

MODEL FOR RANDOM UNION OF INTERACTING DISCS

KATEŘINA HELISOVÁ¹ AND JESPER MØLLER²

¹Charles University in Prague, Faculty of Mathematics and Physics, Department of Probability and Mathematical Statistics, Sokolovská 83, 18675 Prague 8, Czech Republic / Czech Technical University in Prague, Faculty of Electrical Engineering, Department of Mathematics, Technická 2, 16627 Prague 2, Czech Republic, ²Aalborg University, Department of Mathematical Sciences, Frederik Bajers Vej 7G, DK-9220 Aalborg, Denmark
e-mail: ¹helisova@karlin.mff.cuni.cz, helisova@math.feld.cvut.cz, ²jm@math.aau.dk

ABSTRACT

The contribution presents statistical analysis of data given by a digital image of heather growth. The bushes are modeled as a process of interacting discs. Different ways of estimating the parameters of the model are compared.

Keywords: Boolean model, interaction, random closed set, simulation-based maximum likelihood.

INTRODUCTION

Many phenomena in nature can be described by a union of discs. The objects represented by discs do not need to be independent, but they may mutually interact.

The contribution concerns a model of a random closed set given by a finite union of interacting discs with centers in a bounded set $S \subset \mathbf{R}^2$. This model is described by a density (with respect to a Boolean model), which depends on geometrical characteristics (e.g. area or perimeter) of the given set. The fact, that in applications, only the union and not the discs themselves can be observed, is taken into account so that all the characteristics depend only on the whole union. More details about this model can be found in (Møller and Helisová, 2008a)

Here, the focus is given on estimating the parameters of the model by MCMC simulation-based maximum likelihood approach (MCMC MLE), see (Møller and Waagepetersen, 2004), applied to heather data first presented in (Diggle, 1981). Problems with edge effects are solved by two different ways - conditional and unconditional MLE. The main aim is to compare these two ways for three different reference processes. For this comparing, some summary statistics (see (Stoyan et al, 1988)) and shape characteristics (see (Ripley, 1988)) are used.

The presented results are obtained mainly by using known methods from the theory of point processes. For more details about this analysis, see (Møller and Helisová, 2008b).

BASIC DEFINITIONS

Consider a point process \mathbf{X} defined on \mathbf{R}^d as a measurable mapping from some probability space (Ω, \mathcal{F}, P) to (N, \mathcal{N}) , where N is the system of locally finite subsets of \mathbf{R}^d with the σ -algebra $\mathcal{N} = \sigma(\{\mathbf{x} \in N : \#(\mathbf{x} \cap A) = m\} : A \in \mathcal{B}, m \in \mathbf{N}_0)$. The distribution $P_{\mathbf{X}}$ of \mathbf{X} is given by $P_{\mathbf{X}}(F) = P(\{\omega \in \Omega : X(\omega) \in F\})$ for $F \in \mathcal{N}$. We say that the point process \mathbf{X} is absolutely continuous with respect to the point process \mathbf{Y} if the distribution of \mathbf{X} is absolutely continuous with respect to the distribution of \mathbf{Y} .

Let \mathbf{Y} be the Poisson process with an intensity measure μ (i.e. the process satisfying that (a) for any finite collection $\{A_n\}$ of disjoint sets in \mathbf{R}^d , the numbers of points in these sets, $\mathbf{Y}(A_n)$, are independent random variables and (b) for each $A \subset \mathbf{R}^d$ such that $\mu(A) < \infty$, $\mathbf{Y}(A)$ has Poisson distribution with parameter $\mu(A)$) and denote $\Pi(F) = P(Y \in F)$ for $F \in \mathcal{N}$. A point process \mathbf{X} is given by density f with respect to the Poisson process \mathbf{Y} if

$$P(X \in F) = \int_F f(\mathbf{x}) \Pi(d\mathbf{x}).$$

MODEL

For the construction of the model, denote $b = b(z, r)$ a disc with a center $z \in \mathbf{R}^2$ and radius $r \in (0, \infty)$ and identify b with a point $x = (z, r) \in \mathbf{R}^2 \times (0, \infty)$. Then the union of discs $\cup_{i \in I} b_i = \cup_{i \in I} b(z_i, r_i)$, $I \subseteq \mathbf{N}$, can be identified with a point process on $\mathbf{R}^2 \times (0, \infty)$.

