

# A stochastic–geometric model of the variability of soil formed in Pleistocene patterned ground.

R.M. Lark<sup>1\*</sup>, E. Meerschman<sup>2</sup>, M. Van Meirvenne<sup>2</sup>.

<sup>1</sup>*British Geological Survey, Keyworth, Nottinghamshire NG12 5GG, U.K.*, <sup>2</sup>*Research Group Soil Spatial Inventory Techniques, Department of Soil Management, Faculty of Bioscience Engineering, Ghent University, Coupure 653, 9000 Gent, Belgium.*

1

---

## 2 Abstract

3 In this paper we develop a model for the spatial variability of apparent electrical  
4 conductivity,  $EC_a$ , of soil formed in relict patterned ground. The model is based on the  
5 continuous local trend (CLT) random processes introduced by Lark (2012) (Geoderma,  
6 189–190, 661–670). These models are non-Gaussian and so their parameters cannot be  
7 estimated just by fitting a variogram model. We show how a plausible CLT model, and  
8 parameters for this model, can be found by the structured use of soil knowledge about  
9 the pedogenic processes in the particular environment and the physical properties of the  
10 soil material, along with some limited descriptive statistics on the target variable. This  
11 approach is attractive to soil scientists in that it makes the geostatistical analysis of soil  
12 properties an explicitly pedological procedure, and not simply a numerical exercise. We  
13 use this approach to develop a CLT model for  $EC_a$  at our target site. We then develop  
14 a test statistic which measures the extent to which soils on this site with small values  
15 of  $EC_a$ , which are coarser and so more permeable, tend to be spatially connected in the  
16 landscape. When we apply this statistic to our data we get results which indicate that  
17 the CLT model is more appropriate for the variable than is a Gaussian model, even after  
18 transformation of the data. The CLT model could be used to generate training images of  
19 soil processes to be used for computing conditional distributions of variables at unsampled  
20 sites by multiple point geostatistical algorithms.

21 *Keywords:* Patterned ground; Apparent electrical conductivity; Electromagnetic induc-  
22 tion; Stochastic geometry; Voronoi tessellation; Multiple point geostatistics; Pedometrics.

---

\*Corresponding author: *E-mail address:* mlark@nerc.ac.uk (R.M. Lark).

## 23 1. Introduction

24 *‘Mais surtout nous insisterons sur la nécessité d’incorporer au maximum la physique*  
25 *du problème et le contexte géologique de la zone étudiée’.* Chilès and Guillen (1984).

26 In most geostatistical analyses of soil the data are assumed to be a realization of a  
27 multi-Gaussian random function, perhaps after they have been transformed so that their  
28 histogram represents a Gaussian distribution. Furthermore, the random function com-  
29 monly has a spatial covariance function drawn from a limited subset of models (Webster  
30 and Oliver, 2007), which are used because of their convenient mathematical properties. In  
31 some of the earth sciences there has been progress in the development of random functions  
32 with parameters that are determined, or at least constrained, by parameters of underlying  
33 processes which have a physical meaning (e.g. Kolvos et al., 2004; Chilès and Guillen,  
34 1984). This has advantages (Lark, 2012a), for example, the efficiency of spatial sampling  
35 to model the spatial covariance function could be improved if prior distributions for co-  
36 variance parameters could be specified from process knowledge. However, this has not  
37 been achieved in soil science. Lark (2012a) suggested that this is probably because the  
38 variables that soil scientists study are commonly influenced by a more complex set of fac-  
39 tors at more diverse spatial scales than is the case for the variables where it has proved  
40 possible to specify the covariance function from process information. For example, the  
41 covariance function for diffusion processes is well-established (Whittle, 1954; 1962), and  
42 diffusion is a source of spatial variation in the concentration of nutrients in soil, but it is  
43 just one of many sources of spatial variation, and is of limited importance at the spatial  
44 scales most generally studied for practical purposes.

45 Lark (2012a,2012b) suggested that progress might be made by recognizing a number  
46 of distinct *modes* of soil variation, simple and generalizable rules that capture how the  
47 effects of factors of soil variation vary laterally, and which map naturally on to particular  
48 spatial random functions. For example, in conditions where soil variation is strongly  
49 determined by differences between discrete domains in the landscape (such as geological  
50 units, topographic units, fields etc.) then a subdivision of space into random sets such as  
51 Poisson Voronoi polygons may be appropriate (Lark, 2009) and properties of the spatial  
52 model (such as the mean chord length of the polygons) may be given a physical meaning.

