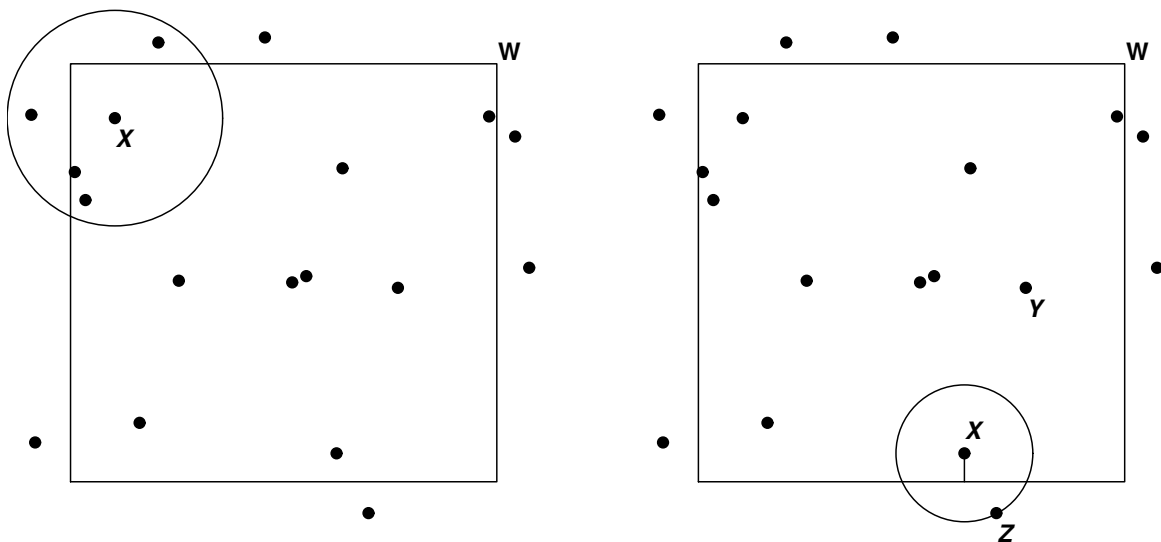


Figure 1: An illustration of edge effects in the case of estimating  $K$ -function (left) and  $G$ -function (right).

### 1.1.2 Estimation of intensity function

Let  $\Phi$  be a stationary point process on  $\mathbb{R}^d$  with intensity  $\lambda$ . It follows directly from the definition that

$$\hat{\lambda} = \frac{\Phi(W)}{|W|}$$

is an unbiased estimator of  $\lambda$ . In the package `spatstat`, this estimator can be computed using `summary.ppp`. If  $\Phi$  is a homogeneous Poisson point process,  $\hat{\lambda}$  is moreover the maximum likelihood estimator. In fact, we may understand  $\Phi$  on  $W$  as a finite point process with density w.r.t. the distribution of unit Poisson process on  $W$  and we know that its density has the form

$$p_\lambda(\varphi) = \lambda^{\varphi(W)} e^{(1-\lambda)|W|}.$$

It is easy to verify that the likelihood function  $L(\lambda) = p_\lambda(\varphi)$  is maximized for  $\lambda = \varphi(W)/|W|$  (see Exercise class).

For a non-stationary point process, non-parametric *kernel* estimators of its intensity function are often used. One possibility is to consider the estimator (`density.ppp`)

$$\hat{\lambda}(x) = \frac{1}{c_{W,b}(x)} \sum_{Y \in \Phi \cap W} k_b(x - Y), \quad x \in W,$$

where  $k_b$  is a kernel function with *bandwidth*  $b > 0$ , i.e.  $k_b(x) = \frac{k(x/b)}{b^d}$  for some probability density  $k$ , and

$$c_{W,b}(x) = \int_W k_b(x - y) dy$$

is the *edge correction factor*. Another possibility is to use a more exact but computationally more demanding estimator (`density.ppp` with `diggle=TRUE`)

$$\hat{\lambda}(x) = \sum_{Y \in \Phi \cap W} \frac{k_b(x - Y)}{c_{W,b}(Y)}, \quad x \in W. \quad (1)$$

The quality of the estimator  $\hat{\lambda}(x)$  is usually sensitive to the choice of bandwidth while the choice of kernel function doesn't play such an important role. For small values of  $b$ , the estimator is too concentrated around the observed points. On the other hand, larger  $b$  leads to oversmoothing. Density functions of the uniform distribution on a ball or the Gaussian distribution are one of the most common choices for the kernel function  $k$ . Also,  $k$  is often chosen as the product of one-dimensional densities:  $k(x) = k_1(x_1) \cdots k_d(x_d)$  for  $x = (x_1, \dots, x_d) \in \mathbb{R}^d$ . A popular example of a one-dimensional kernel function is the *Epanechnikov kernel*:

$$e(u) = \frac{3}{4}(1 - u^2), \quad u \in [-1, 1].$$

Note that if  $k$  is symmetric, then (1) is globally unbiased in the sense that

$$\int_W \hat{\lambda}(x) dx = \Phi(W),$$

i.e.

$$\mathbb{E} \int_W \hat{\lambda}(x) dx = \int_W \lambda(x) dx.$$

For a non-stationary Poisson process, we can express the likelihood function using

$$p_\theta(\varphi) = \exp\left\{|W| - \int_W \lambda_\theta(x) dx\right\} \prod_{x \in \varphi} \lambda_\theta(x).$$

If the intensity function  $\lambda_\theta(x)$  has a parametric form, then the task of finding the maximum likelihood estimator of parameter  $\theta$  has to be solved by some numerical method. We will return to this problem in Subsection 1.3.

### 1.1.3 Estimation of $K$ -function

Recall that the  $K$ -function  $K(r)$  of a stationary point process  $\Phi$  is defined through the relation

$$\lambda K(r) = \mathbb{E}_o^! \Phi(b(o, r)), \quad r > 0,$$

where  $\mathbb{E}_o^!$  is the expectation w.r.t. the reduced Palm distribution at the origin  $o$ . This definition could be equivalently rewritten as

$$\lambda K(r) = \mathbb{E} \sum_{X \in \Phi \cap A} \frac{\Phi(b(X, r) \setminus \{X\})}{\lambda |A|} = \mathbb{E} \sum_{X, Y \in \Phi}^{\neq} \frac{\mathbf{1}_{[X \in A, \|X-Y\| \leq r]}}{\lambda |A|}, \quad (2)$$

where  $A$  is an arbitrary Borel set with finite and positive Lebesgue measure ( $0 < |A| < \infty$ ).

The following estimators can be obtained by the function `Kest` in the package `spatstat`.

*0. uncorrected estimate:* The equation (2) offers a theoretical unbiased estimator of  $\lambda^2 K(r)$  in the form

$$\widehat{\lambda^2 K(r)} = \sum_{X \in \Phi \cap W} \frac{\Phi(b(X, r) \setminus \{X\})}{|W|} = \sum_{X, Y \in \Phi}^{\neq} \frac{\mathbf{1}_{[X \in W, \|X-Y\| \leq r]}}{|W|}.$$

This estimator could only be used if we have additional information about the point pattern outside the window  $W$ . This is because we put no restriction on points  $Y$  except that they have to be closer than the given distance  $r$  to some point of  $\Phi$  lying in  $W$ . For one particular point  $X$ , the distance  $r$  may be larger than the distance to the boundary of  $W$  and thus we need to count points lying outside  $W$ . This situation is known as the *plus sampling*. The problem of the inapplicability of the estimate rests in edge effects. We are unable to determine  $\Phi(b(X, r) \setminus \{X\})$  only from the information about points lying inside the window  $W$ , see Figure 1 (left). If we ignore the edge effects and consider only the points inside  $W$ , we get a negatively biased estimator

$$\lambda^2 K_u(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W|}.$$

The package `spatstat` enables to calculate this estimator by setting `correction="none"` in `Kest`. However, it is only for instructional reasons. This estimator should not be used in practice. All the following estimators try to compensate for the problem of edge effects by including the edge correction factor  $e_{W,r}(X, Y)$ . The estimators of  $\lambda^2 K(r)$  then have the form

$$\sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W|} e_{W,r}(X, Y).$$

For  $\widehat{\lambda^2 K_u(r)}$  there is no edge correction and the factor  $e_{W,r}(X, Y)$  is identically equal to 1.

*1. border estimate, correction="border":* The simplest way to avoid the edge effects is to consider  $\Phi$  in a smaller window

$$W_{\ominus r} = W \ominus b(o, r) = \{y \in W : b(y, r) \subseteq W\} = \{y \in W : d(y, \partial W) \geq r\},$$

where  $\partial W$  denotes the boundary of  $W$ . We only consider the points for which we are able to determine the number of neighbours at distance  $r$ . This procedure is known as the *minus sampling*. Then

$$\widehat{\lambda^2 K_b(r)} = \sum_{X \in \Phi \cap W_{\ominus r}} \frac{\Phi(b(X, r) \setminus \{X\})}{|W_{\ominus r}|} = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[X \in W_{\ominus r}, \|X-Y\| \leq r]}}{|W_{\ominus r}|}$$

is an unbiased estimator  $\lambda^2 K(r)$ , as it follows from (2). The edge correction factor is

$$e_{W,r}(X, Y) = \frac{\mathbf{1}_{[X \in W_{\ominus r}]} |W|}{|W_{\ominus r}|}.$$

The estimator  $\widehat{\lambda^2 K_b}(r)$  is defined for  $r < r_b = \sup\{s > 0 : |W_{\ominus s}| > 0\}$ . For example,  $r_b = 0.5$  for  $W = [0, 1]^2$ .

2. *translation correction, correction="translate"*: Another possibility is to let the edge correction factor be the function of both  $X$  and  $Y$ . The *translation correction factor* is

$$e_{W,r}(X, Y) = \frac{|W|}{|W \cap (W + X - Y)|},$$

which leads to the estimator

$$\widehat{\lambda^2 K_t}(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W \cap (W + X - Y)|}.$$

Using the Campbell theorem it can be shown that this estimator is unbiased (see Exercise class). The estimator is well defined for  $r < r_t = \sup\{s > 0 : |W \cap (W + x)| > 0 \forall x : \|x\| \leq s\}$ . For example,  $r_t = 1$  for  $W = [0, 1]^2$ . Similarly, we can define the estimator of the reduced second-order moment measure  $\mathcal{K}$  using

$$\widehat{\lambda^2 \mathcal{K}_t}(B) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[X-Y \in B]}}{|W \cap (W + X - Y)|}.$$

The smoothed kernel estimate of the density of  $\mathcal{K}$  can be obtained by `Kmeasure`.

3. *Ripley's isotropic correction, correction="isotropic" or correction="Ripley"*: Another correction factor was suggested by B. D. Ripley [12]. It has the form

$$e_{W,r}(X, Y) = \frac{|\partial b(X, \|X - Y\|)|}{|\partial b(X, \|X - Y\|) \cap W|}$$

and yields the estimator

$$\widehat{\lambda^2 K_R}(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W|} \cdot \frac{|\partial b(X, \|X - Y\|)|}{|\partial b(X, \|X - Y\|) \cap W|}.$$

The notation  $|\partial b(x, r)|$  means the  $(d-1)$ -Hausdorff measure of  $\partial b(x, r)$ . If the process is motion-invariant, it can be shown that  $\widehat{\lambda^2 K_R}(r)$  is an unbiased estimator of  $\lambda^2 K(r)$  for  $r < r_0 = \inf\{t > 0 : |W^{(t)}| < |W|\}$ , where  $W^{(t)} = \{x \in W : \partial b(x, t) \cap W \neq \emptyset\}$ . Ohser's modification [9] is given by

$$\widehat{\lambda^2 K_O}(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W^{(\|X-Y\|)}|} \cdot \frac{|\partial b(X, \|X - Y\|)|}{|\partial b(X, \|X - Y\|) \cap W|}.$$

It extends the definition of  $\widehat{\lambda^2 K_R}(r)$  to  $r < r^* = \sup\{s > 0 : |W^{(s)}| > 0\}$ . For  $r < r_0$  we have  $\widehat{\lambda^2 K_R}(r) = \widehat{\lambda^2 K_O}(r)$ . In case of the planar unit square  $W = [0, 1]^2$ ,  $r_0 = \sqrt{2}/2$ ,  $r^* = \sqrt{2}$ , and  $W^{(s)} = W$  for all  $s \leq r_0$ .

In order to get the estimator of  $K(r)$  itself, we have to divide the estimator of  $\lambda^2 K(r)$  by the estimator of  $\lambda^2$ . This violates the unbiasedness property. The bias and variance are typically increasing with increasing  $r$ . For a planar rectangular window, it is recommended to determine the estimators only for  $r$  smaller than a  $1/4$  of the shorter side length of the rectangle. The estimator of  $\lambda^2$  often has the following form

$$\widehat{\lambda^2} = \frac{\Phi(W)(\Phi(W) - 1)}{|W|^2}.$$

The motivation is that  $\widehat{\lambda^2}$  is an unbiased estimator of  $\lambda^2$  in the case of Poisson point process  $\Phi$ .

The border estimator of the  $K$ -function does not have to be a monotone function in  $r$  (as opposed to the theoretical function). With increasing  $r$  and dimension  $d$ , the loss of information could be essential. The estimators based on the translation or isotropic correction factors have statistically better properties. On the other hand, the computation of  $\widehat{K}_b$  is faster.

### 1.1.4 Estimation of inhomogeneous $K$ -function

Let  $\Phi$  be a second-order intensity reweighted stationary point process. Then we define the inhomogeneous  $K$ -function as

$$K_{\text{inhom}}(r) = \mathbb{E} \sum_{X, Y \in \Phi}^{\neq} \frac{\mathbf{1}_{[X \in A, \|X - Y\| \leq r]}}{\lambda(X)\lambda(Y)|A|},$$

where  $A$  is an arbitrary Borel set with finite and positive Lebesgue measure ( $0 < |A| < \infty$ ). Its estimators could be obtained similarly as in the case of stationary processes. For example, the estimator with translation correction has the form

$$\widehat{K}_{\text{inhom}}(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X - Y\| \leq r]}}{\widehat{\lambda}(X)\widehat{\lambda}(Y)|W \cap (W + X - Y)|},$$

where  $\widehat{\lambda}(x)$  is the estimator of intensity function  $\lambda(x)$ . In the package `spatstat`, we would use `Kinhom` with the choice `correction="translate"`.

### 1.1.5 Estimation of pair correlation function

For a motion-invariant point process, the pair correlation function  $g$  is related to the  $K$ -function by

$$g(r) = \frac{K'(r)}{\sigma_d r^{d-1}}, \quad r > 0,$$

where  $\sigma_d = |\partial b(o, 1)|$  is the surface area of the unit ball in  $\mathbb{R}^d$ . The kernel estimator of  $g$  is

$$\widehat{g}(r) = \frac{1}{\widehat{\lambda}^2} \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{k_b(r - \|X - Y\|)}{\sigma_d r^{d-1} |W|} e_{W,r}(X, Y),$$

where  $k_b$  is a suitable kernel function with bandwidth  $b$  and  $e_{W,r}(X, Y)$  is the edge correction factor. The estimator is obtained by calling the function `pcf` in `spatstat`. The choices `correction="translate"` and `correction="ripley"` correspond to the translation and isotropic edge correction, respectively. In the case of a second-order intensity reweighted stationary point process,  $\widehat{\lambda}^2$  is replaced by  $\widehat{\lambda}(X)\widehat{\lambda}(Y)$  in the denominator of each summand. The computation could be performed by the function `pcfinhom`.

Another possibility would be to use some estimator of the  $K$ -function and approximate its derivative by numerical methods (e.g., using splines). This is usually demanding because the estimators of  $K$ -function are piecewise constant functions.

### 1.1.6 Estimation of nearest neighbour distance distribution function

Recall that for a stationary point process we define the nearest neighbour distance distribution function as

$$G(r) = P_o^!(\{\varphi \in \mathcal{N} : \varphi(b(o, r)) > 0\}), \quad r > 0,$$

where  $P_o^!$  is the reduced Palm distribution.

For the computation of the following estimators of  $G$  in  $\mathbb{R}$  we can use the function `Gest`.

*0. uncorrected estimate, correction="none"*: If we know the nearest neighbour distance for each observed point of the process, we can estimate  $G$  classically by the empirical distribution function

$$\widehat{G}(r) = \frac{1}{\Phi(W)} \sum_{X \in \Phi \cap W} \mathbf{1}_{[e(X) \leq r]},$$

where  $e(x) = d(x, \Phi \setminus \{x\})$  is the distance from the point  $x$  to its nearest neighbour. From the Campbell-Mecke theorem it follows that

$$\mathbb{E} \sum_{X \in \Phi \cap W} \mathbf{1}_{[e(X) \leq r]} = \lambda \int_W \int_{\mathcal{N}} \mathbf{1}_{[d(o, \varphi) \leq r]} P_o^!(d\varphi) dx = \lambda |W| G(r).$$

Hence, we see that  $\hat{G}(r)$  is a so-called *ratio-unbiased* estimator of  $G$ . It means that the ratio of expectations of the numerator and denominator of  $\hat{G}(r)$  gives  $G(r)$ , i.e.

$$\frac{\mathbb{E} \sum_{X \in \Phi \cap W} \mathbf{1}_{[e(X) \leq r]}}{\mathbb{E} \Phi(W)} = \frac{\lambda |W| G(r)}{\lambda |W|} = G(r).$$

Again, we are unable to get  $e(X)$  for each  $X \in \Phi \cap W$  due to the edge effects, see Figure 1 (right). We can replace  $e(X)$  by the distance  $e^*(X) = d(X, (\Phi \setminus \{X\}) \cap W) \geq e(X)$ , which we are able to observe in the observation window. Then we obtain the naive estimator

$$\hat{G}_u(r) = \frac{1}{\Phi(W)} \sum_{X \in \Phi \cap W} \mathbf{1}_{[e^*(X) \leq r]}.$$

Since  $e^*(X) \leq r$  implies  $e(X) \leq r$ , we have  $\hat{G}_u(r) \leq \hat{G}(r)$ . The estimator  $\hat{G}_u(r)$  is not used for practical purposes. However, it can be obtained by the choice `correction="none"`.

1. *border estimate*, `correction="border"` or `"rs"`: By restricting to the points  $X$  lying in the eroded window  $W_{\ominus r}$ , the following ratio-unbiased estimator is obtained,

$$\hat{G}_b(r) = \frac{1}{\Phi(W_{\ominus r})} \sum_{X \in \Phi \cap W_{\ominus r}} \mathbf{1}_{[e(X) \leq r]} = \frac{1}{\Phi(W_{\ominus r})} \sum_{X \in \Phi \cap W_{\ominus r}} \mathbf{1}_{[e^*(X) \leq r]}.$$

2. *Kaplan–Meier estimate*, `correction="km"`: Edge effects could be understood as a type of censoring (see Subsection 5.1). Therefore, we can introduce the Kaplan–Meier type estimator,

$$\hat{G}_{KM}(r) = 1 - \prod_{s \leq r} \left( 1 - \frac{\#\{X \in \Phi \cap W : e(X) = s, e(X) \leq c(X)\}}{\#\{X \in \Phi \cap W : e(X) \geq s, c(X) \geq s\}} \right),$$

where  $c(x) = d(x, \partial W)$  is the distance from  $x$  to the window boundary and  $\#$  stands for the number of elements. If  $e(X) \leq c(X)$ , we are sure that we observe true distance to the nearest neighbour of  $x$ . In the opposite case, the distance  $e(X)$  is censored by  $c(X)$ . We only know that  $e(X)$  is larger than  $c(X)$ . However, we don't have information about the exact value of  $e(X)$ . Note that the information contained in the observation window  $W$  is sufficient for evaluating the estimator  $\hat{G}_{KM}(r)$ . As opposed to the classical situation of random censoring, we may not expect the independence of data and censors. Hence, the optimality of the Kaplan–Meier estimator is violated. Nevertheless, it usually gives better results than the border estimator.

For an absolutely continuous distribution function  $H(t)$  with density  $h(t)$ , the *hazard rate* is defined as  $\lambda_h(t) = h(t)/(1 - H(t))$ . The spatial Kaplan–Meier method enables us to estimate the hazard rate of the distribution function  $G(r)$ . However, we have to be cautious because  $G$  does not necessarily have a density. In `spatstat`, this estimator can be obtained together with the Kaplan–Meier estimator of  $G$ -function.