Consider a Poisson point process \mathbf{Y} on $\mathbf{R}^2 \times (0, \infty)$ with an intensity measure $\rho(z) dz Q(dr)$, where ρ is an intensity function of a Poisson point process on \mathbf{R}^2 and Q is a probability measure on $(0, \infty)$. Then the

process of discs corresponding to the point process \mathbf{Y} is a Boolean model with germs given by the Poisson point process with the intensity function ρ and grains are random discs with radii distribution Q .

Our model for a random set given by a union of interacting discs is the union corresponding to a point process \mathbf{X} which is absolutely continuous with respect to the Poisson process \mathbf{Y} and given by a density $f(\mathbf{x})$ with respect to \mathbf{Y} for any finite configuration $\mathbf{x} = \{x_1, \dots, x_n\}$.

We assume that \mathbf{X} is a finite point process defined on $S \times (0, R)$, where $S \subset \mathbf{R}^2$ is a bounded set such that $\int_S \rho(z) dz > 0$, and $R < \infty$. Further, we assume the density in the form

$$f_\theta(\mathbf{x}) = \frac{1}{c_\theta} \exp(\theta \cdot T(\mathcal{U}_\mathbf{x})) = \frac{\exp(\theta_1 A(\mathcal{U}_\mathbf{x}) + \theta_2 L(\mathcal{U}_\mathbf{x}) + \theta_3 N_{cc}(\mathcal{U}_\mathbf{x}) + \theta_4 N_h(\mathcal{U}_\mathbf{x}))}{c_\theta}, \quad (1)$$

where c_θ is a normalizing constant, $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)$ is a vector of parameters, \cdot denotes the inner product and $T(\mathcal{U}_\mathbf{x}) = (A(\mathcal{U}_\mathbf{x}), L(\mathcal{U}_\mathbf{x}), N_{cc}(\mathcal{U}_\mathbf{x}), N_h(\mathcal{U}_\mathbf{x}))$ is the vector of geometrical characteristics of the union $\mathcal{U}_\mathbf{x}$ of discs corresponding to the configuration \mathbf{x} , where $A(\mathcal{U}_\mathbf{x})$ denotes the area, $L(\mathcal{U}_\mathbf{x})$ the perimeter, $N_{cc}(\mathcal{U}_\mathbf{x})$ the number of connected components and $N_h(\mathcal{U}_\mathbf{x})$ the number of holes in the union $\mathcal{U}_\mathbf{x}$.

The interpretation of the density is following: If $\theta_1 = \dots = \theta_4 = 0$, there are no interactions among the discs and the model corresponds to the Boolean model. Else the configurations of the model have different geometrical characteristics than the reference Boolean model, e.g. if $\theta_1 > 0$, then the unions of discs with (in average) larger area than the area of the unions of discs of the reference process are more probable.

Because the interactions in the model depend on the vector of geometrical characteristics T , we call the model T -interaction process.

DATA ANALYSIS

We applied the model to the data of heather growth observed in a region 10×20 m in Jädraås (Sweden), see upper left picture in Figure 1.

Since the bushes are observed in a bounded window W , which does not include the whole growth, the problem with edge effects occurs. We solve this problem by two different ways of using MCMC MLE method. These ways are described in the following subsections.

CONDITIONAL MLE

Split \mathbf{X} into $\mathbf{X}^{(a)}$, $\mathbf{X}^{(b)}$, $\mathbf{X}^{(c)}$ corresponding to discs belonging to connected components of $\mathcal{U}_\mathbf{x}$ which are respectively

- (a) whole contained in the window W ,
- (b) intersecting both W and its complement W^c ,
- (c) whole contained in W^c .

Let $\mathbf{x}^{(b)}$ denote a realization of $\mathbf{X}^{(b)}$, i.e. $\mathbf{x}^{(b)}$ corresponds to a finite configuration of discs such that every component of $\mathbf{x}^{(b)}$ intersects both W and W^c . By (Møller and Helisová, 2008a), Proposition 5, we have that conditionally on $\mathbf{X}^{(b)} = \mathbf{x}^{(b)}$, the processes $\mathbf{X}^{(a)}$ and $\mathbf{X}^{(c)}$ are independent, and the conditional distribution of $\mathbf{X}^{(a)}$ depends on $\mathbf{x}^{(b)}$ only through $V = W \cap \mathcal{U}_{\mathbf{x}^{(b)}}$.