53 Lark (2012b) proposed a mode of soil variation: continuous local trends. Under this  
54 mode of variation soil varies laterally in space, changing continuously rather than in a  
55 step-wise fashion; and these trends are local and repeating, so that they are essentially  
56 unpredictable (in contrast to a large-scale trend in a variable that might be observed across  
57 a study area). Examples of continuous local trends would be concentration gradients  
58 around the rhizosphere, or around individual plants, and catenary variation at landscape  
59 scale. Lark (2012b) proposed a general family of random functions to describe continuous  
60 local trends (CLT random functions). The value of a CLT variable at some location is given  
61 by a distance function, whose argument is the distance from the location of interest to the  
62 nearest event in a realization of a spatial point process. This makes the CLT a random  
63 function. The CLT variables considered by Lark (2012b), and in this paper, are Poisson  
64 CLT (PCLT) variables because the spatial point process is completely spatially random.  
65 Lark (2012b) estimated parameters of a PCLT process from data on a soil variable. It  
66 was also pointed out that the PCLT process might differ from a comparable Gaussian  
67 random function with respect to its multiple point statistics (Strebelle, 2002). This raises  
68 the possibility that PCLT models, as well as mapping closely on to a particular mode of  
69 soil variation, might be practically useful for applications where spatial connectivity plays  
70 a major role controlling processes in soil and so the multiple point statistics of the variable  
71 are important.

72 In this paper we use a PCLT random function to model the variation of apparent  
73 electrical conductivity,  $EC_a$ , of soil at a site where this variable is strongly influenced  
74 by spatial patterns in the parent material. These patterns arose from the development  
75 of ice wedges in Eocene clay under permafrost conditions, and subsequent infilling by  
76 coarser material which leads to strong textural contrasts in the soil. The objective is  
77 to show how soil knowledge: general knowledge about soil formation in the particular  
78 environment and its relationship to  $EC_a$ , and some simple descriptive statistics of the  
79 data (summary statistics and empirical variograms), allow us to select and fit a PCLT  
80 model. We then compare the PCLT model with a trans-Gaussian (TG) model of the  
81 data, i.e. a model fitted by conventional geostatistical analysis after the data have been  
82 transformed to approximate normality. Specifically we compare the models with respect  
83 to a statistic that summarizes the spatial connectivity of the coarser material, which might

84 be relevant to simulations of transport processes in the soil. We then evaluate which model  
85 appears best to represent the spatial pattern in the data.

## 86 **2. Case Study**

### 87 *2.1 The study area and data collection.*

88 We surveyed an area of Pleistocene patterned ground in the sandy silt region of  
89 Belgium. The patterned ground was identified by polygonal crop marks on an aerial  
90 photograph and interpreted to be the result of ice wedge formation during the last glacial  
91 period. The study area and data collection were discussed in detail by Meerschman et  
92 al. (2011), therefore we limit ourselves here to a brief presentation of it. More general  
93 information on ice-wedge polygons constitutes part of the soil-knowledge base that we use  
94 in this study, and is presented in section 2.3.2 below as it is required.

95 The study area (0.6 ha) was located in an agricultural field in Deinze, Belgium  
96 (central coordinates: 51° 01'16"N, 3°29'41"E). Excavation of a small part of the study  
97 area (6×6-m) to a depth of 0.9 m uncovered an ice-wedge pseudomorph with a diameter  
98 of about 6 m. The wedges were formed in clay-rich Tertiary marine sediments that were  
99 covered with a 0.6 m layer of silty-sand Quaternary deposits. Texture analysis on 94  
100 subsoil samples (0.6 - 0.8 m) showed a clear contrast between the Eocene host material  
101 (on average 21% clay) and the superficial material (on average 6% clay).

102 Previous studies (Saey et al., 2009; Cockx et al., 2006) have shown that  $EC_a$  is a  
103 useful covariate to study textural variability at profile and polygon-scale in soils formed  
104 in these conditions. The study area was surveyed with a mobile proximal soil sensor  
105 measuring  $EC_a$  ( $mS\ m^{-1}$ ) of an underlying soil volume down to approximately 1.5 m. The  
106 sensor was mounted on a sled pulled by an all terrain vehicle. The vehicle drove along  
107 parallel lines with an in-between distance of on average 0.75 m. The within-line distance  
108 between sensor response registrations was 0.15 m.

### 109 *2.2 Initial data analysis.*

110 Meerschman et al. (2011) noted that the  $EC_a$  measurements clearly reflected the  
111 polygonal patterns: small  $EC_a$  values indicated the former ice wedges filled with lighter  
112 material. In addition to the short-range variation in  $EC_a$ , there were large values of  $EC_a$

113 near an old field track in the north-east of the surveyed region. To avoid any assumptions  
114 about the form of this trend we decided to restrict our analyses to the lower left quadrant  
115 of the surveyed area, a region of approximately 40×40-m, with 17 792 observations, which  
116 excludes this area with elevated EC<sub>a</sub>. Figure 1 shows a post-plot of these data.