3. *Hanisch estimate*, `correction="han"` or `"Hanisch"`: Another improvement of the border estimate is obtained by the following edge correction:

$$\hat{G}_H(r) = \frac{1}{\hat{\lambda}} \sum_{X \in \Phi \cap W} \frac{\mathbf{1}_{[e(X) \leq c(X)]}}{|W_{\ominus e(X)}|} \mathbf{1}_{[e(X) \leq r]},$$

where

$$\hat{\lambda} = \sum_{X \in \Phi \cap W} \frac{\mathbf{1}_{[e(X) \leq c(X)]}}{|W_{\ominus e(X)}|}.$$

Note that only points with a smaller distance to their neighbour than to the boundary of the observation window are used to compute this estimator. In Figure 1 (right), point  $X$  is not involved in the estimator because its distance to the boundary of  $W$  is smaller than the distance to its nearest neighbour. The Hanisch estimator is ratio-unbiased as we can easily verify by the Campbell–Mecke theorem if we realize that  $\mathbf{1}_{[e(X) \leq c(X)]} = \mathbf{1}_{[X \in W_{\ominus e(X)}]}$ .

It is important to realize that mentioned estimators of  $G$  may not have the properties of a distribution function:  $\hat{G}_b$  is not necessarily monotone,  $\hat{G}_{KM}$  is non-decreasing but its maximal value could be strictly smaller than 1.



### 1.1.7 Estimation of contact distribution function

The contact distribution function of a stationary point process  $\Phi$  is defined as

$$F(r) = \mathbb{P}(\Phi(b(o, r)) > 0) = \mathbb{P}(D \leq r), \quad r > 0,$$

where  $D = d(o, \Phi)$  is the distance from the origin to the nearest point of  $\Phi$ .

Let us choose a regular grid  $I_a$  in  $\mathbb{R}^d$ :

$$I_a = y + a\mathbb{Z}^d = \{(y_1 + a_1 z_1, \dots, y_d + a_d z_d) \in \mathbb{R}^d : z_i \in \mathbb{Z}\},$$

where  $y = (y_1, \dots, y_d) \in \mathbb{R}^d$  and  $a = (a_1, \dots, a_d) \in \mathbb{R}_+^d$ , i.e.  $a_i > 0$  for  $i = 1, \dots, d$ . The R function `Fest` can be used for the calculation of the following estimators of  $F$ .

0. *uncorrected estimate, correction="none"*: For every point of the grid in the window  $W$ , we find the nearest point of the process. However, this nearest point may lie outside  $W$ . If we only consider points of the process  $\Phi$  that lie in  $W$ , we get

$$\hat{F}_u(r) = \frac{1}{|I_a \cap W|} \sum_{x \in I_a \cap W} \mathbf{1}_{[d(x, \Phi \cap W) \leq r]},$$

where  $|I_a \cap W|$  is the number of points of a finite set  $I_a \cap W$ . This estimator is negatively biased, i.e.  $\mathbb{E}\hat{F}_u(r) \leq F(r)$ , because  $\mathbf{1}_{[d(x, \Phi \cap W) \leq r]} \leq \mathbf{1}_{[d(x, \Phi) \leq r]}$  and  $\mathbb{P}(d(x, \Phi) \leq r) = F(r)$  from stationarity. The bias is caused by edge effects.

1. *border estimate, correction="border" or "rs"*:

$$\hat{F}_b(r) = \frac{1}{|I_a \cap W_{\ominus r}|} \sum_{x \in I_a \cap W_{\ominus r}} \mathbf{1}_{[d(x, \Phi) \leq r]}$$

is an unbiased estimator of  $F(r)$  because  $\mathbb{P}(d(x, \Phi) \leq r) = F(r)$  by stationarity. The continuous version of this estimator (as  $a \rightarrow 0$ ) has the form

$$\hat{F}_b(r) = \frac{|W_{\ominus r} \cap \Phi_r|}{|W_{\ominus r}|},$$

where  $\Phi_r = \{x \in \mathbb{R}^d : d(x, \Phi) \leq r\} = \cup_{X \in \Phi} b(X, r)$ . Again it is an unbiased estimator.

2. *Kaplan–Meier estimate, correction="km"*: Let  $d(x) = d(x, \Phi)$  be the distance from  $x$  to the nearest point of  $\Phi$ . Then we define

$$\hat{F}_{KM}(r) = 1 - \prod_{s \leq r} \left( 1 - \frac{\#\{x \in I_a \cap W : d(x) = s, d(x) \leq c(x)\}}{\#\{x \in I_a \cap W : d(x) \geq s, c(x) \geq s\}} \right),$$

where  $c(x) = d(x, \partial W)$  is the distance from  $x$  to the boundary of  $W$ . The contact distribution function  $F(r)$  of a stationary point process is absolutely continuous and the hazard rate  $\lambda_h(r)$  exists. Its estimator is based on the Kaplan–Meier estimator  $\hat{F}_{KM}(r)$ .

3. *Chiu–Stoyan estimate, correction="cs" or "Hanisch"*: Using the same correction as in the Hanisch estimator of  $G$  we obtain

$$\hat{F}_{CS}(r) = \frac{1}{C_a} \sum_{x \in I_a \cap W} \frac{\mathbf{1}_{[d(x) \leq c(x)]}}{|W_{\ominus d(x)}|} \mathbf{1}_{[d(x) \leq r]},$$

where

$$C_a = \sum_{x \in I_a \cap W} \frac{\mathbf{1}_{[d(x) \leq c(x)]}}{|W_{\ominus d(x)}|}.$$

The continuous version of this estimator is

$$\hat{F}_{CS}(r) = \frac{1}{C} \int_W \frac{\mathbf{1}_{[d(x) \leq c(x)]}}{|W_{\ominus d(x)}|} \mathbf{1}_{[d(x) \leq r]} dx,$$

where

$$C = \int_W \frac{\mathbf{1}_{[d(x) \leq c(x)]}}{|W_{\ominus d(x)}|} dx.$$

We emphasize that the estimators of  $F$  may not have the properties of a distribution function:  $\hat{F}_b$  is not necessarily monotone,  $\hat{F}_{KM}$  is non-decreasing but its maximal value could be strictly smaller than 1. The border estimator is less efficient than the Kaplan–Meier estimator or the Chiu–Stoyan estimator.

### 1.1.8 Estimation of $J$ -function

In `spatstat` it is possible to estimate the  $J$ -function

$$J(r) = \frac{1 - G(r)}{1 - F(r)}, \quad r > 0 : F(r) < 1,$$

using the function `Jest`.

The estimator of  $J$  arises from its definition:

$$\hat{J}(r) = \frac{1 - \hat{G}(r)}{1 - \hat{F}(r)}.$$

We may distinguish the following estimators (depending on the type of estimator of  $G$  and  $F$ ):

- *uncorrected* (correction="none"),
- *border* (correction="border" or "rs"),
- *Kaplan–Meier* (correction="km"),
- *Hanisch* (correction="Hanisch").

Even if the uncorrected estimators  $\hat{G}_u$  and  $\hat{F}_u$  are substantially biased, taking their ratio gives an approximately unbiased estimator (at least when the point process is close to the Poisson process). The advantage of this estimator is that it is insensitive to edge effects. Therefore, it could be used when the edge effects are severe. The other three estimators are slightly biased (ratio of two approximately unbiased estimators). The logarithm of the Kaplan–Meier estimator is an unbiased estimator of  $\log J$ .

The package `spatstat` enables us to estimate four basic summary characteristics (functions  $F$ ,  $G$ ,  $J$ ,  $K$ ) at the same time by `allstats`.

### 1.1.9 Estimation of aggregation index

The expectation of an arbitrary non-negative random variable  $T$  can be expressed using its distribution function  $H(t)$  as follows (e.g., [6], Lemma 5.7)

$$\mathbb{E}T = \int_0^\infty (1 - H(t)) dt.$$

Having the estimator  $\hat{G}(t)$  of the nearest neighbour distance distribution function  $G(t)$ , we can estimate the Clark–Evans index

$$\text{CE} = \frac{d(\lambda\omega_d)^{1/d}}{\Gamma(1/d)} \mathbb{E}_o^! D$$

as

$$\widehat{\text{CE}} = \frac{d(\hat{\lambda}\omega_d)^{1/d}}{\Gamma(1/d)} \int_0^\infty (1 - \hat{G}(t)) dt.$$

In `spatstat` the function `clarkevans` can be used for the estimation of CE.

## 1.2 Hypothesis testing

The next statistical task consists of testing the hypothesis that the observed point pattern corresponds to a given point process model. The most important case is testing the hypothesis of complete spatial randomness. If we do not reject this hypothesis, then the observed data can be modelled by a Poisson point process, and it is unnecessary to consider more complicated processes. A complete spatial randomness test is one of the basic steps of exploratory data analysis.

First, divide the observation window  $W$  to  $k$  mutually disjoint regions (so-called *quadrats*) of the same volume and count the number of points in each of these quadrats. Denote these counts by  $n_1, \dots, n_k$ . We consider the null hypothesis that the data corresponds to a homogeneous Poisson point process. Under this hypothesis, the counts should form a realization of a random sample from the Poisson distribution with parameter  $\lambda|W|/k$ . Moreover, all  $n = n_1 + \dots + n_k$  points are i.i.d. and have the uniform distribution in  $W$ . Hence, we can use the well-known Pearson's  $\chi^2$  goodness-of-fit test. The test statistic is given as

$$\sum_{i=1}^k \frac{(n_i - n/k)^2}{n/k}$$

and it is equal to the *dispersion index*

$$I = \frac{(k-1)s^2}{\bar{n}},$$

where  $\bar{n} = \frac{1}{k} \sum_{i=1}^k n_i = n/k$  is the average number of points per quadrat and  $s^2 = \frac{1}{k-1} \sum_{i=1}^k (n_i - \bar{n})^2$  is the sample variance. The index  $I$  has approximately  $\chi^2$ -distribution with  $k-1$  degrees of freedom. We reject the null hypothesis if  $I \leq \chi_{k-1}^2(\alpha/2)$  or  $I \geq \chi_{k-1}^2(1-\alpha/2)$ , where  $\chi_{k-1}^2(p)$  is the  $p$ -quantile of  $\chi^2$ -distribution with  $k-1$  degrees of freedom. In order to get a reasonable approximation by the asymptotic  $\chi^2$ -distribution, the practical recommendation on the number of points of the investigated point pattern is  $\bar{n} > 5$ . Small values of  $I$  correspond to smaller variability than for the Poisson process. Thus, small values of  $I$  indicate regularity of the investigated point pattern. On the other hand, larger values  $I$  show bigger variability in the point pattern. This situation corresponds to the clustering of the observed point pattern. In *spatstat*, the test can be performed using `quadrat.test`.

The test based on the dispersion index is one of few cases in spatial statistics where the (asymptotic) distribution of a test statistic is tractable. In most situations, the test statistic has a very complicated distribution. Therefore, simulation (Monte Carlo) tests are used. First, we explain their general idea.

Suppose that we want to test the hypothesis  $H_0$  that the data corresponds to a given model. Consider a suitable test statistic  $T$ . Usually,  $T$  has unknown or intractable distribution. We perform  $M$  independent simulations from the null model  $H_0$  and determine the corresponding test statistics  $T_1, \dots, T_M$ . We rank them from the smallest to the largest and obtain the ordered sample  $T_{(1)} \leq \dots \leq T_{(M)}$ . Under the null hypothesis,  $T$  and  $T_1, \dots, T_M$  are i.i.d. and hence by symmetry every ranking has the same probability. For simplicity, assume that there are no ties in  $T, T_1, \dots, T_M$  almost surely. Then the probability that  $T$  is smaller than  $T_{(q)}$  equals  $q/(M+1)$ . We want to test  $H_0$  on the prescribed significance level  $\alpha$ . If we consider a two-sided test (both small and large values of the test statistic serve against the null hypothesis), we determine  $q$  such that

$$\alpha = \frac{2q}{M+1}.$$

We choose  $M$  so that  $\alpha(M+1)/2$  is an integer. The hypothesis is then rejected if  $T \notin [T_{(q)}, T_{(M-q+1)}]$ . This test is referred to as *Barnard's Monte Carlo test*. We can also determine the  $p$ -value of the test,  $p = 2 \min(p_+, p_-)$ , where

$$p_+ = \frac{1 + \sum_{j=1}^M \mathbf{1}\{T_j < T\}}{M+1} \quad \text{and} \quad p_- = \frac{1 + \sum_{j=1}^M \mathbf{1}\{T_j > T\}}{M+1}.$$

We reject  $H_0$  if  $p \leq \alpha$ .

In point processes, we rather work with functional than numerical characteristics. Let  $T(r)$  be some functional test statistic. Often,  $T(r) = \hat{S}(r)$  is an estimator of some summary characteristic  $S(r)$ . For fixed  $r_0$ , chosen in advance independently on data, we may carry out Barnard's Monte Carlo test with  $T = T(r_0)$ . However, then we use only part of the information given by  $T(r)$ .

Consider an estimator of the summary characteristic  $S(r)$  evaluated on the interval  $[r_1, r_2]$ , where  $0 \leq r_1 < r_2$  are prescribed real constants. Denote by  $\hat{S}(r)$  this estimator computed from data and by  $\hat{S}_1(r), \dots, \hat{S}_M(r)$  the estimators computed from  $M$  independent simulations. For each  $r \in [r_1, r_2]$ , we would be able to determine  $\hat{S}_{(q)}(r)$  and  $\hat{S}_{(M-q+1)}(r)$ . By joining the values  $\hat{S}_{(q)}(r)$  for different  $r$ , we get a so-called *pointwise lower envelope*, while values  $\hat{S}_{(M-q+1)}(r)$  form a *pointwise upper envelope*. In `spatstat`, we can draw these envelopes using the function `envelope` with parameter `global=FALSE`. They could be useful to reveal the deviations from the null hypothesis. However, we have to realize that it would be incorrect to use pointwise envelopes for testing. We are dealing with the problem of multiple testing. When the curve  $\hat{S}(r)$  reaches outside the envelopes for some  $r$ , it doesn't mean that  $H_0$  has to be rejected.

On the other hand, a *global envelope test* rejects the null hypothesis  $H_0$  if  $T(r) = \hat{S}(r)$  is not completely inside the envelope, i.e. if  $\hat{S}(r) \notin [T_{\text{low}}(r), T_{\text{upp}}(r)]$  for some  $r \in [r_1, r_2]$ . The functions  $T_{\text{low}}(r)$  and  $T_{\text{upp}}(r)$  define lower and upper envelopes, respectively. In order to get the exact envelope test, assume that we know the theoretical form of  $S(r) = S_0(r)$  under the null hypothesis. Let us determine the maximal absolute differences from the theoretical function:

$$D = \sup_{r_1 \leq r \leq r_2} |\hat{S}(r) - S_0(r)|, \quad D_i = \sup_{r_1 \leq r \leq r_2} |\hat{S}_i(r) - S_0(r)|, \quad i = 1, \dots, M.$$

Ordering the values  $D_1, \dots, D_M$ , we obtain rank statistics  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(M)}$ . Now we have a one-sided test, only large values of the deviations give evidence against the null hypothesis. We reject the null hypothesis if  $D > D_{(M-q+1)}$ , where  $q$  is chosen according to the required significance level,  $\alpha = \frac{q}{M+1}$ . The testing procedure could be represented in the following way. Construct a band of width  $2D_{(M-q+1)}$  around the function  $S_0(r)$ . If  $\hat{S}(r)$  lies outside of this band (outside envelopes) for some  $r \in [r_1, r_2]$ , the hypothesis is rejected. It means that we have  $T_{\text{low}}(r) = S_0(r) - D_{(M-q+1)}$  and  $T_{\text{upp}}(r) = S_0(r) + D_{(M-q+1)}$ . The probability that  $\hat{S}(r)$  wanders outside the envelopes is exactly  $\alpha$ . These envelopes are called *simultaneous*. In `spatstat`, they could be obtained by the function `envelope` with `global=TRUE`. The corresponding test is known as the *maximum absolute deviation test* and can be performed by the function `mad.test`.

Another possibility is to consider integral deviation measure instead of the supremum one. For data and for each simulation, we determine the integral square deviations from the theoretical function,

$$D = \int_{r_1}^{r_2} (\hat{S}(r) - S_0(r))^2 dr, \quad D_i = \int_{r_1}^{r_2} (\hat{S}_i(r) - S_0(r))^2 dr, \quad i = 1, \dots, M.$$

The null hypothesis is rejected if  $D > D_{(M-q+1)}$ , where  $q$  is chosen according to the required significance level,  $\alpha = \frac{q}{M+1}$ . In this case we speak about the *integral deviation test*. It can be performed by the function `dclf.test`.

The idea of Barnard's Monte Carlo test can be extended to the case where the test statistic  $T$  belongs to a more general space provided that ordering could be defined on this space. An example of the test obtained by ordering the functions is the *global rank envelope test* in [8].

To test complete spatial randomness, we can use any of the simulation tests. As a characteristic  $S(r)$  usually  $F$ ,  $G$ ,  $J$ ,  $K$  or  $L$ -function is used. In the same way, we can perform a goodness-of-fit test for an arbitrary model from which we are able to simulate.

### 1.3 Estimation of model parameters

Another statistical problem is to find a suitable model that describes our data well. Consider that the distribution of a point process  $\Phi$  is specified except for a vector  $\theta$  of unknown parameters. We write  $P_\theta$  for the corresponding distribution of  $\Phi$  and  $\mathbb{E}_\theta$  for the expectation with respect to  $P_\theta$ . Our aim is to find an estimate of  $\theta$  based on the realization of  $\Phi$  in a bounded window  $W \subseteq \mathbb{R}^d$ .

#### 1.3.1 Method of moments

We consider statistics  $T_1, \dots, T_k$ . Let  $t_{j,\theta} = \mathbb{E}_\theta T_j$  be the theoretical expected value of  $T_j$  for given  $\theta$ . Suppose that  $\theta$  is a  $k$ -dimensional vector and we have an analytic expression of  $t_{j,\theta}$  as a function of  $\theta$  for

each  $j = 1, \dots, k$ . Then the estimator  $\hat{\theta}$  by the method of moments is given as the solution (if it exists and is unique) of the equations

$$T_j = t_{j,\theta}, \quad j = 1, \dots, k.$$

That is, the observed values  $T_j$  coincide with the theoretical expected values.

*Example:* Let  $\Phi$  be a stationary point process with intensity  $\lambda$ . We choose  $T_1 = \Phi(W)$ . Then  $t_{1,\lambda} = \lambda|W|$  and we obtain the natural estimator  $\hat{\lambda} = \Phi(W)/|W|$ .

### 1.3.2 Minimum contrast method

For point processes, we often deal with functional summary characteristics. Let  $S(r)$  be a chosen functional characteristic and  $\hat{S}(r)$  be its estimator. Assume that the theoretical form of  $S(r)$  is known and can be expressed as a function of model parameters, say  $S_\theta(r)$ . Then we may consider the estimator evaluated for at least  $k$  distinct values  $r_1, \dots, r_k$  and apply the method of moments with  $T_j = \hat{S}(r_j)$  and  $t_{j,\theta} = S_\theta(r_j)$ ,  $j = 1, \dots, k$ . Here, the relation  $t_{j,\theta} = \mathbb{E}_\theta T_j$  holds only if the estimator  $\hat{S}(r)$  is unbiased.

The theoretical form of  $S(r)$  is known in several cases. Examples include statistics of a Poisson process or pair correlation function of a stationary Neyman-Scott process, which is given by

$$g(x) = 1 + \frac{1}{\lambda_p} \int_{\mathbb{R}^d} p(y)p(y-x) dy, \quad x \in \mathbb{R}^d,$$

where  $\lambda_p$  is the intensity of the parent point process and  $p$  is the density of the displacement vector of a daughter point from its parent point. If  $p$  has a parametric form and is radially symmetric (as for Thomas or Matérn cluster process), then  $g(r) = g(\|x\|)$  is expressed as a function of the model parameters.

We may try to exploit more information given by the estimator  $\hat{S}(r)$ . We look for  $\theta$  that minimizes the deviation of  $\hat{S}(r)$  from  $S_\theta(r)$  over some interval  $[a, b]$ . Define

$$D(\theta) = \int_a^b \left| \hat{S}(r)^q - S_\theta(r)^q \right|^p w(r) dr,$$

where  $0 \leq a < b$  and  $p, q > 0$  are given constants and  $w(r)$  is a weight function. The estimator of  $\theta$  is attained by minimizing the function  $D(\theta)$ . This method is called the *method of minimum contrast*. When the analytic expression of  $S_\theta(r)$  is unknown, we can approximate  $S_\theta(r)$  for fixed  $\theta$  by many simulations from the model. To compute the parameter estimate by the method of minimum contrast (with weight function identically 1), we can use the function `mincontrast` in the package `spatstat`. There, the default choice of parameters  $p$  and  $q$  is  $p = 2$  and  $q = 1/4$ . If  $S(r)$  is the  $K$ -function, `spatstat` enables to find the estimators for some particular models of point processes using the special functions `lgcp.estK` (log-Gaussian Cox process), `matclust.estK` (Matérn cluster process), and `thomas.estK` (Thomas point process). Similar functions `lgcp.estpcf`, `matclust.estpcf`, and `thomas.estpcf` exist for  $S(r)$  being the pair correlation function.