Denote \mathbf{Z} the whole data set, $\tilde{\mathbf{Z}} = \mathbf{Z} \cap W$ the data we can observe, $\mathbf{Z}^{(b)}$ the components intersecting the boundary of W and $\mathbf{Z}^{(a)} = \tilde{\mathbf{Z}} \setminus \mathbf{Z}^{(b)}$. Then we have the (conditional) log likelihood function in the form

$$L_c(\theta) = \theta \cdot T(\mathbf{Z}^{(a)}) - \log c_\theta. \quad (2)$$

Here, the edge effects are omitted. However, if W is not large enough, many components intersect the boundary and we can lose much data (in the worst case $\mathbf{x}^{(a)}$ is empty).

UNCONDITIONAL MLE

This way is based on ignoring everything outside the observation window W . Considering $S = W$, we can approximate the log likelihood function by

$$L_u(\theta) = \theta \cdot T(\tilde{\mathbf{Z}}) - \log c_\theta. \quad (3)$$

Here we have no data loss, but on the other hand, the method is less exact, since \mathbf{Z} (and hence $\mathcal{U}_\mathbf{x}$) may expand outside W .

Notice that the normalizing constant in (3) is different from the one in (2) - while c_θ in (3) corresponds to the density of the whole \mathbf{X} , c_θ in (2) is normalizing constant in conditional density of $\mathbf{X}^{(a)}$. Moreover, both these constants have no explicit expression, and therefore they need to be approximated using MCMC simulations.

RESULTS

From the data, we have for conditional MLE

$$A(\mathbf{Z}^{(a)}) = 45.6, L(\mathbf{Z}^{(a)}) = 190,$$

$$N_{cc}(\mathbf{Z}^{(a)}) = 32, N_h(\mathbf{Z}^{(a)}) = 2$$

and for unconditional MLE

$$A(\tilde{\mathbf{Z}}) = 100.3, L(\tilde{\mathbf{Z}}) = 382.8,$$

$$N_{cc}(\tilde{\mathbf{Z}}) = 56, N_h(\tilde{\mathbf{Z}}) = 6.$$

In order to have more comparing, we produce the estimates for three different reference processes (the choices of their parameters values are based on previous analyses):

(R1) $\rho = 2.45$ and Q is the restriction of distribution $N(0.26, 0.16^2)$ to the interval $[0, 0.50]$;

(R2) $\rho = 2.45$ and Q is the uniform distribution on $[0, 0.53]$;

(R3) $\rho = 1.16$ and Q is the uniform distribution on $[0, 0.53]$.

Comparing the data with the reference processes is in the following figure.

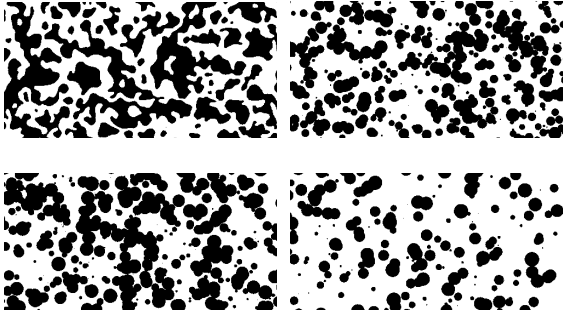


Fig. 1. Comparing the data (upper left) with simulation of reference Boolean models (R1) (upper right), (R2) (lower left) and (R3) (lower right).

The estimated parameters are further tested by Wald test (see (Møller and Waagepetersen, 2004)) if they can be considered to be equal to zero.