117 Figure 2 shows the histogram of the data. Summary statistics are presented in Table  
118 1. Note that the data are mildly skewed. In the analyses reported below the PCLT model  
119 was fitted in all cases to the raw data, and all analyses with the TG model were done with  
120 the data after a transformation which is described in section 2.3.1 below.

### 121 *2.3 Spatial analysis.*

122 In this section we describe the analysis of the EC<sub>a</sub> data to fit a TG model and a PCLT  
123 model. The first task (section 2.3.1) was straightforward after a data transformation, which  
124 is described. In section 2.3.2 we describe how soil knowledge was used to fit the PCLT  
125 model.

#### 126 *2.3.1. Trans-Gaussian model*

127 The objective of the case study is to compare a continuous local trend (PCLT) model  
128 of the data with a trans-Gaussian (TG) model, as might be used in standard geostatistical  
129 analysis. Although the data are only mildly skew, since the objective of this exercise is to  
130 compare a Gaussian or Trans-Gaussian model with a stochastic geometric alternative, it  
131 was decided to transform the data so that the histogram and summary statistics were as  
132 close as possible to those expected for data drawn from a Gaussian random variable. We  
133 therefore used a Box-Cox transformation of the data to normality for the TG modelling:

$$\begin{aligned} y &= \frac{z^\zeta - 1}{\zeta} \quad \zeta \neq 1, \\ &= \log_e(z) \quad \zeta = 1, \end{aligned} \tag{1}$$

134 where  $z$  is a value on the original scale and  $y$  is a transformed value. We used the BOXCOX  
135 procedure from the MASS package (Venables and Ripley, 2002) for the R platform (R  
136 Development Core Team, 2012) to find the likelihood profile of the  $\zeta$  parameter, and  
137 selected the value with maximum likelihood. The data were then transformed with the  
138 maximum likelihood estimate of  $\zeta$ , substituted into Eq. (1) and then standardized to zero  
139 mean and unit variance. The estimate of  $\zeta$  and summary statistics for the data after

140 transformation, and standardization, are presented in Table 2.

141 An isotropic empirical variogram of the transformed and standardized data was  
142 then computed using the method of moments estimator due to Matheron (1962) as im-  
143 plemented in the FVARIOGRAM directive in GenStat (Payne, 2010). An authorized model  
144 was then fitted to the estimated variogram by weighted least squares (Cressie, 1985) using  
145 the MVARIOGRAM procedure in GenStat (Harding et al., 2010). Alternative models were  
146 considered and the stable or powered exponential model was selected on the basis of the  
147 Akaike information criterion (McBratney and Webster, 1986). This variogram model takes  
148 the form

$$\gamma(r) = c_0 + c_1 (1 - \exp(-\{r/a\}^\kappa)), \quad (2)$$

149 where  $c_0$  and  $c_1$  are, respectively, the variances of the nugget and spatially correlated  
150 components of the variable,  $r$  is lag distance,  $a$  is a distance parameter and  $\kappa$  is a shape  
151 parameter where  $0 < \kappa \leq 2$ . The estimates of these parameters are presented in Table 2,  
152 and the estimates of the variogram of the TG variable, and the fitted model are shown in  
153 Figure 3.

### 154 2.3.2. Stochastic Geometric model

155 Estimates of the isotropic variogram of the raw data on  $EC_a$  were obtained using the  
156 method of moments estimator due to Matheron (1962) as described for the transformed  
157 data in section 2.3.1. (these are the solid symbols in Figure 6). The identification and  
158 fitting of an appropriate stochastic geometric model for the soil variable will allow us to  
159 plot a continuous variogram function for these estimates.

160 When a TG model is fitted it is assumed that, after any transformation, the data  
161  $\mathbf{y} = \{y(\mathbf{x}_1), y(\mathbf{x}_2), \dots, y(\mathbf{x}_n)\}$  from the  $n$  locations  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  can be regarded as a  
162 realization of an  $n$ -variate Gaussian random variable,  $\mathbf{Y}$ . Under this assumption the  
163 variogram of  $\mathbf{Y}$  entirely summarizes the information that the data contains about its  
164 spatial variability, and the task of estimating model parameters, under the assumption of  
165 a stationary mean, reduces to the task of estimating variogram parameters. This is not the  
166 case with models for random variables, such as the PCLT models, which have non-zero odd  
167 moments of order three or larger, and therefore are not Gaussian. The fitting of a PCLT  
168 model cannot, therefore, simply reduce to the computation of parameters which minimize

169 the weighted sum of squared residuals between the empirical and fitted variogram.