*Example:* Let  $\Phi$  be a Thomas point process with parameters  $\lambda_p$  (intensity of the parent Poisson process),  $\lambda_c$  (mean number of cluster points), and  $\sigma^2$  (variance of the normal distribution describing the displacement of a cluster point from its parent point). Then the pair correlation function is

$$g(r) = 1 + \frac{1}{\lambda_p(4\pi\sigma^2)^{d/2}} \exp\left\{-\frac{r^2}{4\sigma^2}\right\}, \quad r > 0.$$

We can estimate  $g(r)$  by the kernel estimator with some edge correction factor. Having such estimator  $\hat{g}(r)$  defined on the interval  $[a, b]$ , we can define the contrast function as

$$D(\lambda_p, \sigma^2) = \int_a^b \left( \hat{g}(r) - 1 - \frac{1}{\lambda_p(4\pi\sigma^2)^{d/2}} \exp\left\{-\frac{r^2}{4\sigma^2}\right\} \right)^2 dr.$$

The method of minimum contrast requires the minimization of this integral which has to be done by numerical methods. Notice that the parameter  $\lambda_c$  is not appearing in the contrast function so we have to estimate it by other approaches. The most natural approach would rely on the relation  $\lambda = \lambda_c \lambda_p$  and the estimation of the intensity  $\lambda$ .

### 1.3.3 Maximum likelihood method

Another approach for the estimation of model parameters is based on the maximum likelihood. Assume that  $\Phi$  is a finite point process with density  $p$  w.r.t. the distribution  $\Pi$  of unit Poisson process on the bounded Borel set  $B \subseteq \mathbb{R}^d$ . The density is parameterized by the vector  $\theta$  of unknown parameters,  $p(\varphi) = p_\theta(\varphi)$ . For simplicity, we consider that the observation window  $W$  coincides with  $B$ . The maximum likelihood estimator of  $\theta$  is obtained by maximizing the likelihood function  $L(\theta) = p_\theta(\varphi)$ , where  $\varphi$  is an observed realization of the process  $\Phi$ . Often it is more advantageous to maximize the log-likelihood function,  $l(\theta) = \log L(\theta)$ . The log-likelihood function is known for the Poisson process with intensity function  $\lambda_\theta$ :

$$l(\theta) = |W| - \int_W \lambda_\theta(y) dy + \sum_{x \in \varphi} \log \lambda_\theta(x).$$

As we already mentioned in Subsection 1.1, in homogeneous case ( $\lambda_\theta(x) = \lambda$ ) the argument of maxima is  $\lambda = \varphi(W)/|W|$ . For the inhomogeneous Poisson point process, the maximum likelihood estimator is not analytically tractable and we have to use numerical algorithms (e.g., the Newton-Raphson method) for maximization of the log-likelihood function.

For other processes than Poisson, the normalizing constant is typically given by a complicated integral, which is impossible to compute explicitly. In that case, we can use the Monte Carlo methods. Let the density of the point process have the form  $p_\theta(\varphi) = h_\theta(\varphi)/c_\theta$ , where  $h_\theta$  is a known function and  $c_\theta = \mathbb{E}h_\theta(\Phi_P)$  is the unknown normalizing constant ( $\Phi_P$  is the Poisson point process on  $B$  with unit intensity). Then  $l(\theta) = \log h_\theta(\varphi) - \log c_\theta$ . It will be more advantageous to maximize the likelihood ratio w.r.t. some fixed parameter value  $\theta_0$ ,

$$l(\theta) - l(\theta_0) = \log \frac{p_\theta(\varphi)}{p_{\theta_0}(\varphi)} = \log \frac{h_\theta(\varphi)}{h_{\theta_0}(\varphi)} - \log \frac{c_\theta}{c_{\theta_0}}.$$

For the first term, we have an analytic expression, while the second term may be approximated by MCMC (Markov Chain Monte Carlo) methods. The ratio of normalizing constants could be written as

$$\begin{aligned} \frac{c_\theta}{c_{\theta_0}} &= \frac{1}{c_{\theta_0}} \int h_\theta(\varphi) \Pi(d\varphi) = \int \frac{h_\theta(\varphi)}{h_{\theta_0}(\varphi)} \frac{h_{\theta_0}(\varphi)}{c_{\theta_0}} \Pi(d\varphi) \\ &= \int \frac{h_\theta(\varphi)}{h_{\theta_0}(\varphi)} p_{\theta_0}(\varphi) \Pi(d\varphi) = \int \frac{h_\theta(\varphi)}{h_{\theta_0}(\varphi)} P_{\theta_0}(d\varphi) = \mathbb{E}_{\theta_0} \frac{h_\theta(\Phi)}{h_{\theta_0}(\Phi)}, \end{aligned}$$

where  $P_{\theta_0}$  is the distribution of  $\Phi$  with density  $p_{\theta_0}$  (i.e. the true parameter is  $\theta_0$ ) and  $\mathbb{E}_{\theta_0}$  is the expectation w.r.t. this distribution. Here, we assume that  $h_{\theta_0}(\varphi) = 0$  implies  $h_\theta(\varphi) = 0$  and use the convention  $0/0 = 1$ . There exist different MCMC algorithms for generating the process with distribution  $P_{\theta_0}$ . They are based on the construction of a Markov chain  $\{\Phi^{(n)}\}$  whose limiting distribution is given by the density  $p_{\theta_0}$  w.r.t. the distribution  $\Pi$  of the point process  $\Phi_P$ . Replacing the expectation  $\mathbb{E}_{\theta_0} \frac{h_\theta(\Phi)}{h_{\theta_0}(\Phi)}$  by its sample mean, we get the approximation of the log-likelihood ratio

$$l_{\theta_0, n}(\theta) = \log \frac{h_\theta(\varphi)}{h_{\theta_0}(\varphi)} - \log \frac{1}{n} \sum_{i=0}^{n-1} \frac{h_\theta(\Phi^{(i)})}{h_{\theta_0}(\Phi^{(i)})}.$$

This approximation is called the *importance sampling approximation*. The maximization of  $l_{\theta_0, n}(\theta)$  gives the MCMC approximation  $\hat{\theta}_n$  of the maximum likelihood estimator  $\hat{\theta}$  of the parameter  $\theta$ . This approximation is usable if  $\theta_0$  is close to  $\hat{\theta}$ . As  $\theta_0$  we usually take some rough estimator obtained by a simpler and less efficient method. The whole procedure can be iteratively repeated. There exist also alternative approximations, the details can be found in [7, Subsections 8.2.4. and 8.2.5].

### 1.3.4 Maximum pseudolikelihood method

Since the likelihood function is often complicated, another strategy to estimate the model parameters is based on approximation of the likelihood function by some simpler variant.

**Definition 1.1.** Let  $\Phi$  be a finite point process on a bounded Borel set  $B$  with Papangelou conditional intensity  $\lambda_\theta^*(x, \varphi)$ , where  $\theta$  is the vector of unknown parameters. A realization  $\varphi$  of  $\Phi$  is observed in the window  $W$ . We assume that  $W$  coincides with  $B$ . We define the *pseudolikelihood* by the relation

$$\text{PL}(\theta) = \exp \left\{ |W| - \int_W \lambda_\theta^*(y, \varphi) dy \right\} \prod_{x \in \varphi} \lambda_\theta^*(x, \varphi \setminus \{x\}).$$

The estimator  $\hat{\theta}$  that maximizes  $\text{PL}(\theta)$  is called the *maximum pseudolikelihood estimator* of  $\theta$ .

**Remark 1.1.** For a Poisson point process,  $\lambda^*(x, \varphi) = \lambda(x)$  and thus the pseudolikelihood coincides with the likelihood.

*Example:* A Strauss point process has the Papangelou conditional intensity

$$\lambda^*(x, \varphi) = \beta \gamma^{t_R(x, \varphi)},$$

where  $t_R(x, \varphi) = \sum_{y \in \varphi} \mathbf{1}_{[0 < \|x-y\| \leq R]}$ . Unknown parameters are  $\beta > 0$ ,  $0 \leq \gamma \leq 1$  and  $R > 0$ . The logarithm of the pseudolikelihood is

$$\begin{aligned} \log \text{PL}(\beta, \gamma, R) &= |W| - \int_W \beta \gamma^{t_R(y, \varphi)} dy + \sum_{x \in \varphi} (\log \beta + t_R(x, \varphi \setminus \{x\}) \log \gamma) \\ &= |W| - \int_W \beta \gamma^{t_R(y, \varphi)} dy + \varphi(W) \log \beta + 2S_R(\varphi) \log \gamma, \end{aligned}$$

where  $S_R(\varphi) = \sum_{\{x, y\} \subseteq \varphi} \mathbf{1}_{[0 < \|x-y\| \leq R]}$ . If we put the derivatives w.r.t.  $\beta$  and  $\gamma$  equal to zero, we get the equations

$$\begin{aligned} \varphi(W) &= \beta \int_W \gamma^{t_R(y, \varphi)} dy, \\ 2S_R(\varphi) &= \beta \int_W t_R(y, \varphi) \gamma^{t_R(y, \varphi)} dy. \end{aligned}$$

The parameter  $R$  is considered to be known and we search for the solution numerically. In this way the estimators of  $\beta$  and  $\gamma$  are obtained. We realize that  $t_R(y, \varphi)$  takes only non-negative integer values and denote

$$m_k = \int_W \mathbf{1}_{[t_R(y, \varphi)=k]} dy, \quad k \in \mathbb{N}_0.$$

Then our system of equations has the form

$$\begin{aligned} \varphi(W) &= \beta \sum_{k=0}^{\varphi(W)} \gamma^k m_k, \\ 2S_R(\varphi) &= \beta \sum_{k=0}^{\varphi(W)} k \gamma^k m_k, \end{aligned}$$

because  $t_R(y, \varphi)$  is at most  $\varphi(W)$ . The advantage of assuming  $R$  to be known is that the Papangelou conditional intensity has a log-linear form. We can choose several different values  $R_1, \dots, R_K$  of a parameter  $R$ . For each value we calculate the maximum pseudolikelihood estimates of  $\beta$  and  $\gamma$ . Then we determine such  $R_i$ ,  $i = 1, \dots, K$ , for which the pseudolikelihood is the largest. This value is taken as the estimate of  $R$ .

### 1.3.5 Second-order composite likelihood

The maximum pseudolikelihood method belongs to a more general class of statistical methods that are based on so-called *composite likelihood*. These methods are used when the maximum likelihood method is computationally very demanding or inaccessible. The composite likelihood is a function obtained by multiplying a collection of likelihoods of simpler components. These components may not be independent. The particular form depends on the context. In the setting of point processes it was suggested to consider the product over the contributions of individual points or pairs of points.

Let  $\Phi$  be a stationary point process on  $\mathbb{R}^d$  with a second-order product density  $\lambda_\theta^{(2)}$  that is parameterized by a vector  $\theta$ . From stationarity it follows that  $\lambda_\theta^{(2)}(x, y) = \lambda_\theta^{(2)}(x - y)$ . Again we assume that  $\Phi$  is observed in the window  $W$ . Then the density of pairs of points in  $W$  is

$$f_\theta(x, y) = \frac{\lambda_\theta^{(2)}(x - y)}{\int_W \int_W \lambda_\theta^{(2)}(u - v) du dv}, \quad x, y \in W.$$

Of course, the distinct pairs of points are not independent. However, we consider the product of the densities  $f_\theta(x, y)$  over all observed pairs. After taking the logarithm, we have

$$\log \text{CL}(\theta) = \sum_{X, Y \in \Phi \cap W}^{\neq} \left[ \log \lambda_\theta^{(2)}(X - Y) - \log \int_W \int_W \lambda_\theta^{(2)}(u - v) du dv \right].$$

For practical purposes we disregard pairs at larger distances because for them the interactions are typically weak. Therefore, we do not lose much information if we omit them. Moreover, in this way we reduce computational complexity and variability of the resulting estimator. Let us choose  $R > 0$  and work with pairs of points in the distance smaller than  $R$ . We get the density

$$f_\theta^R(x, y) = \frac{\lambda_\theta^{(2)}(x - y) \mathbf{1}\{\|x - y\| < R\}}{\int_W \int_W \lambda_\theta^{(2)}(u - v) \mathbf{1}\{\|u - v\| < R\} du dv}, \quad x, y \in W,$$

and the logarithm of the composite likelihood

$$\log \text{CL}^R(\theta) = \sum_{X, Y \in \Phi \cap W: \|X - Y\| < R}^{\neq} \left[ \log \lambda_\theta^{(2)}(X - Y) - \log \int_W \int_W \lambda_\theta^{(2)}(u - v) \mathbf{1}\{\|u - v\| < R\} du dv \right].$$

The estimator of  $\theta$  is obtained by maximizing this function. Note that in the expression of  $f_\theta^R$  or  $\log \text{CL}^R(\theta)$  we are allowed to replace the product density  $\lambda_\theta^{(2)}$  by the pair correlation function  $g_\theta = \lambda_\theta^{(2)}/\lambda^2$ , where  $\lambda$  is the intensity of  $\Phi$ .

As opposed to previous two subsections we are now working with stationary point processes. The composite likelihood method is used mainly for Cox point processes where we often have an analytic form of the second-order product density. If  $\Phi$  is a stationary Cox point process with driving intensity function  $Z$  having the distribution depending on  $\theta$ , then  $\lambda_\theta^{(2)}(x - y) = \mathbb{E}Z(x)Z(y)$ . Next we will demonstrate another method suitable for Cox point processes. This method is based on another second-order characteristics.

### 1.3.6 Palm likelihood

Let  $\Phi$  be a stationary point process on  $\mathbb{R}^d$  with intensity  $\lambda$  and second-order product density  $\lambda^{(2)}$ . Then we can write  $\lambda^{(2)}(y - x) = \lambda \lambda_o(y - x)$ , where  $\lambda_o$  is called the *Palm intensity*. The second-order factorial moment measure can be expressed from the Campbell theorem as

$$\alpha^{(2)}(A \times B) = \mathbb{E} \sum_{X, Y \in \Phi}^{\neq} \mathbf{1}_{[X \in A, Y \in B]} = \int_A \int_B \lambda^{(2)}(y - x) dy dx = \lambda \int_A \int_{B-x} \lambda_o(u) du dx.$$

On the other hand, by the Campbell–Mecke theorem we have

$$\alpha^{(2)}(A \times B) = \mathbb{E} \sum_{X \in \Phi \cap A} \Phi(B \setminus \{X\}) = \lambda \int_A \int \varphi(B - x) P_o^!(d\varphi) dx.$$



Comparing these two expressions we find out that

$$E_o^! \Phi(B) = \int_B \lambda_o(u) du, \quad B \in \mathcal{B}^d,$$

i.e.  $\lambda_o$  is the intensity function of the reduced Palm distribution of  $\Phi$ . Now it is also clearer why  $\lambda_o$  is called the Palm intensity. Realize that it is a second-order characteristic. For a Poisson point process this function is constant. However, generally  $\lambda_o$  is not constant.

We will consider the point process of differences of observed points of  $\Phi$  in  $W$  with the distance smaller than  $R$ , i.e.

$$\Phi_R = \{Y - X : X \neq Y \in \Phi \cap W, \|Y - X\| < R\}.$$

Obviously, it is a point process in the ball  $b(o, R)$ . Its intensity measure is (by Campbell's theorem)

$$\begin{aligned} \mathbb{E} \Phi_R(A) &= \mathbb{E} \sum_{X, Y \in \Phi \cap W}^{\neq} \mathbf{1}_{[Y-X \in A]} = \int_W \int_W \mathbf{1}_{[y-x \in A]} \lambda^{(2)}(y-x) dy dx \\ &= \int \int \mathbf{1}_{[x \in W, x+u \in W]} \mathbf{1}_{[u \in A]} \lambda^{(2)}(u) dx du = \lambda \int_A |W \cap (W-u)| \lambda_o(u) du. \end{aligned}$$

Therefore, the intensity function of  $\Phi_R$  is

$$\lambda_R(u) = \lambda \lambda_o(u) |W \cap (W-u)|, \quad u \in b(o, R).$$

We assume that a parametric form  $\lambda_o^\theta(u)$  of the Palm intensity is given. We want to estimate the vector  $\theta$  of unknown parameters. To do this we consider  $\Phi_R$  as an inhomogeneous Poisson process with the intensity function  $\lambda_R(u)$  which is approximated so that the unknown true intensity is replaced by the observed intensity  $\Phi(W)/|W|$  and the term  $|W \cap (W-u)|$  is replaced by  $|W|$ , which is a reasonable approximation if  $R$  is substantially smaller than the window side. Altogether we approximate  $\lambda_R(u)$  as  $\Phi(W) \lambda_o(u)$ . The likelihood function is approximated by the likelihood function for the Poisson process with this approximated intensity function. Such likelihood function is referred to as the *Palm likelihood*. It means that the logarithm of the Palm likelihood is

$$\log L_P(\theta) = \sum_{X, Y \in \Phi \cap W}^{\neq} \mathbf{1}_{[\|Y-X\| < R]} \log \Phi(W) \lambda_o^\theta(Y-X) + |b(o, R)| - \Phi(W) \int_{b(o, R)} \lambda_o^\theta(u) du.$$

An alternative way how to get the Palm likelihood is to consider the point processes

$$\Phi_X = \{Y - X : Y \in \Phi \setminus \{X\}\}, \quad X \in \Phi \cap W,$$

which are inhomogeneous point processes with intensity function  $\lambda_o$ . We ignore interactions in the processes  $\Phi_X \cap b(o, R)$  and approximate them by inhomogeneous Poisson processes whose log-likelihoods are

$$\sum_{Y \in \Phi \cap W} \mathbf{1}_{[0 < \|X-Y\| < R]} \log \lambda_o^\theta(Y-X) + |b(o, R)| - \int_{b(o, R)} \lambda_o^\theta(u) du.$$

Now we regard  $\Phi_X$ ,  $X \in \Phi \cap W$ , as independent identically distributed point processes, and we ignore edge effects. Then the logarithmic Palm likelihood has the form

$$\log L_P(\theta) = \sum_{X, Y \in \Phi \cap W} \mathbf{1}_{[0 < \|X-Y\| < R]} \log \lambda_o^\theta(Y-X) + \Phi(W) |b(o, R)| - \Phi(W) \int_{b(o, R)} \lambda_o^\theta(u) du,$$

which differs from the previous expression only by a constant.

### 1.3.7 Takacs–Fiksel method

In order to find the maximum likelihood estimate one solves the equation  $l'(\theta) = 0$ . The method of moments is based on the relation  $S_\theta(r) - \hat{S}(r) = 0$ . Both of these approaches could be included in the framework of estimating equations.

**Definition 1.2.** Let  $\Phi$  be a point process with distribution  $Q_\theta$  depending on an unknown parameter  $\theta \in \Theta$ . Consider a function  $\psi : \Theta \times \mathcal{N} \rightarrow \mathbb{R}$  such that  $\mathbb{E}_\theta \psi(\theta, \Phi) = 0$  for each  $\theta \in \Theta$ . Here,  $\mathbb{E}_\theta$  denotes the expectation w.r.t.  $Q_\theta$ . For a given realization  $\varphi$  the equation  $\psi(\theta, \varphi) = 0$  is called an *unbiased estimating equation*. By different choices of  $\psi$  we obtain a system of equations. Its solution  $\hat{\theta}$  is used as an estimator of  $\theta$  based on  $\varphi$ .

Besides the method of moments and maximum likelihood (or pseudolikelihood) method, another example of estimating equations for point process models is given by the Takacs–Fiksel method. This method is based on the Georgii–Nguyen–Zessin identity

$$\mathbb{E} \sum_{X \in \Phi} h(X, \Phi \setminus \{X\}) = \int_{\mathbb{R}^d} \mathbb{E} h(x, \Phi) \lambda^*(x, \Phi) dx, \quad (3)$$

where  $\lambda^*$  is the conditional intensity of  $\Phi$ .

In the case of a finite point process  $\Phi$  with density  $p$  w.r.t. the distribution of unit Poisson point process  $\Phi_P$  on the bounded Borel set  $B$ ,  $\lambda^*$  is the Papangelou conditional intensity. If  $\Phi$  is the Poisson point process on  $\mathbb{R}^d$  with intensity function  $\lambda$ , then  $\lambda^*(x, \Phi) = \lambda(x)$ .