The estimates obtained by conditional likelihood (MLEc) and unconditional likelihood (MLEu) together with the corresponding p -values obtained from Wald test are shown in the following tables:

| MLEc | θ_1 | θ_2 | θ_3 | θ_4 |
|------------|------------|------------|------------|------------|
| (R1) | -2.14 | 0.89 | -1.78 | -1.01 |
| p -value | 0.0063 | 2.3e-05 | 2.6e-12 | 0.1435 |
| (R2) | -4.81 | 1.17 | -2.26 | -0.69 |
| p -value | 1.2e-09 | 4.9e-08 | < e-16 | 0.3149 |
| (R3) | -3.67 | 1.62 | -2.25 | -0.13 |
| p -value | 3.7e-05 | 8.4e-12 | < e-16 | 0.8415 |

| MLEu | θ_1 | θ_2 | θ_3 | θ_4 |
|------------|------------|------------|------------|------------|
| (R1) | -0.52 | -0.10 | -1.11 | -0.91 |
| p -value | 0.3149 | 0.5071 | 2.1e-09 | 0.0200 |
| (R2) | -3.32 | 0.72 | -1.62 | -0.49 |
| p -value | 1.0e-11 | 3.9e-07 | < e-16 | 0.2207 |
| (R3) | -1.79 | 1.04 | -1.64 | 0.01 |
| p -value | 0.0022 | 8.6e-10 | 2.2e-16 | 1 |

Since in the most cases, the parameter θ_4 has very large p -value, N_h seems to be irrelevant while the remaining characteristics with very low p -value are important. Therefore we omit N_h from the model and estimate again the parameters for the reduced model with the density

$$f_{\theta}(\mathbf{x}) = \frac{1}{c_{\theta}} \exp(\theta_1 A(\mathcal{U}_{\mathbf{x}}) + \theta_2 L(\mathcal{U}_{\mathbf{x}}) + \theta_3 N_{cc}(\mathcal{U}_{\mathbf{x}})).$$

Then we obtain

| MLEc | θ_1 | θ_2 | θ_3 |
|------------|------------|------------|------------|
| (R1) | -2.33 | 0.92 | -1.77 |
| p -value | 0.0020 | 4.6e-06 | 7.5e-12 |
| (R2) | -4.91 | 1.18 | -2.25 |
| p -value | 7.0e-10 | 1.3e-08 | < e-16 |
| (R3) | -3.71 | 1.64 | -2.25 |
| p -value | 3.7e-05 | 6.7e-12 | < e-16 |

| MLEu | θ_1 | θ_2 | θ_3 |
|------------|------------|------------|------------|
| (R1) | -0.91 | -0.02 | -1.13 |
| p -value | 0.0512 | 0.8875 | 3.7e-10 |
| (R2) | -1.75 | 1.02 | -1.63 |
| p -value | 0.0018 | 3.08e-10 | < e-16 |
| (R3) | -3.45 | 0.74 | -1.63 |
| p -value | 2.1e-12 | 2.2e-07 | < e-16 |

In five of the six models, all the estimates are considered not to be equal to zero. Hence we consider the values in the table to be the parameters of the final (A, L, N_{cc}) -interaction models.

MODEL CONTROL

Let $\mathbf{A} \subset \mathbf{R}^2$ be a set observed in a (bounded) window $W \subset \mathbf{R}^2$ (in our case, \mathbf{A} represents either the data $\tilde{\mathbf{Z}}$ or the set $\mathcal{U}_{\mathbf{X}}$ corresponding to the simulated disc process \mathbf{X} observed inside W) and G is a set of pixels in digital image of the set \mathbf{A} .

In order to compare the results of the methods, we construct the following pictures of simulations and plots of summary statistics and shape characteristics, where the three rows correspond from above to results for models with respect to (R1), (R2) and (R3), and the two columns correspond to models obtained by MLEc (left) and MLEu (right).

- Figure 2: simulations of the model for visual compare with the data;
- Figure 3: function $\hat{T}(r) = -\frac{1}{r} \log(1 - \hat{H}(r))$, $r > 0$, where \hat{H} is the estimate of spherical contact distribution function (see (Stoyan et al , 1988)) given by

$$\hat{H}(r) = \frac{\sum_{u \in G} \mathbf{1}[u \notin \mathbf{A}, u + b_r \subset W, (u + b_r) \cap \mathbf{A} \neq \emptyset]}{\sum_{u \in G} \mathbf{1}[u \notin \mathbf{A}, u + b_r \subset W]},$$

where $b_r = b(0, r)$ denotes a disc with center in $0 \in \mathbf{R}^2$ and radius r .