170 In this study our approach to the selection and estimation of a PCLT model is  
171 to constrain it by soil knowledge. Soil knowledge consists of general understanding of  
172 the underlying processes that influence soil formation and so the variation of the target  
173 variable, and also of general quantitative information about the variable in the study site  
174 or a homologous site, represented by summary statistics, empirical variograms or similar  
175 information. In the following sections we go through a semi-formal process of model  
176 identification based on inferences from soil knowledge and culminating in the estimation  
177 of parameters for an appropriate model. Each subsection is headed with a question, and  
178 with the general source of soil knowledge used to address it. The individual elements  
179 of soil knowledge are then summarized in brief labelled sentences, expanded in a short  
180 paragraph. Inferences from this soil knowledge are then set out.

181 *2.3.2.1. Question: ‘What mode of soil variation?’ Soil knowledge about the underlying*  
182 *pedogenetic process.*

183 The identification of a general mode of soil variation is based on two items of soil  
184 knowledge which are listed below.

185 **SK1.** *The dominant source of soil variation at metre scales in this landscape is the*  
186 *presence of Pleistocene ice-wedge polygons.* These are described in more detail by  
187 Meerschman et al (2011). Ice-wedge polygons form in periglacial conditions on sur-  
188 faces with slopes less than a critical value. Over much of central Europe ice-wedge  
189 polygons formed in periglacial conditions during the Quaternary, they are detectable  
190 at the study site from airphotography. It has been shown (Cresto Aleina et al., 2012)  
191 that the comparable polygonal patterns in ground of contemporary tundra can be  
192 modelled as a Poisson Voronoi Tessellation (PVT), that is to say one may postulate  
193 an underlying homogeneous spatial point process of completely spatially random  
194 seed points, and any one polygon consists of all locations nearest to one associated  
195 seed point than to any of the others. See Lark (2009) for a summary of some of the  
196 properties of PVT spatial processes and Okabe et al. (2000) for a more complete  
197 account. Note, in particular, that the polygons generated by this process are not of  
198 uniform size or shape. By analogy we infer that a PVT model would be a plausible

199 descriptor of the ice-wedge polygons at the study site.

200 **SK2.** *We may expect more or less continuous variation in depth-integrated soil proper-*  
201 *ties from the centre to the edge of any polygon.* Much of the polygonal patterned  
202 ground formed in Europe and North America during the Quaternary was covered  
203 by aeolian or glacio-fluvial sand or silty deposits. These have an important role in  
204 subsequent pedogenesis (Catt, 1979; Walters, 1994) imposing local lateral trends.  
205 At the centre of a polygon there is typically a relatively thin layer of sandy or silty  
206 superficial material over the host material in which the ice wedges originally formed.  
207 After thawing, the space previously occupied by ice in the wedges that delineate  
208 the polygons was typically filled with the superficial material. Any depth-integrated  
209 soil property, such as  $EC_a$ , can therefore be expected to vary laterally (although not  
210 necessarily linearly) from the centre of the polygon to its edge if there is a texture  
211 contrast between the host material and the superficial material. There is such a con-  
212 trast at the Deinze study site where the overlying material is silty-sand Quaternary  
213 deposits, and the host material is Eocene sandy clay (Meerschman et al., 2011).

214 From these two elements of soil knowledge we may infer that the spatial variation of  
215 a depth integrated soil property such as  $EC_a$ , in these conditions, can plausibly be regarded  
216 as a Poisson Continuous Local Trend random process as defined by Lark (2012b). In the  
217 next section we consider what distance function might be proposed.

218 *2.3.2.2 Question: ‘What type of distance function is plausible?’ Soil knowledge about*  
219 *pedogenetic processes and summary statistics.*

220 **SK3.** *We may expect  $EC_a$  to decline from the polygon centre to the rim.* It is generally  
221 found that measurements of  $EC_a$  made by electromagnetic induction are positively  
222 correlated with the clay content of the soil (e.g. Kachanoski et al., 2002; Saey et  
223 al., 2009). For this reason we should expect  $EC_a$ , as a depth-integrated variable,  
224 to decline from the polygon centre, where the thickness of sandy and silty material  
225 over the heavier host material is thinner, to the edge of the polygon where the  
226 former ice wedge is filled with the lighter material. This was found to be the case  
227 by Meerschman et al. (2011).

228 **SK4.** *The data on  $EC_a$  are mildly positively skewed.* This can be seen in Table 1.



229 The simplest PCLT model, as used by Lark (2012b), has a linear distance function  
 230  $\mathcal{D}(k) \propto k$ . If the distance function has a positive slope, i.e.  $\{k' > k\} \rightarrow \{\mathcal{D}(k') > \mathcal{D}(k)\}$ ,  
 231 then it can be seen that the corresponding PCLT random function has a moderate posi-  
 232 tive skewness (about 0.65). A linear distance function with a negative slope, needed for  
 233 consistency with **SK3**, would therefore give rise to a random function with a moderately  
 234 negative skewness. This is not compatible with **SK4**.