Assume that we know the parametric form of the conditional intensity  $\lambda_\theta^*(x, \varphi)$ . We define

$$\psi_h(\theta, \varphi) = \sum_{x \in \varphi \cap W} h(x, \varphi \setminus \{x\}) - \int_W h(x, \varphi) \lambda_\theta^*(x, \varphi) dx$$

for an arbitrary function  $h$ . Then the Georgii–Nguyen–Zessin identity implies  $\mathbb{E}_\theta \psi_h(\theta, \Phi) = 0$ . The estimator of  $\theta$  is obtained as the solution of the unbiased estimating equation  $\psi_h(\theta, \varphi) = 0$ . Similarly as in the method of minimum contrast, we can choose more functions  $h$  than the number of unknown parameters. For example, if we have  $k$  function  $h_1, \dots, h_k$ , we may search for  $\theta$  which minimizes

$$\sum_{i=1}^k \psi_{h_i}(\theta, \varphi)^2.$$

**Remark 1.2.** When we obtain the estimator  $\hat{\theta}$  as the solution of an unbiased estimating equation, it doesn't mean that it is an unbiased estimator of  $\theta$ .

For some natural choices of  $h$  it may be impossible to determine  $\psi_h(\theta, \varphi)$  only from the observation  $\varphi$  in a bounded window  $W$ . The problems with edge effects may arise. Then instead of  $\psi_h(\theta, \varphi)$  we can take the estimate  $\hat{\psi}_h(\theta, \varphi)$  which considers corrections of edge effects. For example, in the case of a stationary point process put  $h(x, \varphi) = \varphi(b(x, r)) \mathbf{1}_{[x \in W]} / |W|$ . Then the expectation of the first term in  $\psi_h(\theta, \Phi)$ , i.e. the left-hand side in (3), is equal to  $\lambda^2 K(r)$ . However, we are not able to determine the first term in  $\psi_h(\theta, \varphi)$  only from the information inside  $W$ . The solution is to replace it by some edge-corrected estimator of  $\lambda^2 K(r)$ , e.g., the estimator with translation edge correction.

*Example:* Consider a Strauss point process and assume that the parameter  $R$  is known. Our goal is to estimate  $\theta = (\beta, \gamma)$ . The choice  $h_1(x, \varphi) = 1$  gives

$$\psi_{h_1}(\theta, \varphi) = \varphi(W) - \beta \int_W \gamma^{t_R(x, \varphi)} dx.$$

Taking  $h_2(x, \varphi) = t_R(x, \varphi)$  we get

$$\psi_{h_2}(\theta, \varphi) = 2S_R(\varphi) - \beta \int_W t_R(x, \varphi) \gamma^{t_R(x, \varphi)} dx.$$

Note that we have obtained the same two equations with two unknown parameters as in the case of the maximum pseudolikelihood method.

## 1.4 Model diagnostics

In order to verify that the fitted parametric model is appropriate we can exploit the idea of Monte Carlo tests. If we are able to simulate from our fitted model, then we can determine some summary characteristics for each simulated realization of the model. The results from simulations could be compared with the characteristics estimated from the data. This comparison should not show any substantial deviations when the parameters of the model are determined correctly. The problems with this approach will begin to exhibit for more general inhomogeneous point processes where we would need appropriate summary characteristics.

Now let us examine how we can use the generalization of residuals from the classical linear regression models in the context of point processes. Generally, the residuals are differences between the observed and fitted values. If the fitted model is correct, the residuals should fluctuate around zero. On the contrary, large deviations from zero may indicate what is wrong with the fitted model (e.g., incorrectly estimated trend or interactions).

**Definition 1.3.** Let  $\Phi$  be a point process with the conditional intensity  $\lambda^*$ . For some non-negative measurable function  $h$  we define *h-weighted innovation* as the signed random measure

$$I(B, h, \lambda^*) = \sum_{X \in \Phi \cap B} h(X, \Phi \setminus \{X\}) - \int_B h(x, \Phi) \lambda^*(x, \Phi) dx.$$

According to the Georgii–Nguyen–Zessin identity (3) we have  $\mathbb{E}I(B, h, \lambda^*) = 0$  for any  $B \in \mathcal{B}^d$ .

*Example:* Let  $\Phi$  be the Poisson point process with intensity function  $\lambda$  and consider the following three choices of  $h$ :  $h(x, \varphi) = 1$ ,  $h(x, \varphi) = 1/\lambda^*(x, \varphi)$  and  $h(x, \varphi) = 1/\sqrt{\lambda^*(x, \varphi)}$ . Since  $\lambda^*(x, \varphi) = \lambda(x)$ , we get

$$\begin{aligned} I(B, 1, \lambda) &= \Phi(B) - \int_B \lambda(x) dx, \\ I(B, 1/\lambda^*, \lambda) &= \sum_{X \in \Phi \cap B} \frac{1}{\lambda(X)} - |B|, \\ I(B, 1/\sqrt{\lambda^*}, \lambda) &= \sum_{X \in \Phi \cap B} \frac{1}{\sqrt{\lambda(X)}} - \int_B \sqrt{\lambda(x)} dx. \end{aligned}$$

It is easy to verify that  $\mathbb{E}I(B, h, \lambda) = 0$  by direct computation using the Campbell theorem. For the variances we have

$$\begin{aligned} \text{var } I(B, 1, \lambda) &= \int_B \lambda(x) dx, \\ \text{var } I(B, 1/\lambda^*, \lambda) &= \int_B \frac{1}{\lambda(x)} dx, \\ \text{var } I(B, 1/\sqrt{\lambda^*}, \lambda) &= |B|. \end{aligned}$$

These relations follow directly from the following lemma.

**Lemma 1.1.** Let  $\Phi$  be the Poisson point process on  $\mathbb{R}^d$  with intensity measure  $\Lambda$ . Then for an arbitrary non-negative measurable function  $f$ ,

$$\text{var} \sum_{X \in \Phi} f(X) = \int f(x)^2 \Lambda(dx).$$

*Proof.* From the Campbell theorem,

$$\mathbb{E} \sum_{X \in \Phi} f(X) = \int f(x) \Lambda(dx).$$

The second moment could be rewritten using the second-order Campbell theorem. Moreover, we make use of the fact that the second-order factorial moment measure of the Poisson point process is  $\Lambda \times \Lambda$ ,

$$\begin{aligned} \mathbb{E} \left( \sum_{X \in \Phi} f(X) \right)^2 &= \mathbb{E} \sum_{X, Y \in \Phi} f(X)f(Y) = \mathbb{E} \sum_{X \in \Phi} f(X)^2 + \mathbb{E} \sum_{X, Y \in \Phi}^{\neq} f(X)f(Y) \\ &= \int f(x)^2 \Lambda(dx) + \int \int f(x)f(y) \Lambda(dx) \Lambda(dy) = \int f(x)^2 \Lambda(dx) + \left( \int f(x) \Lambda(dx) \right)^2. \end{aligned}$$

From this we already obtain the assertion.  $\square$

Assume that the conditional intensity  $\lambda_{\theta}^*(x, \varphi)$  depends on the parameter  $\theta$ . Furthermore, suppose that we found the estimator  $\hat{\theta}$  (e.g., by one of the methods from Subsection 1.3) based on the observation of  $\Phi$  in the bounded window  $W$ . Then the estimator of the conditional intensity is  $\widehat{\lambda}^*(x, \varphi) = \lambda_{\hat{\theta}}^*(x, \varphi)$ . In the definition of innovation we admit that  $h$  depends on  $\hat{\theta}$ . Put  $\hat{h}(x, \varphi) = h_{\hat{\theta}}(x, \varphi)$ .

**Definition 1.4.** A random signed measure

$$R(B, \hat{h}, \hat{\theta}) = I(B, \hat{h}, \widehat{\lambda}^*) = \sum_{X \in \Phi \cap B} h_{\hat{\theta}}(X, \Phi \setminus \{X\}) - \int_B h_{\hat{\theta}}(x, \Phi) \lambda_{\hat{\theta}}^*(x, \Phi) dx$$

is called the *h-weighted residual measure*.

Since  $\mathbb{E}I(B, h, \lambda^*) = 0$ , we expect that  $R(B, \hat{h}, \hat{\theta})$  is around zero when the model with  $\hat{\theta}$  is correct. Note that the expectation  $\mathbb{E}R(B, \hat{h}, \hat{\theta})$  does not have to be zero but it should be approximately zero when the model is true. Regions  $B$  with extreme values of  $R(B, \hat{h}, \hat{\theta})$  may indicate regions of irregularity. The usual choices for  $h$  include  $h = 1$  (*raw residuals*),  $h = 1/\lambda^*$  (*inverse-lambda residuals*) and  $h = 1/\sqrt{\lambda^*}$  (*Pearson residuals*). Residuals for these three choices may be computed by `residuals.ppm` in package `spatstat`. For  $h = 1$  we have

$$R(B, 1, \hat{\theta}) = \Phi(B) - \int_B \lambda_{\hat{\theta}}^*(x, \Phi) dx.$$

It means that the raw residual measure is given as the difference of an atomic measure with atoms in the observed points and the measure with density  $\lambda_{\hat{\theta}}^*(x, \Phi)$  w.r.t. the Lebesgue measure.

*Example:* Consider a stationary Poisson point process  $\Phi$  with intensity  $\lambda$ , which could be estimated as  $\hat{\lambda} = \Phi(W)/|W|$ . Then (provided that  $\Phi(W) > 0$ )

$$\begin{aligned} R(B, 1, \hat{\lambda}) &= \Phi(B) - \Phi(W) \frac{|B|}{|W|}, \\ R(B, 1/\hat{\lambda}^*, \hat{\lambda}) &= |W| \frac{\Phi(B)}{\Phi(W)} - |B|, \\ R(B, 1/\sqrt{\hat{\lambda}^*}, \hat{\lambda}) &= \Phi(B) \sqrt{\frac{|W|}{\Phi(W)}} - |B| \sqrt{\frac{\Phi(W)}{|W|}}. \end{aligned}$$

It can be shown that the expectations of these three  $h$ -weighted residual measures are 0. We can also notice that  $R(W, 1, \hat{\lambda}) = R(W, 1/\hat{\lambda}^*, \hat{\lambda}) = R(W, 1/\sqrt{\hat{\lambda}^*}, \hat{\lambda}) = 0$ . This corresponds to the situation in the classical linear regression where the sum of all residuals is 0.

For the graphical representation of the residuals it is convenient to perform kernel smoothing.

**Definition 1.5.** Let  $k$  be a probability density on  $\mathbb{R}^d$ . The realization of a point process  $\Phi$  is observed in a bounded window  $W$ . We have constructed the estimator  $\hat{\theta}$  of  $\theta$ . Define the *smoothed residual field* by the relation

$$S(x) = e(x) \int_W k(x-y) R(dy, \hat{h}, \hat{\theta}),$$

where

$$e(x) = \left( \int_W k(x-y) dy \right)^{-1}$$

is the edge correction.

**Remark 1.3.** For  $h = 1$  we have

$$S(x) = e(x) \sum_{Y \in \Phi \cap W} k(x-Y) - e(x) \int_W k(x-y) \lambda_\Phi^*(y, \Phi) dy.$$

## 2 Marked point processes

Let  $\Phi_m$  be a marked point process on  $\mathbb{R}^d$  with the mark space  $\mathbb{M}$ . The corresponding unmarked point process is denoted by  $\Phi$ .

### 2.1 Estimation of summary characteristics

We will assume that a stationary marked point process  $\Phi_m$  is observed in a bounded window  $W$ . The estimators of summary characteristics are mostly either straightforward from the definition or it is enough to suitably modify the estimators that were defined for point processes (Subsection 1.1).

First consider the case of qualitative marks  $\mathbb{M} = \{1, \dots, k\}$ . It means that  $\Phi_m$  is a  $k$ -dimensional point process  $(\Phi_1, \dots, \Phi_k)$ . The intensities  $\lambda_i$  may be estimated by  $\Phi_i(W)/|W|$ . The stationary mark distribution  $\mathbb{Q}$  is an atomic measure on  $\mathbb{M}$  and natural estimators of  $p_i = \mathbb{Q}(\{i\})$  are  $\Phi_i(W)/\Phi(W)$ .

For the cross  $G$ -function  $G_{ij}(r) = P_o^{i,j}(\{\varphi_m : d(o, \varphi_j) \leq r\})$  and the condensed  $G$ -function  $G_i(r) = P_o^{i,i}(\{\varphi_m : d(o, \varphi) \leq r\})$ ,  $r \geq 0$ , we may use, for example, the Kaplan–Meier estimator:

$$\begin{aligned} \widehat{G}_{ij}(r) &= 1 - \prod_{s \leq r} \left( 1 - \frac{\#\{X \in \Phi_i \cap W : e_j(X) = s, e_j(X) \leq c(X)\}}{\#\{X \in \Phi_i \cap W : e_j(X) \geq s, c(X) \geq s\}} \right), \\ \widehat{G}_i(r) &= 1 - \prod_{s \leq r} \left( 1 - \frac{\#\{X \in \Phi_i \cap W : e(X) = s, e(X) \leq c(X)\}}{\#\{X \in \Phi_i \cap W : e(X) \geq s, c(X) \geq s\}} \right), \end{aligned}$$

where  $c(x) = d(x, \partial W)$  is the distance of  $x$  from the window boundary,  $e_j(x) = d(x, \Phi_j \setminus \{x\})$  is the distance of  $x$  from the nearest point of  $\Phi_j$  and  $e(x) = d(x, \Phi \setminus \{x\})$  is the distance of  $x$  from the nearest point of the process (regardless of the mark). In `spatstat` we can obtain these estimates using `Gcross` and `Gdot`.

The cross  $K$ -function  $K_{ij}$  is defined by the relation

$$\lambda_j K_{ij}(r) = \mathbb{E}_o^{i,j} \Phi_j(b(o, r))$$

and the condensed  $K$ -function  $K_i$  by the relation

$$\lambda K_i(r) = \mathbb{E}_o^{i,i} \Phi(b(o, r)).$$

Here  $\lambda_i$  is the intensity of  $\Phi_i$  and we already noted that its natural estimator is  $\widehat{\lambda}_i = \Phi_i(W)/|W|$ . We mention the estimators of the cross (`Kcross`) and condensed (`Kdot`)  $K$ -function that use the edge correction factor  $e_{W,r}(X, Y)$ ,

$$\begin{aligned} \widehat{K}_{ij}(r) &= \frac{1}{\widehat{\lambda}_i \widehat{\lambda}_j} \sum_{X \in \Phi_i \cap W, Y \in \Phi_j \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W|} e_{W,r}(X, Y), \\ \widehat{K}_i(r) &= \frac{1}{\widehat{\lambda}_i \widehat{\lambda}} \sum_{X \in \Phi_i \cap W, Y \in \Phi \cap W}^{\neq} \frac{\mathbf{1}_{[\|X-Y\| \leq r]}}{|W|} e_{W,r}(X, Y). \end{aligned}$$

Particular choices of  $e_{W,r}(X, Y)$  are considered in Subsection 1.1. Notice that this estimator of the cross  $K$ -function satisfies  $\widehat{K}_{ij}(r) = \widehat{K}_{ji}(r)$ .

If  $\Phi_m$  is motion-invariant then the pair correlation function  $g_{ij}(x, y) = g((x, i), (y, j))$ ,  $x, y \in \mathbb{R}^d$ ,  $i, j \in \{1, \dots, k\}$  depends only on  $r = \|x - y\|$ . Its kernel estimator has the form

$$\widehat{g}_{ij}(r) = \frac{1}{\widehat{\lambda}_i \widehat{\lambda}_j} \sum_{X \in \Phi_i \cap W, Y \in \Phi_j \cap W}^{\neq} \frac{k_b(r - \|X - Y\|)}{\sigma_d r^{d-1} |W|} e_{W,r}(X, Y),$$

where  $k_b$  is a kernel function with bandwidth  $b$ .

For marked point processes with quantitative marks we first deal with the estimation of a non-normalized  $f$ -mark correlation function. It is defined as  $\kappa_f(r) = \frac{\lambda_f^{(2)}(r)}{\lambda^{(2)}(r)}$ , where  $\lambda_f^{(2)}(r)$  is the density of the second-order  $f$ -weighted factorial moment measure

$$\alpha_f^{(2)}(B_1 \times B_2) = \mathbb{E} \sum_{(X_1, M_1), (X_2, M_2) \in \Phi_m}^{\neq} \mathbf{1}_{[X_1 \in B_1, X_2 \in B_2]} f(M_1, M_2), \quad B_1, B_2 \in \mathcal{B}^d.$$

Here, the process  $\Phi_m$  is assumed to be motion-invariant. The kernel estimator of  $\lambda_f^{(2)}(r)$  is

$$\widehat{\lambda}_f^{(2)}(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{f(M(X), M(Y)) k_b(\|X - Y\| - r)}{\sigma_d r^{d-1} |W|} e_{W,r}(X, Y),$$

while an analogous estimator of the second-order product density  $\lambda^{(2)}(r)$  is

$$\widehat{\lambda}^{(2)}(r) = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{k_b(\|X - Y\| - r)}{\sigma_d r^{d-1} |W|} e_{W,r}(X, Y).$$

Then the non-normalized  $f$ -mark correlation function can be estimated as

$$\widehat{\kappa}_f(r) = \frac{\widehat{\lambda}_f^{(2)}(r)}{\widehat{\lambda}^{(2)}(r)}, \quad r > 0.$$

The function  $\kappa_f(r)$  can be also expressed using the two-point mark distribution,

$$\kappa_f(r) = \mathbb{E}_{or} f(M(o), M(r)).$$

This conditional expectation could be estimated by arithmetic mean of  $f$ -values for marks corresponding to the points at distance  $r$ . The number of pairs exactly at distance  $r$  will be usually small. Therefore, we take  $\varepsilon > 0$  and put

$$\widehat{\kappa}_f(r) = \frac{1}{N_f(\varepsilon)} \sum_{X, Y \in \Phi \cap W: \|\|X - Y\| - r\| < \varepsilon/2}^{\neq} f(M(X), M(Y)),$$

where  $N_f(\varepsilon) = \#\{X, Y \in \Phi \cap W : \|\|X - Y\| - r\| < \varepsilon/2\}$ . Note that this estimator corresponds to the former one with constant kernel function on  $[-\varepsilon/2, \varepsilon/2]$  and no edge corrections.

Furthermore, we estimate the normalized  $f$ -mark correlation function  $k_f(r) = \frac{\kappa_f(r)}{c_f}$  by

$$\widehat{k}_f(r) = \frac{\widehat{\kappa}_f(r)}{\widehat{c}_f}, \quad r > 0,$$

where

$$\widehat{c}_f = \frac{1}{\Phi(W)^2} \sum_{X, Y \in \Phi \cap W} f(M(X), M(Y))$$

is the estimator of  $c_f = \int \int f(m_1, m_2) \mathbb{Q}(dm_1) \mathbb{Q}(dm_2)$ . Denote  $\mu = \mathbb{E}M_0 = \int m \mathbb{Q}(dm)$  the mean typical mark. Then for  $f(m_1, m_2) = c(m_1, m_2) = m_1 m_2$  we have  $c_f = \mu^2$  and  $k_f = k_c$  is known as Stoyan's

mark correlation function. For  $f(m_1, m_2) = e(m_1, m_2) = m_1$ ,  $c_f = \mu$  and  $k_f = k_e$  is referred to as the  $r$ -mark function. The values  $k_c(r)$  or  $k_e(r)$  larger than 1 indicate mutual stimulation in the distance  $r$ . On the other hand, the values smaller than 1 correspond to inhibition. The estimates of  $\kappa_f(r)$  and  $k_f(r)$  could be calculated using the function `markcorr` in the package `spatstat`.

The functions  $\kappa_f(r)$  and  $k_f(r)$  are examples of non-cumulative summary characteristics. A cumulative analogue is the mark-weighted  $K$ -function  $K_f(r)$  which generalizes the  $K$ -function for point processes in the following way:

$$\lambda K_f(r) = \frac{1}{c_f} \mathbb{E}_o^! \sum_{(X, M(X)) \in \Phi_m} f(M(o), M(X)) \mathbf{1}_{[X \in b(o, r)]},$$

where  $\lambda$  is the intensity of a stationary process  $\Phi_m$ . An unbiased estimator of  $\lambda^2 c_f K_f(r)$  has the form

$$\lambda^2 \widehat{c_f K_f(r)} = \sum_{X, Y \in \Phi \cap W}^{\neq} \frac{f(M(X), M(Y)) \mathbf{1}_{\{\|X - Y\| \leq r\}}}{|W|} e_{W, r}(X, Y).$$

In order to estimate  $K_f(r)$  we have to divide by estimators of  $\lambda^2$  and  $c_f$ . For particular choices  $c(m_1, m_2) = m_1 m_2$ ,  $e(m_1, m_2) = m_1$ ,  $e^*(m_1, m_2) = m_2$ ,  $\gamma(m_1, m_2) = \frac{1}{2}(m_1 - m_2)^2$  we get the functions  $K_c$ ,  $K_e$ ,  $K_{e^*}$ ,  $K_\gamma$ , respectively.