- Figure 4: estimate of covariance function $C(r) = P(u \in \mathbf{A}, v \in \mathbf{A})$ for any two points $u, v \in \mathbf{R}^2$ with distance $\|u - v\| = r$ given by

$$\hat{C}(r) = \frac{\sum_{u, v \in G} \mathbf{1}[\|u - v\| = r, \{u, v\} \subset \mathbf{A}]}{\sum_{u, v \in G} \mathbf{1}[\|u - v\| = r]}.$$

- Figure 5: dilatation (for more details, see (Ripley , 1988))

$$d(r) = \frac{|\mathbf{A}_{\oplus r} \cap W_{\ominus r}|}{|W_{\ominus r}|},$$

where $|\cdot|$ denotes the area, $\mathbf{A}_{\oplus r} = \cup_{u \in \mathbf{A}} b(u, r)$ is enlarging and $\mathbf{A}_{\ominus r} = \{u : b(u, r) \subseteq \mathbf{A}\}$ is contracting of the set \mathbf{A} by a disc with radius r .

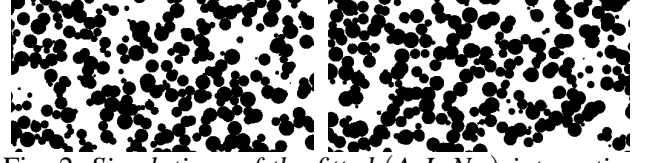
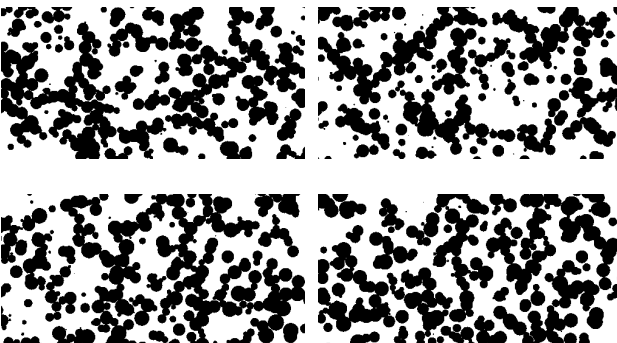


Fig. 2. Simulations of the fitted (A, L, N_{cc}) -interaction models.

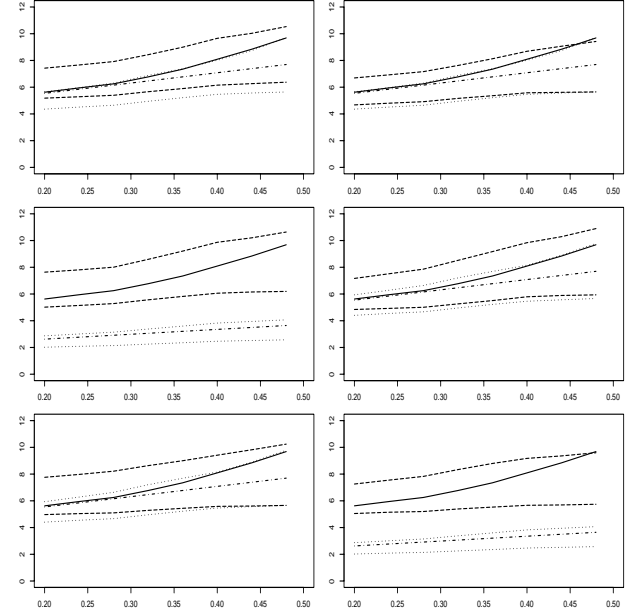


Fig. 3. Comparing the theoretical functions $T(r)$ for the reference Boolean models (dot-dashed lines) with $\hat{T}(r)$ based on the data (solid lines) and its simulated 2.5% and 97.5% envelopes obtained under the Boolean model (R1), (R2), or (R3) (dotted lines) and the corresponding (A, L, N_{cc}) -interaction model (dashed lines).