235 Of the distance functions examined by Lark (2012b) one in which the distance  
 236 function is proportional to the reciprocal of distance is compatible with **SK3** and **SK4**.  
 237 The reciprocal of distance declines with distance (**SK3**), and the example of such a random  
 238 function given by Lark (2012b) has mild positive skewness (**SK4**). On this basis it was  
 239 decided to proceed with further analysis on the assumption that the data on  $EC_a$  could be  
 240 regarded as realizations of a PCLT process with a distance function linearly proportional  
 241 to

$$\mathcal{D}(k) = \frac{1}{k + \alpha}, \quad (3)$$

242 where  $k$  is distance to the nearest event of the underlying spatial point process, and  $\alpha$   
 243 is a parameter which must take some value  $\alpha > 0$  to ensure that the distance function  
 244 is defined for all positive  $k$ . We refer to this PCLT as the inverse-distance PCLT in the  
 245 remainder of this paper. Note that the distance function in Eq. (3) defines what we shall  
 246 call the standard PCLT variable. The random variable that models the target soil variable  
 247 is linearly proportional to the standard PCLT variable, so fitting the model entails the  
 248 estimation of parameters of the standard PCLT along with a scale parameter which is the  
 249 a priori variance of the random variable.

250 The inverse-distance function was selected because it was seen to be a simple func-  
 251 tion, at least potentially compatible with available soil knowledge. In due course its  
 252 parameters are estimated and this gives some further indication of its plausibility, and in  
 253 section 2.3.3 we evaluate statistics to compare its plausibility with the TG model.

254 We call the standard inverse-distance PCLT random function  $Z_{id}$ . We shall model  
 255 the  $EC_a$  data as a realization of a random function  $Y$  where

$$Y = \beta Z = \beta(Z_n + Z_{id}), \quad (4)$$

256 where  $\beta$  is a constant of proportionality and  $Z_n$  is an independently and identically dis-

257 tributed Gaussian nugget component of mean zero. This nugget component is included  
 258 in the random model for the target variable to account for any variation spatially corre-  
 259 lated at scales finer than the sampling interval. This is common practice in geostatistical  
 260 modelling with standard covariance models such as the spherical, exponential or Matérn.

261 We now obtain the cumulative distribution and density functions of  $Z_{\text{id}}$ . We first  
 262 define the inverse of the distance function in Eq.(3),  $\dot{\mathcal{D}}(z_{\text{id}})$ , such that

$$\left\{ z_{\text{id}} = \mathcal{D}(k) = \frac{1}{k + \alpha} \right\} \Leftrightarrow \left\{ \dot{\mathcal{D}}(z_{\text{id}}) = k \right\}.$$

263 Then

$$\dot{\mathcal{D}}(z_{\text{id}}) = \frac{1}{z_{\text{id}}} - \alpha. \quad (5)$$

264 Since  $\mathcal{D}(k)$  is monotonic and decreasing with increasing  $k$  for admissible (non-  
 265 negative) values of  $k$ , the marginal cumulative distribution function of  $Z_{\text{id}}$ ,  $F_{\text{id}}(z)$  can  
 266 be written as

$$F_{\text{id}}(z_{\text{id}}) = 1 - F_k(\dot{\mathcal{D}}(z)), \quad (6)$$

267 where  $F_k(k)$  is the marginal cumulative distribution function of  $k$ . In Eq. (14) of Lark  
 268 (2012b) it is shown that, for a Poisson point process in 2-D with intensity  $\lambda$ ,

$$F(k) = 1 - \exp \{ -\lambda\pi k^2 \}, \quad (7)$$

269 and so

$$F_{\text{id}}(z_{\text{id}}) = \exp \left\{ -\lambda\pi \left( \frac{1}{z_{\text{id}}} - \alpha \right)^2 \right\}, \quad (8)$$

270 which is defined for  $0 \leq z_{\text{id}} \leq 1/\alpha$ , which shows that the random function  $Z_{\text{id}}$  has an  
 271 upper and a lower bound.

272 By differentiation of  $F_{\text{id}}(z_{\text{id}})$  with respect to  $z_{\text{id}}$  we can obtain a probability density  
 273 function (PDF):

$$\begin{aligned} f_{\text{id}}(z_{\text{id}}) &= \frac{2\lambda\pi \left( \frac{1}{z_{\text{id}}} - \alpha \right)}{z_{\text{id}}^2} \exp \left\{ -\lambda\pi \left( \frac{1}{z_{\text{id}}} - \alpha \right)^2 \right\}, \quad 0 < z_{\text{id}} \leq \frac{1}{\alpha} \\ &= 0, \quad \text{otherwise.} \end{aligned} \quad (9)$$

274 A soil variable modelled as an inverse-distance PCLT random function is assumed to  
 275 have a spatially correlated component that is linearly proportional to  $z_{id}$  for some values  
 276 of the parameters  $\alpha$  and  $\lambda$ . As noted above, the soil variable is assumed to be a realization  
 277 of a random function  $Z$  that includes an independent Gaussian nugget component of mean  
 278 zero. If the PDF of the nugget component is denoted by  $f_n(z_n)$ , then the PDF of  $Z$ ,  $f(Z)$ ,  
 279 can be obtained by the convolution operation

$$f(z) = \int_{-\infty}^{\infty} f_n(x) f_{id}(z-x) dx, \quad (10)$$

280 since  $Z_{id}$  and  $Z_n$  are independent random variables (Dudewicz and Mishra, 1988).