As a representative of numerical characteristics we consider the non-normalized nearest-neighbour correlation index defined as  $\bar{\nu}_f = \mathbb{E}_o^! f(M(o), M(Z_1))$ , where  $Z_1$  is the point of the process that is the closest to the origin and  $M(Z_1)$  is the mark at  $Z_1$ . This index can be naturally estimated by

$$\widehat{\bar{\nu}_f} = \frac{1}{\Phi(W)} \sum_{X \in \Phi \cap W} f(M(X), M(Z_X)),$$

where  $Z_X \in \Phi$  is the nearest neighbour of  $X$ . A natural estimator of the normalized nearest-neighbour correlation index  $\bar{n}_f = \frac{\bar{\nu}_f}{c_f}$  is obtained as

$$\widehat{\bar{n}_f} = \frac{\widehat{\bar{\nu}_f}}{\widehat{c_f}}.$$

For  $f(m_1, m_2) = m_1 m_2$  the values  $\bar{n}_f > 1$  indicate mutual stimulation between neighbours. In `spatstat` the commands `nncorr` ( $f(m_1, m_2) = m_1 m_2$ ), `nnmean` ( $f(m_1, m_2) = m_1$ ) and `nnvario` ( $f(m_1, m_2) = (m_1 - m_2)^2/2$ ) may be used.

## 2.2 Tests of independence

Statistical analysis of a marked point process mostly starts with the test of the hypothesis of independent marks. If the marks may be considered independent, we can use the methods developed for independent data. First, we will deal with the testing of independent marks. Afterwards, we mention some approaches for testing the independence of marks and locations. We pursue a non-parametric approach and use simulation tests whose principle was explained in Subsection 1.2.

### 2.2.1 Testing independent marking

First consider a two-dimensional point process  $\Phi_m = (\Phi_1, \Phi_2)$ . The null hypothesis of independent marks may have two different interpretations:

1. independent marking (random labelling) – to the points of  $\Phi$  independently randomly either mark 1 or mark 2 is assigned,
2. random superposition – two independent point processes  $\Phi_1$  and  $\Phi_2$  form a bivariate point process.

The first situation is an example of *posterior marking* – we describe how the marks were created conditionally on given locations of points. This is an appropriate model when the points are tree locations and the trees could be affected by some disease or catastrophe (mark 1) or not (mark 2). In the second situation we have *prior marking* – a marked point process is formed by certain mechanism, namely union of two independent populations.

For testing the hypothesis that  $\Phi_m$  is independently marked point process the method of *random allocation* is used. We fix the locations and create new marks by random permutation of observed marks (`rlabel` in `spatstat`). We generate  $M$  such permutations and carry out the corresponding Monte Carlo test. Under the hypothesis of independent marking,

$$\begin{aligned} K(r) &= K_{11}(r) = K_{22}(r) = K_{12}(r) = K_{1\cdot}(r), \\ g(r) &= g_{11}(r) = g_{22}(r) = g_{12}(r), \\ G(r) &= G_{1\cdot}(r), \\ J(r) &= J_{1\cdot}(r), \end{aligned}$$

where  $K(r)$ ,  $g(r)$ ,  $G(r)$  and  $J(r)$  are functional characteristics of the unmarked point process  $\Phi$ . Therefore, a useful test statistic could be  $S(r) = K_{1\cdot}(r) - K(r)$ , which is equal to  $S_0(r) = 0$  if the null hypothesis is true.

When we want to test the hypothesis of random superposition of point processes  $\Phi_1$  and  $\Phi_2$ , we can use the method of *random shift*. The locations of points with mark 1 are fixed. We generate  $M$  realizations of the subprocess  $\Phi_2$  so that all its points are simultaneously shifted (`rshift`) by a vector with prescribed length  $R > 0$ . For each of  $M$  realizations of a bivariate point process we calculate estimator of  $S(r)$  and apply simultaneous Monte Carlo test. As a function  $S(r)$  we may use one of the cross functional characteristics. Under the hypothesis of random superposition the following relations hold:

$$\begin{aligned} K_{12}(r) &= \omega_d r^d, \\ g_{12}(r) &= 1, \\ G_{12}(r) &= F_2(r), \\ J_{12}(r) &= 1. \end{aligned}$$

In the case of process with quantitative marks we can again test the hypothesis of independent marking by the method of random allocation. It means that the locations are kept fixed and the marks are assigned by permuting the observed marks (sampling without replacement). In this way all  $M$  simulations lead to the same empirical mark distribution. Another possibility is to generate marks from the empirical mark distribution (sampling with replacement). The test statistic could be one of the  $f$ -mark correlation functions or the  $f$ -weighted  $K$ -function. For motion-invariant independently marked point processes we have  $k_f(r) = 1$  and  $K_c(r) = K_\gamma(r) = K(r)$ , where  $K(r)$  is the  $K$ -function of the corresponding unmarked point process.

### 2.2.2 Independence of marks and locations

We are going to present three methods for testing the independence of marks and locations in marked point processes with quantitative marks. If the marks and locations are independent, we may investigate both components separately, which simplifies the statistical inference. The geostatistical marking is an appropriate marking model where the marks are independent of locations.

The first method is based on the summary characteristics  $K_f(r)$ , in particular special cases  $K_e(r)$  and  $K_{e^*}(r)$ . Both these functions are equal to the  $K$ -function  $K(r)$  of the unmarked point process  $\Phi$  if the process is geostatistically marked. The test works conditionally on the locations and it is based on the random allocation. We will generate  $M$  realizations by random permutations of marks (or by random sample from empirical mark distribution). Since the locations are fixed, the estimators of  $K(r)$  will be the same for all  $M$  simulations as well as for data. We use this estimator of  $K(r)$  as the statistics  $S_0(r)$  in Monte Carlo test. Further, we estimate  $K_e(r)$  or  $K_{e^*}(r)$  from data and also from  $M$  simulations. Under the null hypothesis, all these  $M + 1$  functions should look approximately like  $S_0(r)$ . By simultaneous or integral Monte Carlo test, we find out whether the estimate from data significantly differs. This approach ignores correlations among marks. This can cause that the hypothesis is rejected not because of dependence between marks and locations but because of dependencies within marks.

The second method originates from the paper [13]. It is based on the fact that the functions  $E(r) = \kappa_e(r) = \mathbb{E}_{or} M(r)$  and  $V(r) = \kappa_v(r) - \kappa_e(r)^2 = \mathbb{E}_{or} (M(o) - E(r))^2$  are constant for motion-invariant geostatistically marked point process. If the estimates of these functions from the data significantly differ from a constant function, it gives evidence against the null hypothesis. Defining  $E(0) = \mathbb{E}M_0$  and



$V(0) = \text{var } M_0$ , we can perform simultaneous Monte Carlo test with the choice  $S(r) = E(r) - E(0)$  or  $S(r) = V(r) - V(0)$ . Under the null hypothesis,  $S_0(r) = 0$ .

Also the third test is based on the principle of simultaneous tests. It was proposed in the paper [3]. Let  $\varphi_m = \{(x_1, m_1), \dots, (x_n, m_n)\}$  be a realization of a motion-invariant marked point process  $\Phi_m$  observed in the window  $W$ . Assume that the data are recorded in some fixed order. Let  $\delta(x_i) = d(x_i, \{x_{i+1}, \dots, x_n\})$ ,  $i = 1, \dots, n$ , denote the distance of point  $x_i$  to the nearest point of the process with a larger index. For given  $r > 0$  we choose those points with  $\delta(x_i) \leq r$ . The number of selected points will be  $n_r$ . It is recommended to choose  $r$  small in comparison with distance of nearest neighbours. Since the selection of points does not depend on marks, the mean of marks of  $n_r$  points should be under the null hypothesis close to the mean of marks of arbitrary randomly selected  $n_r$  points out of  $n$  points. On the other hand, if a mark is dependent on the presence of other points in the vicinity of its location, then the means of marks selected according to the proposed criterion and the means of randomly selected marks should differ significantly. We generate  $M$  different random samples of marks and for each such sample of size  $n_r$  we compute its mean. The test itself works in the same way as the classical Monte Carlo test described in Subsection 1.2 ( $T$  is the mean of  $n_r$  marks).

### 3 Geostatistics

Geostatistics is a part of spatial statistics dealing with data formed by finitely many observations of a given variable in some fixed spatial locations.

The geostatistical data are modelled by a random field  $\{Z(x) : x \in D\}$ , where  $D \subseteq \mathbb{R}^d$  has positive  $d$ -dimensional Lebesgue measure. Recall that an intrinsically stationary random field satisfies the conditions  $\mathbb{E}(Z(x) - Z(y)) = 0$  and  $\text{var}(Z(x) - Z(y)) = 2\gamma(x - y)$ . The function  $2\gamma(h) = \text{var}(Z(x + h) - Z(x)) = \mathbb{E}(Z(x + h) - Z(x))^2$  is called the variogram. Our first aim is to estimate the variogram from the observations  $Z(x_1), \dots, Z(x_n)$ , where  $x_1, \dots, x_n \in D$  are fixed deterministic points.

#### 3.1 Variogram estimation

##### 3.1.1 Non-parametric estimators

To get the first impression about the variogram, we can plot the squares of differences of observed values  $(Z(x_i) - Z(x_j))^2$  against  $x_i - x_j$  or  $\|x_i - x_j\|$ . Such graph is called the *empirical variogram cloud* and can be obtained in the package `geoR` [11] by `variog` with `option="cloud"`. This graph often does not give a clear picture because the number of possible pairs  $\{x_i, x_j\}$  of distinct points could be quite large, specifically, it is  $\binom{n}{2}$ . More useful information is obtained by averaging the values corresponding to the same difference  $x_i - x_j$ . Then we have the following unbiased estimator of the variogram,

$$2\hat{\gamma}(h) = \frac{1}{|N(h)|} \sum_{N(h)} (Z(x_i) - Z(x_j))^2, \quad (4)$$

where  $N(h) = \{(x_i, x_j) : x_i - x_j = h, i, j = 1, \dots, n\}$  and  $|N(h)|$  is the cardinality of  $N(h)$ . It is in fact the estimator obtained by the method of moments. The following properties of the estimator can be easily seen:  $\hat{\gamma}(h) \geq 0$ ,  $\hat{\gamma}(0) = 0$  and  $\hat{\gamma}(h) = \hat{\gamma}(-h)$ . Thus, the estimator preserves the basic theoretical properties of the variogram. Note that  $N(h)$  and  $N(-h)$  could be different sets but they have the same cardinality and  $(x_i, x_j) \in N(h)$  if and only if  $(x_j, x_i) \in N(-h)$ . Therefore, the symmetry of the estimator follows. For small sample size or irregularly scattered points  $x_1, \dots, x_n$ , where the measurements are taken, the number of pairs in  $N(h)$  will be very small and the estimator of  $2\gamma(h)$  will have large variance. The practical recommendation is to use  $h$  for which  $|N(h)| \geq 30$ . If we are unable to assure this condition, we divide (similarly as in the construction of histogram) pairs of points into several groups with similar differences  $x_i - x_j$ . We calculate the mean of variables  $(Z(x_i) - Z(x_j))^2$  in each group. In the package `RandomFields` we can perform this by calling the function `EmpiricalVariogram`, in the package `gstat` [10] by `variogram`. Another possibility is to use kernel smoothing with a kernel function  $k_b$  and bandwidth  $b$ :

$$2\hat{\gamma}(h) = \frac{\sum_{i \neq j} (Z(x_i) - Z(x_j))^2 k_b(x_i - x_j - h)}{\sum_{i \neq j} k_b(x_i - x_j - h)}.$$

Both smoothed and histogram-based estimator can be computed in the package `geoR` using `variog`.

The estimators based on the squared differences  $(Z(x_i) - Z(x_j))^2$  are very sensitive to outlying observations because for them large values are squared making them even larger. Assume that  $\{Z(x) : x \in D\}$  is a Gaussian random field. Then  $(Z(x+h) - Z(x))^2$  has distribution  $2\gamma(h) \cdot \chi_1^2$ , which is very skewed. The fourth root is a suitable transformation that creates a distribution “close” to the normal distribution, see Figure 2. Instead of  $(Z(x_i) - Z(x_j))^2$  we can thus work with  $|Z(x_i) - Z(x_j)|^{1/2}$ . This leads us to the robust version of the variogram estimator:

$$2\bar{\gamma}(h) = \left( \frac{1}{|N(h)|} \sum_{N(h)} |Z(x_i) - Z(x_j)|^{1/2} \right)^4 / B(h),$$

where  $B(h) = 0.457 + 0.494/|N(h)|$ . The fourth power is there to preserve the proper scale. This transformation breaks the unbiasedness of the estimator, and so the term  $B(h)$  is added. This term represents the bias correction and ensures approximately unbiased estimator. The robust estimator is computed in the package `geoR` by the choice `estimator.type="modulus"` in `variog`. Except of reducing the influence of outliers another advantage of the robust estimator is that the summands are less correlated than in the case of the classical estimator (4).

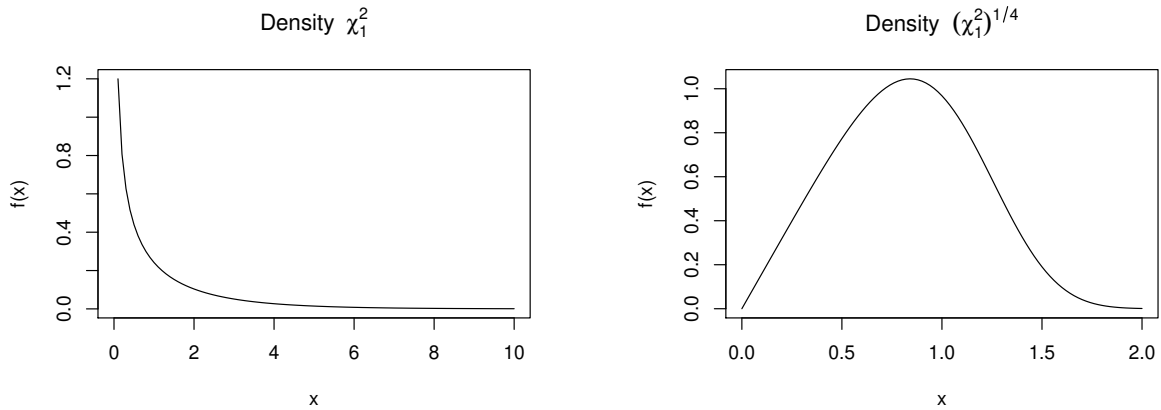


Figure 2: The density of a random variable  $X$  with  $\chi_1^2$ -distribution (left) and the density of a random variable  $X^{1/4}$  (right).

If we assume that the random field is weakly stationary, we may also work with the autocovariance function  $C(h) = \text{cov}(Z(x), Z(x+h))$ . In geostatistics, the term *covariogram* is usually used for the autocovariance function. Then there exists a relation between the semivariogram and covariogram,

$$\gamma(h) = C(o) - C(h). \tag{5}$$

The classical empirical estimator of the autocovariance function is

$$\hat{C}(h) = \frac{1}{|N(h)|} \sum_{N(h)} (Z(x_i) - \bar{Z})(Z(x_j) - \bar{Z}), \tag{6}$$

where the sample mean  $\bar{Z} = \frac{1}{n} \sum_{j=1}^n Z(x_j)$  estimates the mean  $\mu$ . The disadvantage is that we have to estimate  $\mu$  which causes bias of the estimator (6). For this reason, the variogram seems to provide better characterization of dependence than the autocovariance function. Hence, the variogram is often preferred to the covariogram in geostatistics. However, the autocovariance function is much more widely used in the classical statistics. The estimator of the autocovariance function is symmetric ( $\hat{C}(h) = \hat{C}(-h)$ ) and for  $h = o$  we have the estimator of variance:

$$\hat{C}(o) = \frac{1}{n} \sum_{i=1}^n (Z(x_i) - \bar{Z})^2.$$

Rewriting (4) so that we add and subtract  $\bar{Z}$  in each summand, we get

$$2\hat{\gamma}(h) = \frac{1}{|N(h)|} \sum_{N(h)} [(Z(x_i) - \bar{Z})^2 + (Z(x_j) - \bar{Z})^2] - 2\hat{C}(h).$$

Therefore,  $2\hat{\gamma}(h) \neq 2(\hat{C}(o) - \hat{C}(h))$ , and so the relation (5) is not preserved if we switch to the moment estimates. It would be unreasonable to estimate variogram by  $2(\hat{C}(o) - \hat{C}(h))$ , i.e. by plugging the sample covariances into (5) because negative values could be obtained.

Often we assume that the random field is also isotropic. Then the variogram is a function of distance  $\|h\|$ . We can exploit this in the construction of estimators. In the histogram-based estimator we can consider groups of pairs of points with similar mutual distances. In the kernel estimator we use  $k_b(\|x_i - x_j\| - \|h\|)$ , where  $k_b$  is a one-dimensional kernel function.

The disadvantage of non-parametric estimates is their larger variance and also the resulting estimators do not have to be valid variograms or covariograms. We know that every variogram must be conditionally negative definite and every covariogram must be positive semidefinite. However, the estimators  $\hat{\gamma}$  and  $\hat{C}$  do not necessarily have these properties. Therefore, we are now going to study parametric methods for estimation of variogram and autocovariance function.

### 3.1.2 Parametric methods

We select a parametric model for the variogram  $2\gamma_\theta(h)$  or the covariogram  $C_\theta(h)$ , where  $\theta \in \Theta$  is the vector of unknown parameters. For example, we may consider a power model for the variogram,

$$2\gamma_\theta(h) = c_0 + b\|h\|^\alpha, \quad \theta = (c_0, b, \alpha)^\top,$$

where  $c_0 \geq 0$  is the nugget,  $b \geq 0$  and  $0 \leq \alpha < 2$ . Our aim is to estimate  $\theta$  from data.

#### Least squares

The first possibility is a curve-fitting method of some non-parametric estimator computed in several values  $h_k$ ,  $k = 1, \dots, K$ . The simplest approach would be to minimize

$$\sum_{k=1}^K (2\hat{\gamma}(h_k) - 2\gamma_\theta(h_k))^2.$$

This is the ordinary least squares method. It disregards the correlations among the estimates  $2\hat{\gamma}(h_k)$  and their unequal variances. Put  $\mathbf{h} = (h_1, \dots, h_K)^\top$ ,  $2\hat{\gamma}(\mathbf{h}) = (2\hat{\gamma}(h_1), \dots, 2\hat{\gamma}(h_K))^\top$  and  $2\gamma_\theta(\mathbf{h}) = (2\gamma_\theta(h_1), \dots, 2\gamma_\theta(h_K))^\top$ , and consider the statistical model in the form

$$2\hat{\gamma}(\mathbf{h}) = 2\gamma_\theta(\mathbf{h}) + e(\mathbf{h}),$$

where we assume that  $e(\mathbf{h}) = (e(h_1), \dots, e(h_K))^\top$  has zero mean and variance matrix  $\mathbf{V}(\theta)$ , which may depend on  $\theta$ . Now we can apply the method of generalized least squares and minimize

$$(2\hat{\gamma}(\mathbf{h}) - 2\gamma_\theta(\mathbf{h}))^\top \mathbf{V}(\theta)^{-1} (2\hat{\gamma}(\mathbf{h}) - 2\gamma_\theta(\mathbf{h}))$$

w.r.t.  $\theta \in \Theta$ . The problem is how to obtain the matrix  $\mathbf{V}(\theta)$ .

Let  $\{Z(x) : x \in D\}$  be a Gaussian random field. Then

$$\mathbb{E}(Z(x_1 + h_1) - Z(x_1))^2 = 2\gamma(h_1) \quad \text{and} \quad \text{var}(Z(x_1 + h_1) - Z(x_1))^2 = 2(2\gamma(h_1))^2.$$

In order to express the covariance we use that  $\text{cov}(X^2, Y^2) = 2\rho^2$  holds for a random vector  $(X, Y)^\top$  with bivariate normal distribution such that  $\text{var} X = \text{var} Y = 1$  and  $\text{cov}(X, Y) = \rho$ . Hence,

$$\begin{aligned} \text{cov}((Z(x_1 + h_1) - Z(x_1))^2, (Z(x_2 + h_2) - Z(x_2))^2) &= 2(\gamma(x_1 - x_2 + h_1) + \gamma(x_1 - x_2 - h_2) \\ &\quad - \gamma(x_1 - x_2 + h_1 - h_2) - \gamma(x_1 - x_2))^2. \end{aligned}$$

The variance of the estimator (4) is

$$\begin{aligned}\text{var } 2\hat{\gamma}(h_k) &= \frac{1}{|N(h_k)|^2} \text{var} \sum_{N(h_k)} (Z(x_i) - Z(x_j))^2 \\ &= \frac{1}{|N(h_k)|^2} \sum_{i,j} \sum_{l,m} \text{cov}((Z(x_i) - Z(x_j))^2, (Z(x_l) - Z(x_m))^2).\end{aligned}$$

A simple approximation of this variance is

$$\text{var } 2\hat{\gamma}(h_k) \approx \frac{2(2\gamma_\theta(h_k))^2}{|N(h_k)|}. \quad (7)$$

This approximation is precise if  $(Z(x_i) - Z(x_j))^2$  are uncorrelated. We replace the matrix  $\mathbf{V}(\theta)$  by the diagonal matrix  $\mathbf{\Delta}(\theta)$  with elements given by the relation (7). Then we obtain the weighted sum of squares

$$(2\hat{\gamma}(\mathbf{h}) - 2\gamma_\theta(\mathbf{h}))^T \mathbf{\Delta}(\theta)^{-1} (2\hat{\gamma}(\mathbf{h}) - 2\gamma_\theta(\mathbf{h})) = \sum_{k=1}^K \frac{|N(h_k)|}{2\gamma_\theta(h_k)^2} (\hat{\gamma}(h_k) - \gamma_\theta(h_k))^2.$$

The estimator of  $\theta$  by the method of weighted least squares is obtained by the minimization of this sum.