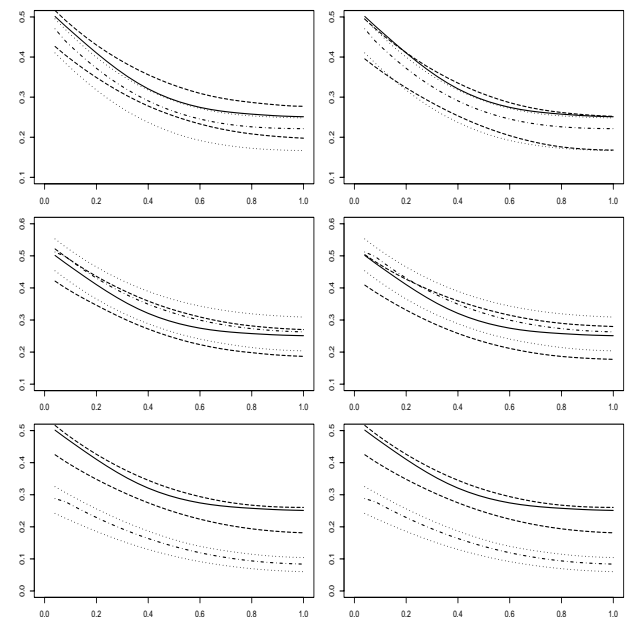


Fig. 4. Comparing the theoretical functions $C(r)$ for the reference Boolean models (dot-dashed lines) with $\hat{C}(r)$ based on the data (solid lines) and its simulated 2.5% and 97.5% envelopes obtained under the Boolean model (R1), (R2), or (R3) (dotted lines) and the corresponding (A, L, N_{cc}) -interaction model (dashed lines).

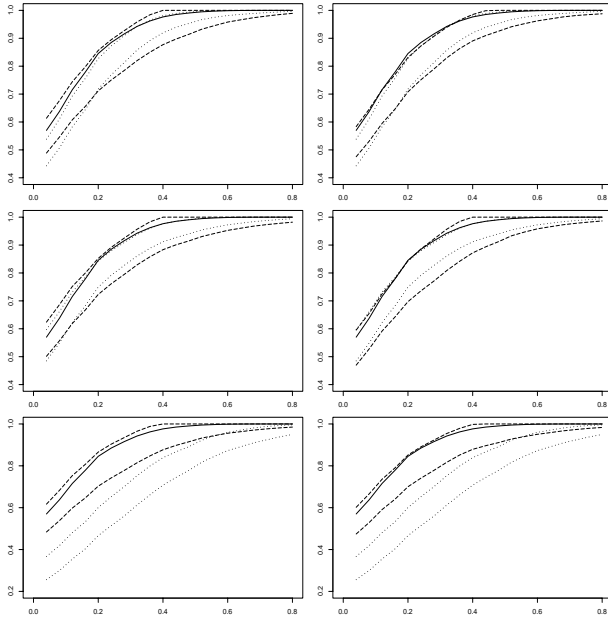


Fig. 5. Comparing $d(r)$ based on the data (solid lines) with simulated 2.5% and 97.5% envelopes obtained under the Boolean model (R1), (R2), or (R3) (dotted lines) and the corresponding (A, L, N_{cc}) -interaction model (dashed lines).

CONCLUSION

In Figures 3 - 5, one can see that in many cases, the plots of summary and shape statistics for the

data lie more in the middle of the envelopes for the models with parameters obtained by MLEc. Due to this observation, MLEc can be considered more exact than MLEu. Nevertheless, there are apparent misfits neither for models obtained by MLEc nor for that obtained by MLEu, and therefore we can conclude that both the methods provide suitable estimates.

ACKNOWLEDGEMENTS

The research was supported by the Danish Natural Science Research Council, grant 272-06-0442 "Point process modelling and statistical inference", by grant IAA 101120604 and by the Czech Government under the research programme MSM 6840770038.

REFERENCES

- Diggle PJ (1981). Binary mosaics and the spatial pattern of heather. *Biometrics* 37:531–539.
- Møller J and Helisová K (2008). Power diagrams and interaction processes for unions of discs. *Advances in Applied Probability* 40:321–347.
- Møller J and Helisová K (2008). Likelihood inference for unions of interacting discs. To appear.
- Møller J and Waagepetersen RP (2004). *Statistical Inference and Simulation for Spatial Point Processes*. Boca Raton: Chapman and Hall/CRC.
- Ripley BD (1988). *Statistical Inference for Spatial Processes*. Cambridge: Cambridge University Press.
- Stoyan D, Kendall WS and Mecke J (1995). *Stochastic Geometry and Its Applications*. Chichester: Wiley.