281 The next question that we consider is a plausible range of values for the  $\alpha$  parameter.

282 *2.3.2.3 Question: ‘What is a plausible range of values for,  $\lambda$ , the intensity of the process*  
 283 *and for the parameter  $\alpha$  of the distance function?’ Soil knowledge from field observations*  
 284 *and an estimate of the proportion of variation of  $EC_a$  that is attributable to the nugget*  
 285 *component*

286 **SK5.** *Meerschman et al. (2011) report a detailed excavation of a polygonal cell with*  
 287 *diameter about 6 m, which they regard as typical from airphoto evidence. If all cells*  
 288 *have a diameter of  $d$  m then the average intensity of an underlying spatial point*  
 289 *process is the reciprocal of the cell area which may be approximated (treating the*  
 290 *cells as circular) by  $4/\pi d^2$ . On the basis of the observation of Meerschman et al.*  
 291 *(2011) it was decided to consider a range of possible values of  $\lambda$  for the spatial point*  
 292 *process in the interval  $[0.02m^{-2}, 0.08m^{-2}]$  which corresponds to a range of polygon*  
 293 *diameters from 4 to 8 m (i.e. 2 m either side of the value proposed as representative).*

294 **SK6.** *The nugget variance of the (untransformed)  $EC_a$  data is about 10% of the correlated*  
 295 *variance. This information is used to calculate moments for the variable  $Z$ , given*  
 296 *values of  $\alpha$  and  $\lambda$ , by evaluation of the PDF in Eq.(10). It should be noted that in*  
 297 *the final model the nugget variance is estimated separately, and is not constrained*  
 298 *by this assumption. To obtain this proportion we fitted a powered exponential*  
 299 *model, Eq.(2), to the empirical variogram of the  $EC_a$  data(not shown here) using the*  
 300 *MVARIOGRAM procedure in GenStat (Harding et al, 2010).*

301 The mean and variance of an inverse-distance PCLT random function,  $Z_{id}$ , for some

302 values of the parameters  $\alpha$  and  $\lambda$  was obtained from the PDF in Eq.(9), the QDAG algorithm  
 303 in the IMSL library (Visual Numerics, 2006) was used for numerical integration. It was  
 304 then possible to compute the variance of an independent Gaussian nugget component,  $Z_n$ ,  
 305 such that the variances of  $Z_{id}$  and  $Z_n$  were in the same ratio as **SK6** suggests pertains for  
 306 the  $EC_a$  data. The coefficient of skewness for the sum of these two random variables could  
 307 then be calculated from moments obtained by numerical integration of the convolution of  
 308 the distributions of  $Z_{id}$  and  $Z_n$ , as described in Eq.(10).

309 Figure 4 is a plot of values of the skewness coefficient of variable  $Z_{id}$  for values of  
 310 the parameters  $\alpha$  and  $\lambda$ , the range for  $\lambda$  obtained from **SK3**. Note that over much of the  
 311 range of values of  $\lambda$  it is  $\alpha$  that has the strongest effect on the skewness. The two contours  
 312 drawn on the Figure bound a region within which the skewness is in the interval  $[0.25, 0.5]$ .  
 313 We regard this as mild positive skewness, compatible with **SK4**, and so we assume that  
 314 jointly plausible values of  $\alpha$  and  $\lambda$  lie within these limits. The Figure shows, for example,  
 315 that values of  $\alpha$  less than 2 m seem unlikely to be compatible with **SK4** since coefficients  
 316 of skewness for such variables are larger than 0.5. Similarly, if  $\lambda = 0.05$  then a plausible  
 317 range of values of  $\alpha$  indicated by the Figure is 2.5–3.8 m.