### Maximum likelihood

The second possibility is to look for an estimator by the maximum likelihood method. For Gaussian random field with mean  $\mu$  and autocovariance function  $C_\theta$ , the log-likelihood based on data  $\mathbf{z}_n = (z(x_1), \dots, z(x_n))^T$  has after multiplying by  $-2$  this form:

$$-2 \log L(\mu, \theta) = n \log 2\pi + \log \det(\mathbf{C}_n(\theta)) + (\mathbf{z}_n - \mu \mathbf{1}_n)^T \mathbf{C}_n(\theta)^{-1} (\mathbf{z}_n - \mu \mathbf{1}_n),$$

where  $\mathbf{1}_n = (1, \dots, 1)^T$  and  $\mathbf{C}_n(\theta)_{ij} = C_\theta(x_i - x_j)$  depends on the vector  $\theta$  of covariance parameters. For given  $\theta$ ,  $L(\mu, \theta)$  is maximized for

$$\tilde{\mu} = (\mathbf{1}_n^T \mathbf{C}_n(\theta)^{-1} \mathbf{1}_n)^{-1} \mathbf{1}_n^T \mathbf{C}_n(\theta)^{-1} \mathbf{z}_n. \quad (8)$$

It is the generalized least squares estimator. Plugging  $\tilde{\mu}$  into  $L(\mu, \theta)$  we get the function of  $\theta$  (so-called *profile likelihood*), which has to be maximized (mostly by numerical methods). The estimator of  $\mu$  is then given by (8) with the estimate of  $\theta$  inserted.

A popular variant of maximum likelihood is REML – estimator by the method of *residual/restricted maximal likelihood*. This method does not apply the likelihood directly to data but to the residuals. It relies on finding an appropriate matrix  $\mathbf{A}$  which linearly transforms data  $\mathbf{Z}_n = (Z(x_1), \dots, Z(x_n))^T$  to  $\mathbf{Z}^* = \mathbf{A}\mathbf{Z}_n$  so that the distribution of  $\mathbf{Z}^*$  does not depend on  $\mu$ . The parameter  $\theta$  is then estimated by the maximum likelihood method applied to the transformed data  $\mathbf{Z}^*$ . The choice of matrix  $\mathbf{A}$  is not unique. For example, for matrix  $\mathbf{A}$  of type  $(n-1) \times n$  with entries  $a_{ij} = \mathbf{1}_{[i=j]} - 1/n$ , we get  $\mathbf{A}\mathbf{Z}_n = (Z(x_1) - \bar{Z}, \dots, Z(x_{n-1}) - \bar{Z})^T$  the vector of  $n-1$  differences from the sample mean  $\bar{Z}$ . In this way we get rid of dependence on  $\mu$ . The estimator of  $\theta$  minimizes the function

$$\log \det(\mathbf{A}\mathbf{C}_n(\theta)\mathbf{A}^T) + \mathbf{z}_n^T \mathbf{A}^T (\mathbf{A}\mathbf{C}_n(\theta)\mathbf{A}^T)^{-1} \mathbf{A}\mathbf{z}_n.$$

Plugging this estimator into (8) we get the estimator of  $\mu$ . In practice, the estimators may be determined by the functions `likfit` and `variofit` in the package `geoR` or `fitvario` in the package `RandomFields`.

### Composite likelihood

The composite likelihood method was already mentioned when dealing with parameter estimation in point processes. Similarly, it can be used for estimation of variogram parameters. Assume that the differences  $Z(x_i) - Z(x_j)$  have normal distribution. Summing the contributions of log-likelihood over pairs of distinct points we get the logarithm of composite likelihood:

$$\log \text{CL}(\theta) = \sum_{i,j=1,\dots,n}^{\neq} \left[ -\frac{1}{2} \log 4\pi\gamma_\theta(x_i - x_j) - \frac{1}{4\gamma_\theta(x_i - x_j)} (z(x_i) - z(x_j))^2 \right].$$

We are looking for  $\theta$  that maximizes  $\text{CL}(\theta)$ . So we differentiate w.r.t. components  $\theta_k$  and put equal to zero:

$$\sum_{i,j=1,\dots,n}^{\neq} \frac{\partial \gamma_\theta(x_i - x_j)}{\partial \theta_k} \frac{1}{4\gamma_\theta(x_i - x_j)^2} [(z(x_i) - z(x_j))^2 - 2\gamma_\theta(x_i - x_j)] = 0.$$

### 3.1.3 Model validation

Once we have chosen a parametric model of variogram and estimated its parameters, we are interested how well the obtained model  $2\gamma_{\hat{\theta}}$  describes the data. In the next subsection we will see how to obtain the prediction  $\hat{Z}(x_0)$  of  $Z(x_0)$  together with the prediction error  $\sigma^2(x_0)$ . It depends on the fitted variogram, data and the locations  $x_0, x_1, \dots, x_n$ . If we are able to get  $Z(x_0)$ , e.g., by additional measurement or from remaining data that we left for model validation, we may compare the difference between  $Z(x_0)$  and  $\hat{Z}(x_0)$ . These values should be close to each other if the variogram is chosen correctly.

If all data were used for variogram fitting and it is impossible to perform additional measurement, we can accomplish the *cross-validation*. We omit the location  $x_j$  and calculate the prediction  $\hat{Z}_{-j}(x_j)$  from the  $n - 1$  remaining observations and the fitted variogram  $2\gamma_{\hat{\theta}}(h)$ . The corresponding prediction error is denoted by  $\sigma_{-j}^2(x_j)$ . We perform this procedure for each  $j = 1, \dots, n$  and calculate the standardized residuals

$$\frac{Z(x_j) - \hat{Z}_{-j}(x_j)}{\sigma_{-j}(x_j)}.$$

Their arithmetic mean has to be around 0 and their sample second moment around 1. From the histogram of standardized residuals we can detect possible extreme values of residuals.

## 3.2 Kriging

Again we assume that the observed geostatistical data form a vector  $\mathbf{Z}_n = (Z(x_1), \dots, Z(x_n))^T$ . Our aim is to find the predictor  $\hat{Z}(x_0)$  of the unobservable value  $Z(x_0)$  that the random field attains at some further location  $x_0 \in D$ . The term *kriging* is used for the methods of spatial prediction based on the mean squared error minimization. It is named after a South African mining engineer D. G. Krige. His paper [5] dealing with mineral resources is a pioneering paper for the field of geostatistics.

### 3.2.1 Simple kriging

Let us assume that the random field has finite second moments. Then it is well-known that the mean squared error  $\mathbb{E}[Z(x_0) - \hat{Z}(x_0)]^2$  is minimized by the conditional expectation  $\mathbb{E}[Z(x_0) | \mathbf{Z}_n]$  and the minimum value is  $\mathbb{E}[Z(x_0) - \hat{Z}(x_0)]^2 = \mathbb{E} \text{var}[Z(x_0) | \mathbf{Z}_n]$ , see e.g., [6], Theorem 7.15. It means that  $\mathbb{E}[Z(x_0) | \mathbf{Z}_n]$  is the best predictor. The prediction error can be written as

$$\mathbb{E}[Z(x_0) - \hat{Z}(x_0)]^2 = \text{var } Z(x_0) - \text{var } \hat{Z}(x_0).$$

In practice, the conditional expectation is difficult to determine. Therefore, for simplicity, we restrict ourselves only to linear predictors of the form  $\hat{Z}(x_0) = \alpha + \beta^T \mathbf{Z}_n$ . Our aim is to determine  $\alpha \in \mathbb{R}$  and  $\beta \in \mathbb{R}^n$  so that the mean squared error is minimal. From the theory of linear models we know that the solution is

$$\beta_0 = \mathbf{C}_n^{-1} \mathbf{c}_n, \quad \alpha_0 = \mu(x_0) - \beta_0^T \mu_n,$$

where  $\mu_n = \mathbb{E} \mathbf{Z}_n = (\mu(x_1), \dots, \mu(x_n))^T$ ,  $\mu(x_0) = \mathbb{E} Z(x_0)$ ,

$$\mathbf{C}_n = (\text{cov}(Z(x_i), Z(x_j)))_{i,j=1,\dots,n}$$

is the variance matrix of  $\mathbf{Z}_n$ , and  $\mathbf{c}_n = (C(x_0, x_1), \dots, C(x_0, x_n))^T$ . Hence,

$$\hat{Z}(x_0) = \mu(x_0) + \mathbf{c}_n^T \mathbf{C}_n^{-1} (\mathbf{Z}_n - \mu_n)$$

and the prediction error is

$$\sigma^2(x_0) = \mathbb{E}(Z(x_0) - \hat{Z}(x_0))^2 = \text{var } Z(x_0) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n.$$

The technique for obtaining this spatial prediction is called the *simple kriging*. Even if we haven't required it, the predictor  $\hat{Z}(x_0)$  is unbiased in the sense that  $\mathbb{E} \hat{Z}(x_0) = \mathbb{E} Z(x_0)$ . Notice that the prediction error does not depend on the data. If  $x_0$  is one of the locations  $x_1, \dots, x_n$ , then  $\hat{Z}(x_0) = Z(x_0)$ , i.e. spatial prediction interpolates the data. To make sure that it holds, note that

$$\hat{Z}(x_j) = \mu(x_j) + (C(x_j, x_1), \dots, C(x_j, x_n))^T \mathbf{C}_n^{-1} (\mathbf{Z}_n - \mu_n), \quad j = 1, \dots, n,$$

which, rewritten for the vectors, becomes

$$(\hat{Z}(x_1), \dots, \hat{Z}(x_n))^T = \mu_n + \mathbf{C}_n \mathbf{C}_n^{-1} (\mathbf{Z}_n - \mu_n) = \mathbf{Z}_n.$$

The simple kriging predictor is optimal for Gaussian random fields.

**Lemma 3.1.** Let  $\{Z(x) : x \in D\}$  be a Gaussian random field. The best linear predictor  $\hat{Z}(x_0) = \mu(x_0) + \mathbf{c}_n^T \mathbf{C}_n^{-1} (\mathbf{Z}_n - \mu_n)$  is the best predictor of  $Z(x_0)$  and

$$Z(x_0) | \mathbf{Z}_n \sim N(\hat{Z}(x_0), \mathbb{E}(Z(x_0) - \hat{Z}(x_0))^2),$$

where  $\mathbb{E}(Z(x_0) - \hat{Z}(x_0))^2 = \text{var } Z(x_0) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n$ .

*Proof.* The joint distribution of  $(Z(x_0), \mathbf{Z}_n)^T$  is  $(n+1)$ -variate normal. Conditional distributions in a multivariate normal distribution are again normal. In our case the conditional distribution of  $Z(x_0) | \mathbf{Z}_n$  is normal with mean  $\mu(x_0) + \mathbf{c}_n^T \mathbf{C}_n^{-1} (\mathbf{Z}_n - \mu_n)$  and variance  $\text{var } Z(x_0) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n$ . The best (not necessarily linear) predictor of  $Z(x_0)$  is the conditional expectation  $\mathbb{E}[Z(x_0) | \mathbf{Z}_n]$ .  $\square$

The best linear predictor is optimal in the case of the Gaussian model. However, it may have bad properties when the assumption of the normal distribution is violated. In statistics, this problem is often settled up with a transformation of data leading to a normal distribution. An example is the so-called *Box-Cox transformation*

$$g_\lambda(z) = \begin{cases} \frac{z^\lambda - 1}{\lambda}, & \lambda \neq 0, \\ \log z, & \lambda = 0. \end{cases}$$

There exist different methods to select the parameter  $\lambda$ . It is also possible to follow the Bayesian approach and consider  $\lambda$  to be random.

We have expressed the best linear predictor. However, it depends on the values  $\mu_n$ ,  $\mu(x_0)$ ,  $\mathbf{c}_n$  and  $\mathbf{C}_n$ , which are unknown in practice. In general we have  $(n+1) + n + \binom{n+1}{2}$  unknown parameters that would have to be estimated from  $n$  observations. Therefore, we add some further assumptions. In next two subsections, we will consider a specific form for the mean. Then in Subsection 3.3 we discuss the influence of the estimation of covariances.

### 3.2.2 Ordinary kriging

Assume now that the random field has constant and finite mean  $\mu$ . We look for the linear predictor in the form

$$\hat{Z}(x_0) = \lambda^T \mathbf{Z}_n, \quad \text{where } \sum_{j=1}^n \lambda_j = \lambda^T \mathbf{1}_n = 1,$$

where the components  $\lambda_1, \dots, \lambda_n$  of the vector  $\lambda$  are unknown real coefficients. The condition that their sum is one ensures that the predictor is unbiased:  $\mathbb{E}\hat{Z}(x_0) = \lambda^T \mu \mathbf{1}_n = \mu = \mathbb{E}Z(x_0)$ . The method for finding the spatial prediction under these assumptions is named the *ordinary kriging*.

For intrinsically stationary random field with semivariogram  $\gamma$ , we can express the variance of a linear combination with zero sum of coefficients as follows:

$$\begin{aligned} \mathbb{E}(Z(x_0) - \hat{Z}(x_0))^2 &= \mathbb{E}(Z(x_0) - \lambda^T \mathbf{Z}_n)^2 = \text{var}(Z(x_0) - \lambda^T \mathbf{Z}_n) \\ &= - \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j \gamma(x_i - x_j) + 2 \sum_{i=1}^n \lambda_i \gamma(x_i - x_0). \end{aligned} \quad (9)$$

It means that we don't need to know  $\mu$  in order to determine the predictor  $\hat{Z}(x_0)$ . To find the minimum of (9) under the condition  $\lambda^T \mathbf{1}_n = 1$ , we can apply the method of Lagrange multipliers. For simpler notation we multiply the multiplier by 2 and minimize

$$Q = \text{var}(Z(x_0) - \lambda^T \mathbf{Z}_n) - 2m(\lambda^T \mathbf{1}_n - 1) = -\lambda^T \mathbf{\Gamma}_n \lambda + 2\lambda^T \boldsymbol{\gamma}_n - 2m(\lambda^T \mathbf{1}_n - 1),$$

where  $\mathbf{\Gamma}_n = (\gamma(x_i - x_j))_{i,j=1,\dots,n}$  and  $\gamma_n = (\gamma(x_1 - x_0), \dots, \gamma(x_n - x_0))^T$ . Differentiate  $Q$  w.r.t.  $\lambda$  and  $m$ , set the derivatives equal to zero, and obtain

$$\begin{aligned}\frac{\partial Q}{\partial \lambda} &= -2\mathbf{\Gamma}_n \lambda + 2\gamma_n - 2m\mathbf{1}_n = 0, \\ \frac{\partial Q}{\partial m} &= -2(\lambda^T \mathbf{1}_n - 1) = 0.\end{aligned}$$

The solution is

$$\begin{aligned}\lambda^T &= \left( \gamma_n + \mathbf{1}_n \frac{1 - \mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \gamma_n}{\mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \mathbf{1}_n} \right)^T \mathbf{\Gamma}_n^{-1}, \\ m &= -\frac{1 - \mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \gamma_n}{\mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \mathbf{1}_n}.\end{aligned}\tag{10}$$

Hence, the predictor has the form

$$\hat{Z}(x_0) = \left( \gamma_n + \mathbf{1}_n \frac{1 - \mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \gamma_n}{\mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \mathbf{1}_n} \right)^T \mathbf{\Gamma}_n^{-1} \mathbf{Z}_n = \lambda_1 Z(x_1) + \dots + \lambda_n Z(x_n).$$

The coefficients  $\lambda_i$  are components of the vector (10) and they are called the *prediction weights*. The prediction weights are typically large for points close to  $x_0$ . Nevertheless, their precise values depend on the locations  $x_i$  and the covariance structure of the data. It can happen that  $\lambda_i$  is negative or larger than 1. If  $x_0$  is one of the observed locations, say  $x_i$ , it is not difficult to see that  $m = 0$  and the prediction weights are  $\lambda_i = 1$  and  $\lambda_j = 0$  for  $j \neq i$ , i.e.  $\hat{Z}(x_0) = Z(x_0)$ . The prediction error is

$$\sigma^2(x_0) = \mathbb{E}(Z(x_0) - \hat{Z}(x_0))^2 = 2\lambda^T \gamma_n - \lambda^T \mathbf{\Gamma}_n \lambda = \gamma_n^T \mathbf{\Gamma}_n^{-1} \gamma_n - \frac{(1 - \mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \gamma_n)^2}{\mathbf{1}_n^T \mathbf{\Gamma}_n^{-1} \mathbf{1}_n}.$$

Similarly we can rewrite  $\hat{Z}(x_0)$  for weakly stationary random field using the autocovariance function:

$$\hat{Z}(x_0) = \left( \mathbf{c}_n + \mathbf{1}_n \frac{1 - \mathbf{1}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n}{\mathbf{1}_n^T \mathbf{C}_n^{-1} \mathbf{1}_n} \right)^T \mathbf{C}_n^{-1} \mathbf{Z}_n.$$

The prediction error is

$$\sigma^2(x_0) = \text{var } Z(x_0) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n + \frac{(1 - \mathbf{1}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n)^2}{\mathbf{1}_n^T \mathbf{C}_n^{-1} \mathbf{1}_n}.$$

We see that this error is larger than in the case of simple kriging because the last term is positive. The larger error is caused by the fact that we don't know the mean  $\mu$ .

### 3.2.3 Universal kriging

In this part we deal with the situation when the mean  $\mu(x) = \mathbb{E}Z(x)$  is not constant. The simplest approach is to use a linear model

$$\mu(x) = \sum_{j=0}^p \beta_j f_j(x),$$

where  $f_0(x), \dots, f_p(x)$  are known observed values of functions  $f_j$  in points  $x \in D$  and  $\beta_0, \dots, \beta_p$  are unknown real parameters. A usual choice for  $f_0$  is the constant function equal to 1, then  $\beta_0$  is an absolute term. For  $f_i(x)$  one can consider a polynomial of spatial coordinates of the location  $x$ . In this way, it is possible to model, for example, the linear trend. Another possibility is that  $f_i(x)$  represents some covariate. Denote  $\mathbf{f} = (f_0(x_0), \dots, f_p(x_0))^T$  and let  $\mathbf{F}$  be the matrix of type  $n \times (p+1)$  with elements  $f_j(x_i)$ ,  $i = 1, \dots, n$ ,  $j = 0, \dots, p$ . If we require the predictor in the form

$$\hat{Z}(x_0) = \lambda^T \mathbf{Z}_n, \quad \text{where } \lambda^T \mathbf{F} = \mathbf{f}^T,$$

we speak about the *universal kriging*. The condition  $\lambda^T \mathbf{F} = \mathbf{f}^T$  ensures that this predictor is unbiased because

$$\mathbb{E}\hat{Z}(x_0) = \lambda^T \mathbb{E}\mathbf{Z}_n = \lambda^T \mathbf{F}\beta = \mathbf{f}^T \beta = \mu(x_0) = \mathbb{E}Z(x_0).$$

The optimal predictor (minimizing the mean squared error) can be again found by applying the method of Lagrange multipliers. Analogously as in the case of ordinary kriging, we can show that the optimal prediction weights have the form

$$\lambda^T = (\gamma_n + \mathbf{F}(\mathbf{F}^T \mathbf{\Gamma}_n^{-1} \mathbf{F})^{-1}(\mathbf{f} - \mathbf{F}^T \mathbf{\Gamma}_n^{-1} \gamma_n))^T \mathbf{\Gamma}_n^{-1}.$$

The corresponding prediction error is

$$\gamma_n^T \mathbf{\Gamma}_n^{-1} \gamma_n + (\mathbf{f} - \mathbf{F}^T \mathbf{\Gamma}_n^{-1} \gamma_n)^T (\mathbf{F}^T \mathbf{\Gamma}_n^{-1} \mathbf{F})^{-1} (\mathbf{f} - \mathbf{F}^T \mathbf{\Gamma}_n^{-1} \gamma_n).$$

Using covariance, the prediction weights could be written as

$$\lambda^T = (\mathbf{c}_n + \mathbf{F}(\mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{F})^{-1}(\mathbf{f} - \mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{c}_n))^T \mathbf{C}_n^{-1}$$

and the prediction error as

$$\sigma^2(x_0) = C(o) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n + (\mathbf{f} - \mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{c}_n)^T (\mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{F})^{-1} (\mathbf{f} - \mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{c}_n).$$

By the generalized least squares, we can also estimate the parameter  $\beta$ :

$$\hat{\beta} = (\mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{Z}_n.$$

The predictor could be also written as

$$\hat{Z}(x_0) = \mathbf{f}^T \hat{\beta} + \mathbf{c}_n^T \mathbf{C}_n^{-1} (\mathbf{Z}_n - \mathbf{F} \hat{\beta}).$$

If  $Z(x_0)$  is uncorrelated with the data, the predictor  $\hat{Z}(x_0)$  coincides with the best linear unbiased estimator of the mean, which is equal to  $\mathbf{f}^T \hat{\beta}$ . However, the predictor of  $Z(x_0)$  is generally distinct from the estimator of  $\mathbb{E}Z(x_0)$ .

### 3.2.4 Other possibilities

Assume that instead of the prediction of  $Z(x_0)$ , we are interested in the prediction of the average value in some region (block)  $B$ ,

$$Z(B) = \frac{1}{|B|} \int_B Z(x) dx.$$

An analogy of the ordinary kriging leads to the so-called *block kriging*. We look for the predictor in the form

$$\hat{Z}(B) = \sum_{i=1}^n \hat{\lambda}_i Z(x_i),$$

where  $\sum_{i=1}^n \lambda_i = 1$ . The optimal prediction weights have the form

$$\lambda^T = \left( \mathbf{c}_B + \mathbf{1}_n \frac{1 - \mathbf{1}_n^T \mathbf{C}_n^{-1} \mathbf{c}_B}{\mathbf{1}_n^T \mathbf{C}_n^{-1} \mathbf{1}_n} \right)^T \mathbf{C}_n^{-1},$$

where  $\mathbf{c}_B = (\text{cov}(Z(B), Z(x_1)), \dots, \text{cov}(Z(B), Z(x_n)))^T$ . The expression using variogram would look analogously.

Similarly we may be interested in predicting  $g(Z(x_0))$ , where  $g$  is a given function. The best predictor is  $\mathbb{E}[g(Z(x_0)) | \mathbf{Z}_n]$ .

Another frequent case is the task of estimating probabilities  $\mathbb{P}(Z(x_0) \leq y | \mathbf{Z}_n)$ , where  $y$  is a given real number. We speak about the *indicator kriging*.

To perform the kriging techniques in R, we can use `krige.conv` in the package `geoR` or `krige` in the package `gstat`.



### 3.3 Influence of covariance parameters estimation

The formulas for spatial prediction derived in the previous subsection depend on the values of a covariogram or a variogram that are typically unknown in practice and so they must be somehow estimated. We have already mentioned basic approaches for the estimation of parameters for variogram or covariogram. We put the estimates of the parameters into the parametric formula of the corresponding function. In this way, we obtain the so-called *plug-in* estimators. The procedure goes in the following steps:

1. we select a parametric model for variogram  $\gamma_\theta(h)$  or covariogram  $C_\theta(h)$ ,
2. we estimate the parameter  $\theta$ ,
3. we adjust the statistical inference to take into account that instead of the constant  $\theta$  we work with the random variable  $\hat{\theta}$ .

The plug-in predictor for ordinary kriging has the form

$$\hat{Z}(x_0) = \left( \mathbf{c}_n(\hat{\theta}) + \mathbf{1}_n \frac{1 - \mathbf{1}_n^\top \mathbf{C}_n(\hat{\theta})^{-1} \mathbf{c}_n(\hat{\theta})}{\mathbf{1}_n^\top \mathbf{C}_n(\hat{\theta})^{-1} \mathbf{1}_n} \right)^\top \mathbf{C}_n(\hat{\theta})^{-1} \mathbf{Z}_n.$$

It is no longer the best linear unbiased predictor (BLUP) of  $Z(x_0)$ . It is just the estimator of this predictor (i.e. EBLUP = estimated best linear unbiased predictor). While the prediction error of  $\hat{Z}(x_0)$  is

$$C(o) - \mathbf{c}_n(\theta)^\top \mathbf{C}_n(\theta)^{-1} \mathbf{c}_n(\theta) + \frac{(1 - \mathbf{1}_n^\top \mathbf{C}_n(\theta)^{-1} \mathbf{c}_n(\theta))^2}{\mathbf{1}_n^\top \mathbf{C}_n(\theta)^{-1} \mathbf{1}_n}, \quad (11)$$

the prediction error of  $\hat{Z}(x_0)$  is unknown. If we plug  $\hat{\theta}$  into (11), we get the estimate of the prediction error of  $\hat{Z}(x_0)$ , i.e. of different predictor than we in fact use. This estimated prediction error has tendency to underestimate the true prediction error of  $\hat{Z}(x_0)$  because we neglect the fact that random  $\hat{\theta}$  introduces further variability into the EBLUP.

Return back to the case of universal kriging, where we consider the model  $Z(x) = \mathbf{F}(x)^\top \beta + e(x)$ , where  $\mathbf{F}(x) = (f_0(x), \dots, f_p(x))^\top$  and  $\{e(x) : x \in D\}$  is an intrinsically stationary random field with the variogram parameterized by  $\theta$ . It would be unreasonable to use empirical estimator of  $\theta$  from the data  $\mathbf{Z}_n = (Z(x_1), \dots, Z(x_n))^\top$  because it is substantially biased. The bias of (4) is caused by the fact that  $Z(x)$  does not have constant mean. Therefore,  $\mathbb{E}(Z(x_i) - Z(x_j))^2 = \text{var}(Z(x_i) - Z(x_j)) + (\mu(x_i) - \mu(x_j))^2$ . We would need a variogram estimator for  $\{e(x) : x \in D\}$ , but the error random field  $\{e(x) : x \in D\}$  is unobservable. If  $\beta$  was known, then  $e(x) = Z(x) - \mathbf{F}(x)^\top \beta$  and we would be able to estimate  $\theta$  from  $\mathbf{e}_n = (e(x_1), \dots, e(x_n))^\top$ . However, the parameter  $\beta$  is unknown. If the field  $\{e(x) : x \in D\}$  is weakly stationary with autocovariance function  $C_\theta(h)$ , we get the estimator of  $\beta$  using the method of generalized least squares,

$$\hat{\beta} = (\mathbf{F}^\top \mathbf{C}_n(\theta)^{-1} \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{C}_n(\theta)^{-1} \mathbf{Z}_n. \quad (12)$$

Nevertheless, this estimator requires knowledge of the parameter  $\theta$ . It means that we are not able to reasonably estimate  $\theta$  without knowledge of  $\beta$ . On the other hand, to estimate  $\beta$  we need an estimator of  $\theta$ . This situation is referred to as the cat-and-mouse-game of universal kriging.

A possible solution is the *iteratively re-weighted generalized least squares* method. It is defined by the following steps.

1. obtain an initial estimator of  $\beta$ , independent of  $\theta$ , e.g., by ordinary least squares method:  $\hat{\beta} = (\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{Z}_n$ ,
2. calculate the residuals  $\mathbf{r} = \mathbf{Z}_n - \mathbf{F} \hat{\beta}$ ,
3. estimate a parametric model of the variogram or covariogram of residuals and obtain  $\hat{\theta}$ ,
4. determine the new estimator  $\hat{\beta}$  as  $\hat{\beta} = (\mathbf{F}^\top \mathbf{C}_n(\hat{\theta})^{-1} \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{C}_n(\hat{\theta})^{-1} \mathbf{Z}_n$ ,
5. repeat steps 2.–4. until relative changes of the estimators  $\beta$  and  $\theta$  are small.

The variogram estimator is biased though the bias is not caused by a non-constant mean but by estimating the variogram of residuals and not the variogram of  $\{e(x) : x \in D\}$ .

A theoretical study of this procedure is intricate. It is not assured that the estimates converge to the theoretical parameters.

Another possibility is to use the maximum likelihood method to estimate both  $\beta$  and  $\theta$  simultaneously. For example, for a Gaussian random field  $\{Z(x) : x \in D\}$  the log-likelihood has the form

$$\log L(\beta, \theta) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \log \det \mathbf{C}_n(\theta) - \frac{1}{2} (\mathbf{z}_n - \mathbf{F}\beta)^T \mathbf{C}_n(\theta)^{-1} (\mathbf{z}_n - \mathbf{F}\beta).$$

For a fixed  $\theta$ , this function is maximized for  $\beta$  given by (12) with the vector  $\mathbf{Z}_n$  replaced by the observed data  $\mathbf{z}_n$ . Substituting this to  $\log L(\beta, \theta)$ , we get a function of  $\theta$  (profile likelihood), which has to be maximized numerically.

### 3.4 Bayesian approach

In the classical approach, the best predictor based on the observed data  $\mathbf{z}_n$  is  $\mathbb{E}[Z(x_0) | \mathbf{Z}_n = \mathbf{z}_n]$  and its error is  $\mathbb{E}\text{var}[Z(x_0) | \mathbf{Z}_n]$ . Often we are rather interested in the whole conditional distribution of  $Z(x_0)$  given  $\mathbf{Z}_n = \mathbf{z}_n$  than only in its mean and variance. This distribution is known as the *predictive distribution*. In the Bayesian approach, the predictive distribution is equal to the posterior distribution of  $Z(x_0)$ .

Recall that in Bayesian statistics, the parameters are considered to be random. It means that there is no difference between prediction and parameter estimation. The Bayesian approach is based on the combination of historical information about the unknown parameters  $\theta$  and observed data  $\mathbf{z}_n$ . Information about the parameters is given in the so-called *prior distribution* with density  $p(\theta)$  w.r.t.  $\sigma$ -finite measure  $\nu$  on the parametric space  $\Theta$ . Let  $\mathbf{Z}_n$  given  $\theta$  have a density  $f(\mathbf{z}_n | \theta)$ . Then the *posterior distribution* of  $\theta$  given  $\mathbf{Z}_n = \mathbf{z}_n$  is given by the Bayes theorem

$$p(\theta | \mathbf{z}_n) = \frac{f(\mathbf{z}_n | \theta)p(\theta)}{\int_{\Theta} f(\mathbf{z}_n | \theta)p(\theta) \nu(d\theta)},$$

provided that the denominator is positive. This relation is shortly written as

$$p(\theta | \mathbf{z}_n) \propto f(\mathbf{z}_n | \theta)p(\theta). \quad (13)$$

The symbol  $\propto$  denotes equality up to a multiplicative constant.

Spatial prediction using Bayesian approach is denoted as the *Bayesian kriging*. For the prediction of  $Z(x_0)$  we get the *predictive density* by integrating over  $\theta$ :

$$f(z_0 | \mathbf{z}_n) = \int_{\Theta} f(z_0, \theta | \mathbf{z}_n) \nu(d\theta) = \int_{\Theta} f(z_0 | \mathbf{z}_n, \theta)p(\theta | \mathbf{z}_n) \nu(d\theta). \quad (14)$$

For known  $\theta$ , the result is the same as in the classical approach. The advantage of the Bayesian approach is that it takes the uncertainty about model parameters into consideration. The form (14) of predictive density is mostly quite complicated. Therefore, the MCMC method are used. They enable to generate a sequence  $\theta^{(1)}, \dots, \theta^{(T)}$  from the posterior distribution with density  $p(\theta | \mathbf{z}_n)$ . Then the average

$$\hat{f}(z_0 | \mathbf{z}_n) = \frac{1}{T} \sum_{i=1}^T f(z_0 | \mathbf{z}_n, \theta^{(i)})$$

gives an approximation of the predictive density  $f(z_0 | \mathbf{z}_n)$ . In practice the calculation of this approximation is usually accomplished in the following way. For each  $\theta^{(i)}$  generate  $z_0^{(i)}$  from the distribution with density  $f(z_0 | \mathbf{z}_n, \theta^{(i)})$ . Then  $z_0^{(1)}, \dots, z_0^{(T)}$  is a sample from the predictive distribution and depiction of the corresponding histogram or the kernel density estimator gives an approximate shape of the predictive density. Another possible approach for the determination of the predictive density is at hand when we are able to calculate the posterior density  $p(\theta | \mathbf{z}_n)$  and  $p(\theta | \mathbf{z}_n, z_0)$ . Then we can exploit the relation

$$f(z_0 | \mathbf{z}_n) = f(z_0 | \mathbf{z}_n, \theta) \frac{p(\theta | \mathbf{z}_n)}{p(\theta | \mathbf{z}_n, z_0)}.$$

*Example:* Consider the linear model

$$Z(x) = \mathbf{F}(x)^T \beta + e(x), \quad x \in D,$$

where  $\mathbf{F}(x) = (f_0(x), \dots, f_p(x))^T$  is the vector of covariates,  $\beta = (\beta_0, \dots, \beta_p)^T$  is the vector of regression parameters with prior distribution  $N_{p+1}(\mathbf{m}, \mathbf{Q})$ , and  $\{e(x) : x \in D\}$  is a weakly stationary zero mean Gaussian random field with autocovariance function  $C(h)$ . Assume that we know the vector  $\mathbf{m}$ , matrix  $\mathbf{Q}$  and function  $C$ . Our aim is to find a spatial prediction of  $Z(x_0)$  based on the data  $\mathbf{Z}_n = (Z(x_1), \dots, Z(x_n))^T$ . Denote  $\mathbf{C}_n$  the matrix with elements  $C(x_i - x_j)$ ,  $i, j = 1, \dots, n$ , and  $\mathbf{F}$  the matrix of type  $n \times (p+1)$  with elements  $f_j(x_i)$ ,  $i = 1, \dots, n$ ,  $j = 0, \dots, p$ . Further assume that both  $\mathbf{Q}$  and  $\mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{F}$  have full rank. Since the normal distribution is a conjugate prior for a normally distributed data, the posterior distribution is a multivariate normal distribution. More precisely,  $\beta \mid \mathbf{Z}_n$  is distributed according to  $N_{p+1}(\mathbf{m}^*, \mathbf{Q}^*)$ , where

$$\mathbf{m}^* = (\mathbf{Q}^{-1} + \mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{F})^{-1} (\mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{Z}_n + \mathbf{Q}^{-1} \mathbf{m}), \quad \mathbf{Q}^* = (\mathbf{Q}^{-1} + \mathbf{F}^T \mathbf{C}_n^{-1} \mathbf{F})^{-1}.$$

The joint distribution of  $(\mathbf{Z}_n, Z(x_0))^T$  is multivariate normal  $N_{n+1}(\mathbf{F}_{n0} \beta, \mathbf{C}_{n0})$ , where

$$\mathbf{F}_{n0} = \begin{pmatrix} \mathbf{F} \\ \mathbf{F}(x_0)^T \end{pmatrix}$$

and

$$\mathbf{C}_{n0} = \begin{pmatrix} \mathbf{C}_n & \mathbf{c}_n \\ \mathbf{c}_n^T & C(o) \end{pmatrix},$$

$\mathbf{c}_n = (C(x_0 - x_1), \dots, C(x_0 - x_n))^T$ . The predictive density may be obtained from the expression

$$f(z_0 \mid \mathbf{z}_n) = \int f(z_0 \mid \mathbf{z}_n, \beta) p(\beta \mid \mathbf{z}_n) d\beta,$$

where  $p(\beta \mid \mathbf{z}_n)$  is the density of  $N_{p+1}(\mathbf{m}^*, \mathbf{Q}^*)$  and  $f(z_0 \mid \mathbf{z}_n, \beta)$  is the normal density with mean  $\mathbf{F}(x_0)^T \beta + \mathbf{c}_n^T \mathbf{C}_n^{-1} (\mathbf{z}_n - \mathbf{F} \beta)$  and variance  $C(o) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n$ , as we know from Lemma 3.1. After straightforward (though somewhat lengthy) calculation, we find out that the predictive distribution is normal with mean

$$(\mathbf{F}(x_0)^T - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{F}) \mathbf{Q}^* \mathbf{Q}^{-1} \mathbf{m} + [\mathbf{c}_n^T \mathbf{C}_n^{-1} + (\mathbf{F}(x_0)^T - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{F}) \mathbf{Q}^* \mathbf{F}^T \mathbf{C}_n^{-1}] \mathbf{Z}_n$$

and variance

$$C(o) - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{c}_n + (\mathbf{F}(x_0)^T - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{F}) \mathbf{Q}^* (\mathbf{F}(x_0)^T - \mathbf{c}_n^T \mathbf{C}_n^{-1} \mathbf{F})^T.$$

In practice, we don't know the function  $C$ . However, we may use some of the parametric models (e.g., Whittle–Matérn). Then we specify an appropriate prior distribution for parameters of the autocovariance function and derive the corresponding posterior distribution.

Geostatistical models that we have considered could be understood as two-stage *hierarchical models*. In the first stage of hierarchy, we model the dependence of the data on random effects. Specifically, a random field  $\mathbf{Z} = \{Z(x) : x \in D\}$  is prescribed conditionally on  $\mathbf{e} = \{e(x) : x \in D\}$ . In the second stage of hierarchy we model the distribution of a random effect, i.e. the unobserved random field  $\mathbf{e}$ .

In the Bayesian approach, we have three basic random objects, apart from  $\mathbf{Z}$  and  $\mathbf{e}$  it is also the vector  $\theta$  of unknown parameters. We get a three-stage hierarchical model:

1.  $\mathbf{Z} \mid \theta, \mathbf{e}$ ,
2.  $\mathbf{e} \mid \theta$ ,
3.  $\theta$ .

A particular example is the model described at the universal kriging and used also in the previous example:

$$Z(x) = \mathbf{F}(x)^T \beta + e(x).$$

Assume that the residual random field  $\{e(x) : x \in D\}$  is a centred stationary Gaussian random field and it can be written as the sum of a spatial component and a white noise:

$$e(x) = W(x) + \epsilon(x),$$

where  $\mathbf{W} = \{W(x) : x \in D\}$  is a centred stationary Gaussian random field with the autocovariance function  $C_W(h; \sigma^2, \phi) = \sigma^2 \rho(h; \phi)$  and  $\{\epsilon(x) : x \in D\}$  are uncorrelated random variables having normal distribution with zero mean and variance  $\tau^2$ . It means that the semivariogram of the random field  $\mathbf{e}$  is

$$\gamma_e(h; \sigma^2, \tau^2, \phi) = \tau^2 \mathbf{1}_{[h \neq 0]} + \sigma^2 (1 - \rho(h; \phi)).$$

It is parameterized by the nugget  $\tau^2$ , the partial sill  $\sigma^2$  and the correlation parameter  $\phi$  which appears in the autocorrelation function  $\rho(h; \phi)$  of the random field  $\mathbf{W}$ . The vector of unknown parameters is thus  $\theta = (\beta^T, \sigma^2, \tau^2, \phi)^T$ . Then  $\mathbf{Z}$  given  $\theta$  and  $\mathbf{W}$  is a Gaussian random field with mean  $\mathbf{F}(x)^T \beta + W(x)$  and autocovariance function  $C_{Z|W}(h; \tau^2) = \tau^2 \mathbf{1}_{[h=0]}$ . Notice that  $\mathbf{Z} | \theta, \mathbf{W}$  does not depend on  $\sigma^2$  and  $\phi$  at all. In the second stage of hierarchy, we specify  $\mathbf{W}$  that conditionally on  $\theta$  is a centred stationary Gaussian random field with the autocovariance function  $C_W(h; \sigma^2, \phi)$ , i.e.  $\mathbf{W}$  does not depend on  $\beta$  and  $\tau^2$ . The third stage requires determination of a suitable prior distribution for  $\theta$ . The graphical illustration of this hierarchical model is shown in Figure 3. Usually the components of  $\theta$  are taken to be a priori independent, i.e. the prior density is

$$p(\theta) = p(\beta)p(\sigma^2)p(\tau^2)p(\phi).$$

Appropriate candidates for the choice of marginal prior distributions are multivariate normal distribution for  $\beta$ , inverse  $\Gamma$ -distribution for  $\sigma^2$  and  $\tau^2$  (i.e.  $1/\sigma^2$  and  $1/\tau^2$  have  $\Gamma$ -distribution). The choice for  $\phi$  certainly depends on the form of the variogram, e.g., for the exponential model  $\rho(h; \phi) = \exp\{-\phi \|h\|\}$  a prior distribution for  $\phi$  is often taken to be  $\Gamma$ . The described model may be also formulated as two-stage. We utilize that  $\mathbf{Z} | \theta$  is a Gaussian random field with mean  $\mathbf{F}(x)^T \beta$  and the autocovariance function

$$C_Z(h; \sigma^2, \tau^2, \phi) = \tau^2 \mathbf{1}_{[h=0]} + \sigma^2 \rho(h; \phi).$$

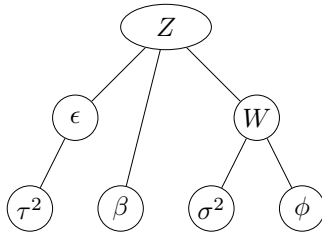


Figure 3: Representation of the three-stage hierarchical model.