#### 318 *2.3.2.4 Model fitting given the soil knowledge*

319 Estimates of the isotropic variogram of the raw data on  $EC_a$  were obtained using the  
 320 method of moments estimator due to Matheron (1962) as implemented in the FVARIOGRAM  
 321 directive in GenStat (Payne, 2010). An inverse-distance PCLT model was then fitted to  
 322 the estimates by weighted least squares, but subject to the condition that  $\alpha$  and  $\lambda$  fall  
 323 jointly within the range defined by the two contours shown in Figure 4. The variogram  
 324 for the standard PCLT variable  $Z_{id}$  variable depends only on the parameters  $\alpha$  and  $\lambda$ .  
 325 In order to fit the PCLT model to the empirical variogram of the soil process it is also  
 326 necessary to estimate the proportionality constant  $\beta$  which scales the standard PCLT  
 327 variable to the variable assumed to be realized in the soil data, as shown in Eq (4). This  
 328 is done indirectly here by direct estimation of the *a priori* (sill) variance of the correlated  
 329 component of the variogram of  $Y$  (defined in Eq (4))

$$c_1 = \beta^2 \text{var}[Z_{id}]$$

330 along with a nugget component

$$c_0 = \beta^2 \text{var}[Z_n]$$

331 where  $\text{var}[Z]$  denotes the *a priori* variance of random variable  $Z$ . The fitted variogram  
 332 for the target random variable,  $Y$ , was specified by:

$$\gamma(r) = c_0 + c_1 g_{\text{id}}(r|\alpha, \lambda), \quad (11)$$

333 where  $g_{\text{id}}(r|\alpha, \lambda)$  is the variogram of the PCLT process with parameters  $\alpha$  and  $\beta$  and the  
 334 *a priori* variance scaled to 1.0 thus:

$$g_{\text{id}|\alpha, \lambda}(r) = 1 - \frac{C_{\text{id}}(r|\alpha, \lambda)}{C_{\text{id}}(0|\alpha, \lambda)}, \quad (12)$$

335 where  $C_{\text{id}}(r|\alpha, \lambda)$  is the covariance function for lag  $r$  for the standard inverse-distance  
 336 PCLT process with parameters  $\alpha$  and  $\lambda$ . The covariance function for a variable in 2-D is  
 337 given by

$$C_{\text{id}}(r|\alpha, \lambda) = \int_{\mathbb{R}^2} \{S(k, k_r) + F(k) + F(k_r) - F(k)F(k_r) - 1\} \left\{ -\frac{1}{(k + \alpha)^2} \right\} dk \left\{ -\frac{1}{(k_r + \alpha)^2} \right\} dk_r, \quad (13)$$

338 where  $S(k, k_r)$  is the joint survival function for the underlying spatial point process, as  
 339 defined by Lark (2012b). This equation is obtained directly from Eq.(20) of Lark (2012b)  
 340 and the reader is referred to that paper for details.

341 The inverse distance model was fitted as follows.

342 i). The value of the parameter  $\alpha$  was set at a fixed value, in turn  $\alpha = 2.0$  m, 2.25 m, 2.50  
 343 m . . .

344 ii). The parameter  $\lambda$  was then set at values over some range  $[\lambda_{\alpha, \min}, \lambda_{\alpha, \max}]$  where  
 345  $0.02 \leq \lambda_{\alpha, \min} < \lambda_{\alpha, \max} \leq 0.08$  such that for specified  $\alpha$  and any  $\lambda \in [\lambda_{\alpha, \min}, \lambda_{\alpha, \max}]$   
 346 the expected value of the skewness coefficient, as read off Figure 4, was within the  
 347 interval  $[0.25, 0.5]$ .

348 iii). For the set values of  $\alpha$  and  $\lambda$  values of  $c_0$  and  $c_1$  were found so that the weighted  
 349 sum of squared deviations of the variogram function in Eq (11) and the empirical  
 350 variogram (Cressie, 1985) were minimized. These values were found with the IMSL  
 351 optimization subroutine BCPOL (Visual Numerics, 2006).

352 iv). Repetition of step (iii) for successive values of  $\lambda \in [\lambda_{\alpha,\min}, \lambda_{\alpha,\max}]$  produced a ‘profile  
353 plot’ of the weighted sum of squares, WSS, against  $\lambda$ . Such plots were produced  
354 for successive values of  $\alpha$ , as designated in step (i). Estimates of  $\alpha$ ,  $\lambda$ ,  $c_0$  and  $c_1$   
355 were found from the profile plot for which the minimum WSS was the smallest of all  
356 observed values.

357 The resulting estimates of  $\alpha$  and  $\lambda$  were 2.5 m and  $0.07 \text{ m}^{-2}$  respectively. The estimates of  
358  $c_0$  and  $c_1$  were 0.49 and 4.03 respectively. Figure 5 shows the profile plot of the weighted  
359 sum of squares with  $\alpha = 2.5 \text{ m}$  and Figure 6 shows the empirical variogram for the un-  
360 transformed data and the fitted inverse-distance PCLT model. In Figure 7 is shown (line)  
361 the corresponding distribution function for the random function  $Z = Z_{\text{id}} + z_n$  standardized  
362 to zero mean and unit variance according to the values of the mean and standard deviation  
363 obtained from the PDF in Eq (10). Also plotted on Figure 7 are points from the empirical  
364 CDF of the standardized  $\text{EC}_a$  data. The theoretical and empirical distribution functions  
365 are in reasonable agreement, although the median of the former seems to be rather smaller  
366 than the latter.