Assume that we observe the vector  $\mathbf{z}_n = (z(x_1), \dots, z(x_n))^T$ . The Bayesian estimate of the parameter is then obtained from the posterior density  $p(\theta | \mathbf{z}_n)$  given by (13), where in this case  $f(\mathbf{z}_n | \theta)$  is the density of  $n$ -variate normal distribution with mean  $\mathbf{F}^T \beta$  and variance matrix  $\tau^2 \mathbf{I}_n + \sigma^2 \mathbf{H}(\phi)$ . Here,  $\mathbf{F} = (f_j(x_i))_{i,j}$  is the matrix of type  $n \times (p+1)$ ,  $\mathbf{I}_n$  is the identity matrix of size  $n$ , and  $\mathbf{H}(\phi)$  is the matrix of type  $n \times n$  with elements  $\rho(x_i - x_j; \phi)$ ,  $i, j = 1, \dots, n$ . We may as well use the three-stage hierarchical model and express the posterior density as

$$p(\theta, \mathbf{w}_n | \mathbf{z}_n) \propto f(\mathbf{z}_n | \theta, \mathbf{w}_n) p(\mathbf{w}_n | \theta) p(\theta),$$

where  $f(\mathbf{z}_n | \theta, \mathbf{w}_n)$  is the density of  $n$ -variate distribution with mean  $\mathbf{F}^T \beta + \mathbf{w}_n$  and variance matrix  $\tau^2 \mathbf{I}_n$  and  $p(\mathbf{w}_n | \theta)$  is the density of  $n$ -variate normal distribution with zero mean and variance matrix  $\sigma^2 \mathbf{H}(\phi)$ . However, in this way the number of parameters increases by  $n$  components of the vector  $\mathbf{w}_n = (w(x_1), \dots, w(x_n))^T$ . In practice, the MCMC methods (in particular, Gibbs sampler) are used. The form (13) is preferable because the variance matrix  $\tau^2 \mathbf{I}_n + \sigma^2 \mathbf{H}(\phi)$  behaves better than the variance

matrix  $\sigma^2\mathbf{H}(\phi)$ . This could be illustrated on the situation when the points  $x_i$  and  $x_j$  are close together. Then the matrix  $\sigma^2\mathbf{H}(\phi)$  is close to a singular matrix while  $\tau^2\mathbf{I}_n + \sigma^2\mathbf{H}(\phi)$  is not.

Estimation of the parameters  $\mathbf{w}_n$  corresponds to the reconstruction of the spatial surface  $\mathbf{W}$  in the measurement points  $x_1, \dots, x_n$ . Similarly, we can be interested in the prediction of  $W(x_0)$  for distinct choices of  $x_0$ . According to the relation

$$p(\mathbf{w}_n | \mathbf{z}_n) = \int \int p(\mathbf{w}_n | \sigma^2, \phi) p(\sigma^2, \phi | \mathbf{z}_n) d\sigma^2 d\phi,$$

we may obtain the posterior distribution of  $\mathbf{W}_n = (W(x_1), \dots, W(x_n))^T$  from the posterior distribution of  $(\sigma^2, \phi)$ . Recall that in our case,  $p(\mathbf{w}_n | \sigma^2, \phi)$  is the density of  $n$ -variate centred normal distribution with variance matrix  $\sigma^2\mathbf{H}(\phi)$ . Let  $((\sigma^2)^{(t)}, \phi^{(t)})$  be the output of MCMC algorithm which generates samples from the distribution with the posterior density  $p(\sigma^2, \phi | \mathbf{z})$ . Then it suffices to generate the vector  $\mathbf{w}_n^{(t)}$  from the distribution with the density  $p(\mathbf{w}_n | (\sigma^2)^{(t)}, \phi^{(t)})$  which gives us the output from the distribution with the density  $p(\mathbf{w}_n | \mathbf{z}_n)$ .

## 4 Lattice data

### 4.1 Modelling and estimation for areal data

By areal data, we mean that the aggregated values associated with some geographical regions (counties, districts, countries, etc.) are recorded. It is convenient to model such data using the random fields on a lattice. The sites of the lattice  $L$  correspond to individual regions. The neighbourhood relation  $\sim$  may be defined in such a way that two regions are in the relation  $\sim$  if and only if they share a common boundary. The areal data are often formed by the counts of a certain event (e.g., reported number of infected people, number of criminal acts). The modelling of discrete spatial data may be based on the generalized linear models.

Let  $\mathbf{Z} = \{Z_i : i \in L\}$  and  $\mathbf{W} = \{W_i : i \in L\}$  be random fields on the lattice  $L$ . Assume that conditionally on  $\mathbf{W}$ , the  $Z_i$  are independent random variables with mean  $\mathbb{E}(Z_i | \mathbf{W}) = \mu_i$ . Next consider the function  $h$  (so-called *link function*) and assume that

$$h(\mu_i) = \mathbf{F}_i^T \beta + W_i,$$

where  $\mathbf{F}_i = (f_{0i}, \dots, f_{pi})^T$  is a vector of region-specific covariates and  $\beta = (\beta_0, \dots, \beta_p)^T$  is a vector of regression parameters. This enables non-linear relationships between data and covariates. For binary data, the usual choice of  $h$  is the logit function  $h(\mu) = \log \frac{\mu}{1-\mu}$ . The random field  $\mathbf{W}$  models spatial variation. It captures spatial dependence present in the data. For example, we can use one of the Gaussian models (CAR, SAR, SMA, SARMA).

The most common approach when modelling the counts is to use a Poisson model. Assume that  $Z_i | \mathbf{W}$  has a Poisson distribution with parameter  $\theta_i E_i$ ,  $i \in L$ . Here,  $E_i$  is supposed to be known and it represents the expected number of events in the region  $i$ . As a link function we can use the logarithm and thus we obtain a linear model for  $\log \theta_i$ :

$$\log \theta_i = -\log E_i + \mathbf{F}_i^T \beta + W_i.$$

Spatial epidemiology is one of the main fields where this model is used. In this case,  $Z_i$  represents the observed number of cases of some disease in region  $i$  and  $E_i$  is the expected number of cases, which can be known from some additional information about the problem or may be some known function of  $n_i$  people at risk of the disease. For example, we can have  $E_i = rn_i$ , where  $r$  is the overall infection rate in the whole population. The rate  $r$  can be estimated by the ratio

$$\frac{\sum_{i \in L} Z_i}{\sum_{i \in L} n_i}.$$

This choice means that we expect the same infection rate in all regions. The value  $\theta_i$  is the region-specific relative risk. It gives the true relative risk of the infection in region  $i$ . As a covariate we can imagine, e.g., the level of air pollution which will have an important contribution when studying respiratory diseases.

We can also view the whole situation as the hierarchical model and use the Bayesian methods to make the statistical inference.

Now we consider a different model that specifies the joint distribution. Let  $\{Z_i : i \in L\}$  be a Markov random field with density  $p(\mathbf{z}; \theta)$  parameterized by a finite-dimensional vector  $\theta$ . For discrete data, the density is equal to the joint probabilities  $\mathbb{P}(Z_i = z_i, i \in L)$ ,  $\mathbf{z} = (z_i, i \in L)$ . For continuous data, it is the joint density w.r.t.  $n$ -dimensional Lebesgue measure.

The maximum likelihood method is one of the most popular statistical methods for estimating the parameters of a model. We find a value  $\hat{\theta}$  at which the likelihood function  $L(\theta) = p(\mathbf{z}; \theta)$  attains its maximum. Here,  $\mathbf{z} = (z_i, i \in L)$  are observed data.

If we have a Markov random field with the Gibbs distribution, then

$$L(\theta) = p(\mathbf{z}; \theta) = \exp \left\{ - \sum_{C \in \mathcal{C}} \Phi_C(\mathbf{z}_C, \theta) \right\} = \frac{\exp \left\{ - \sum_{C \in \mathcal{C}: C \neq \emptyset} \Phi_C(\mathbf{z}_C, \theta) \right\}}{\int \exp \left\{ - \sum_{C \in \mathcal{C}: C \neq \emptyset} \Phi_C(\mathbf{z}_C, \theta) \right\} \nu(d\mathbf{z})},$$

where  $\mathcal{C}$  is the system of cliques (subsets of  $L$  for which any two sites are neighbours). The problem is that the normalizing constant depends on  $\theta$  and it usually has a very complicated form. There exist methods for approximation of the normalizing constant employing simulations (mostly MCMC methods). Then we maximize this approximated likelihood function.

More similar procedure is to consider the so-called *pseudolikelihood*

$$L_P(\theta) = \prod_{i \in L} p(z_i | \mathbf{z}_{\partial i}; \theta) = \prod_{i \in L} \frac{\exp \left\{ - \sum_{C \in \mathcal{C}: C \neq \emptyset, i \in C} \Phi_C(\mathbf{z}_C, \theta) \right\}}{c(\mathbf{z}_{\partial i}, \theta)}.$$

The normalizing constant  $c(\mathbf{z}_{\partial i}, \theta)$  is often easier to express (in the discrete case it is the sum of  $|S|$  terms, where  $S$  is the state space of the random field). If we enumerate the elements of  $L$  by  $1, \dots, n$ , then the likelihood can be written as

$$L(\theta) = p(z_1 | z_2, \dots, z_n; \theta) p(z_2 | z_3, \dots, z_n; \theta) \cdots p(z_{n-1} | z_n; \theta) p(z_n; \theta).$$

Replacing the conditional densities  $p(z_k | z_{k+1}, \dots, z_n; \theta)$  by full conditional densities  $p(z_k | \mathbf{z}_{-k}; \theta)$ , which are equal to  $p(z_k | \mathbf{z}_{\partial k}; \theta)$  thanks to the Markov property, we obtain the pseudolikelihood  $L_P(\theta)$ .

The maximum pseudolikelihood estimator belongs to the class of estimators that are known in statistics as  $M$ -estimators. Generally, an  $M$ -estimator of  $\theta$  is obtained as the maximum of a contrast function  $\varrho(\mathbf{Z}, \theta)$ . In the classical situation of the maximum likelihood estimation for the sequence of i.i.d. random variables, we have

$$\varrho(\mathbf{z}, \theta) = \sum_{i=1}^n \log p(z_i; \theta).$$

In our case, we get

$$\varrho(\mathbf{z}, \theta) = \sum_{i \in L} \log p(z_i | \mathbf{z}_{\partial i}; \theta).$$

## 4.2 Testing of spatial autocorrelation

Recall that for a random field  $\mathbf{Z} = \{Z_i : i \in L\}$  with constant mean  $\mathbb{E}Z_i = \mu$  and constant variance  $\text{var } Z_i = \sigma^2$ , we have defined *Moran's I* by the relation

$$I = \frac{n}{w} \frac{\sum_{i \in L} \sum_{j \in L} w_{ij} (Z_i - \bar{Z})(Z_j - \bar{Z})}{\sum_{i \in L} (Z_i - \bar{Z})^2}$$

and *Geary's c* as

$$c = \frac{n-1}{2w} \frac{\sum_{i \in L} \sum_{j \in L} w_{ij} (Z_i - Z_j)^2}{\sum_{i \in L} (Z_i - \bar{Z})^2},$$

where  $w_{ij}$  are spatial proximity weights (i.e. we require  $w_{ij} = 0$  if  $i = j$  or  $i \not\sim j$ ) and  $w = \sum_{i \in L} \sum_{j \in L} w_{ij}$ . For the computation of Moran and Geary index one can use functions `moran` and `geary`, respectively, in

the R package `spdep`. These characteristics can be used as test statistics for testing the hypothesis of no spatial autocorrelation in data. Denote by  $M$  one of these test statistics (either Moran's  $I$  or Geary's  $c$ ) and by  $M_{obs}$  this test statistic computed from data.

Let us consider two different assumptions that correspond to the null hypothesis:

- a) randomness assumption: all  $n!$  permutations of observed values at  $n$  sites of  $L$  have equal probability  $1/n!$ ,
- b) Gaussianity assumption: random field  $\mathbf{Z}$  is formed by independent random variables with normal distribution  $N(\mu, \sigma^2)$ .

One of the following three approaches is usually used for testing under the randomness assumption.

1. *permutation test*: The null hypothesis  $H_0$  means that the observed values  $Z_i$ ,  $i \in L$ , are assigned completely at random. For  $n$  sites we have  $n!$  possible assignments. If we compute  $M$  for all  $n!$  possibilities, we obtain the distribution of  $M$  under  $H_0$ . Then we can determine the probability that the value  $M_{obs}$  is exceeded. Both large and small values of this probability indicate against  $H_0$  (if we consider two-sided test).
2. *Monte Carlo test*: Even if  $n$  is not very large, the corresponding number of permutations could be huge. Instead of calculating  $M$  for all permutations, we can generate  $k$  random permutations and construct the empirical distribution of  $M$  under  $H_0$ . Larger  $k$  means better approximation of the true distribution under  $H_0$ . We take together  $M_{obs}$  with  $k$  values of  $M$  from generated permutations and order them from the smallest to the largest. For extreme rank values of  $M_{obs}$  the null hypothesis should be rejected. For example, if  $k = 999$  we reject  $H_0$  on the level 5% when the rank of  $M_{obs}$  is between 1 and 25 or between 976 and 1000.
3. *asymptotic test*: Denote by  $\mathbb{E}_r M$  and  $\text{var}_r M$  the expectation and variance of  $M$  under  $H_0$  and randomness assumption, respectively. These first two moments can be determined analytically. Since one can often show the asymptotic normality of  $M$ , it suffices to compare

$$\frac{M_{obs} - \mathbb{E}_r M}{\sqrt{\text{var}_r M}}$$

with quantiles of the standard normal distribution  $N(0, 1)$ .

Under the assumption of Gaussianity, it is not difficult to express the expectation and variance of  $M$  when  $H_0$  holds. Denote these moments by  $\mathbb{E}_g M$  and  $\text{var}_g M$ , respectively. Again we can consider the asymptotic test and compare

$$\frac{M_{obs} - \mathbb{E}_g M}{\sqrt{\text{var}_g M}}$$

with quantiles of the standard normal distribution  $N(0, 1)$ .

It can be shown that  $\mathbb{E}_g I = \mathbb{E}_r I = -\frac{1}{n-1}$  and  $\mathbb{E}_g c = \mathbb{E}_r c = 1$ . The expectations are the same when assuming randomness and Gaussianity. However, the formulas for variances differ (see [2]). Moran's and Geary's statistics can be interpreted as follows: if  $I > \mathbb{E}I$  or  $c < \mathbb{E}c$ , then the lattice site has the tendency to be connected to the site with similar value of the random field. It corresponds to positive spatial autocorrelation. On the contrary, if  $I < \mathbb{E}I$  or  $c > \mathbb{E}c$ , the values at neighbouring sites have tendency to be dissimilar. Therefore, we can also consider one-sided variants of the test against the alternative that the spatial autocorrelation is positive (or negative).

## 5 Appendix

### 5.1 Random censoring

Assume that  $T_1, \dots, T_n$  are independent and identically distributed non-negative random variables with distribution function  $F$ . Our aim is to estimate  $F$ . If we observe all values of  $T_1, \dots, T_n$ , the most natural estimator of  $F$  is the empirical distribution function

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[T_i \leq t]}, \quad t \geq 0.$$

However, in some situations we don't have information about all  $T_i$ . In particular,  $T_i$  could represent the times to some event and their observation is prematurely interrupted. A typical example is the medical study of the influence of some treatment to the survival of a group of patients. Some of the observations are incomplete because the patient moved away or the time reserved for the study expired. Another example comes from the reliability theory where we measure the times to the breakdown of some product. Except of the random variables  $T_i$  (so-called *survival times* or *life times*) we also consider random variables  $C_1, \dots, C_n$  (so-called *censoring times*). We observe the random sample  $(\tilde{T}_1, D_1), \dots, (\tilde{T}_n, D_n)$ , where  $\tilde{T}_i = \min(T_i, C_i)$  are censored survival times and  $D_i = \mathbf{1}_{[T_i \leq C_i]}$  are indicators of non-censoring. For  $D_i = 1$  we observe the true time  $T_i$  while for  $D_i = 0$  the censoring happened and we have only partial information about  $T_i$ , namely  $T_i \geq \tilde{T}_i$ . In the case of random censoring we assume that  $C_1, \dots, C_n$  are i.i.d. random variables, independent of  $T_1, \dots, T_n$ . Then the non-parametric maximum likelihood estimator of  $F$  is the Kaplan–Meier estimator introduced in [4] and defined as

$$\hat{F}_{KM}(t) = 1 - \prod_{s \leq t} \left( 1 - \frac{\#\{i : \tilde{T}_i = s, D_i = 1\}}{\#\{i : \tilde{T}_i \geq s\}} \right).$$

The product effectively consists only of finitely many terms that correspond to the times  $s$  at which some life time is realized. The estimator  $\hat{F}_{KM}(t)$  is always a non-decreasing and right-continuous function. Its limit as  $t \rightarrow \infty$  could be strictly smaller than 1. This happens when the largest observed value is censored.

The intuitive explanation of the Kaplan–Meier estimator is the following. Divide the interval  $[0, t)$  into smaller intervals  $[0, t_1), [t_1, t_2), \dots, [t_k, t)$ . Then

$$1 - F(t) = \mathbb{P}(T_1 > t) = \mathbb{P}(T_1 > t \mid T_1 \geq t_k) \cdot \mathbb{P}(T_1 \geq t_k \mid T_1 \geq t_{k-1}) \cdots \mathbb{P}(T_1 \geq t_2 \mid T_1 \geq t_1) \cdot \mathbb{P}(T_1 \geq t_1),$$

where the conditional probabilities

$$\mathbb{P}(T_1 \geq t_j \mid T_1 \geq t_{j-1}) = 1 - \mathbb{P}(T_1 \in [t_{j-1}, t_j) \mid T_1 \geq t_{j-1})$$

could be estimated by

$$1 - \frac{\#\{i : \tilde{T}_i \in [t_{j-1}, t_j), D_i = 1\}}{\#\{i : \tilde{T}_i \geq t_{j-1}\}}.$$

Making the intervals  $[t_{j-1}, t_j)$  smaller, we get in the limit the expression for  $\hat{F}_{KM}(t)$ .

## References

- [1] A. BADDELEY AND R. TURNER (2005): Spatstat: an R package for analyzing spatial point patterns, *J. Stat. Softw.* **12**, 1–42.
- [2] A. D. CLIFF AND J. K. ORD (1981): *Spatial Processes; Models and Applications*, Pion Limited, London.
- [3] Y. GUAN (2006): Tests for independence between marks and points of a marked point process, *Biometrics* **62**, 126–134.
- [4] E. L. KAPLAN AND P. MEIER (1958): Nonparametric estimation from incomplete observations, *J. Amer. Statist. Assoc.* **53**, 457–481.
- [5] D. G. KRIGE (1951): A statistical approach to some basic mine valuation problems on the Witwatersrand, *J. Chem. Metal. Min. Soc. S. Afr.* **52**, 119–139.
- [6] P. LACHOUT (2004): *Teorie pravděpodobnosti*, second edition, in Czech, Karolinum, Praha.
- [7] J. MØLLER AND R. P. WAAGEPETERSEN (2003): *Statistical Inference and Simulation for Spatial Point Processes*, Chapman & Hall/CRC, Boca Raton.
- [8] M. MYLLYMÄKI, T. MRKVIČKA, P. GRABARNIK, H. SELJO AND U. HAHN (2017): Global envelope tests for spatial processes, *J. R. Statist. Soc. B* **79**, 381–404.



- [9] J. OHSER (1983): On estimators for the reduced second-moment measure of point processes, *Math. Operationsf. Statist., Ser. Statistics* **14**, 63–71.
- [10] E. J. PEBESMA (2004): Multivariable geostatistics in S: the gstat package, *Computers & Geosciences* **30**, 683–691.
- [11] P. J. RIBEIRO JR AND P. J. DIGGLE (2001): geoR: a package for geostatistical analysis, *R-NEWS* **1**, 15–18.
- [12] B. D. RIPLEY (1976): The second-order analysis of stationary point processes, *J. Appl. Probab.* **13**, 255–266.
- [13] M. SCHLATHER, P. J. RIBEIRO JR. AND P. J. DIGGLE (2004): Detecting dependence between marks and locations of marked point processes, *J. R. Statist. Soc. B* **66**, 79–93.