### 367 2.3.3. Comparing the TG and PCLT models

368 It is well known that Gaussian (and trans-Gaussian) models of spatial variation, in which  
369 all information on variability is expressed by two-point statistics such as the covariance  
370 function, are not able to reproduce all important features of natural spatial fields, which  
371 must be represented by higher-order moments (e.g. Guardiano and Srivastava, 1993). This  
372 has been the motivation for the development of multiple point statistics. In this section we  
373 investigate whether the PCLT model allows better characterization of the spatial structure  
374 of the  $\text{EC}_a$  data than does the TG model.

375 One feature of the Gaussian and trans-Gaussian random variables that often limits  
376 their applicability is the fact that large values of the variable tend to be spatially isolated  
377 from other large values, the same holds for small values (e.g. Gómez-Hernández and Wen,  
378 1997; Strebelle, 2002). In this case study we may consider locations with small values of  
379  $\text{EC}_a$ . These locations are likely to be dominated by lighter sandy and silty Quaternary  
380 material, rather than the heavier-textured Eocene host material, and so will have larger  
381 porosity and hydraulic conductivity, than sites where the  $\text{EC}_a$  is larger. If the TG model

382 does not adequately represent the connectivity of such areas then any modelling based  
 383 on TG simulation will fail to represent processes where this lateral connectivity matters.  
 384 This could include processes such as lateral movement of a pollutant plume in saturated  
 385 conditions, the response of the water table to drainage schemes or the lateral spread of  
 386 root pathogens. Figure 8 shows sets of realization of each of the fitted PCLT and TG  
 387 models for  $EC_a$ . The inverse-distance PCLT realizations were generated directly following  
 388 the procedure used by Lark (2012b). The TG realizations were obtained by Sequential  
 389 Gaussian Simulation using the SGSIM subroutine from the GSLIB library (Deutsch and  
 390 Journel, 1997) modified to use the powered exponential variogram function. On visual  
 391 inspection it can be seen that, while some large patches with smaller  $EC_a$  values are seen  
 392 in the TG realization, there are fewer isolated small patches with small  $EC_a$  values in  
 393 the inverse-distance PCLT realization, which has large and connected regions with small  
 394 conductivity around the boundaries of the Voronoi cells of the underlying point process.  
 395 However, this visual inspection is of limited usefulness and a more objective measure is  
 396 needed.

397 To this end we consider a simple test statistic, which can be readily evaluated on the  
 398  $EC_a$  data which are more or less regularly sampled but which do not constitute a compre-  
 399 hensively observed ‘image’. We define the statistic  $P(\tau, \Delta)$  as the expected proportion of  
 400 observations within a square window of width  $\Delta$ , centred at a randomly selected location  
 401  $\mathbf{x}$  which are  $\leq \tau$ , conditional on the value at  $\mathbf{x}$  being  $\leq \tau$ . We may expect these values  
 402 to be smaller for a TG random function than for a function which better-represents the  
 403 spatial structure of a variable in which small values tend to be spatially connected.

404 We estimated  $P(\tau, \Delta)$  for the TG and PCLT random functions fitted to the  $EC_a$   
 405 data by simulation. These are denoted by  $P_{TG}(\tau, \Delta)$  and  $P_{PCLT}(\tau, \Delta)$  respectively. We  
 406 considered windows of width 2 m or larger (because approximately 40  $EC_a$  observations  
 407 occur within a 2-m window). Each simulation program generated a single independent  
 408 realization of the random function at 25 equally-spaced locations in a window of width  $\Delta$   
 409 one of which was at the centre of the window. If the simulated value at the centre was  $\leq \tau$ ,  
 410 the conditioning criterion, then the realization was retained and  $P(\tau, \Delta)$  was estimated as  
 411 the proportion of the observations in the window for which  $\leq \tau$ . This was repeated until  
 412 10 000 independent realizations which met the criterion that the central value was  $\leq \tau$  had

413 been obtained. The PCLT realizations were generated using the procedure described by  
 414 Lark (2012b). The TG realizations were obtained by LU decomposition (Goovaerts, 1997).  
 415 The mean value of  $P_{TG}(\tau, \Delta)$  and the standard deviation of the 10 000 independent values,  
 416 were computed for different values of  $\Delta$  and for  $\tau$  set to the median, first quartile and first  
 417 octile of the  $EC_a$  data. This was also done for  $P_{PCLT}(\tau, \Delta)$ . The difference between the  
 418 mean values of  $P_{PCLT}(\tau, \Delta)$  and  $P_{TG}(\tau, \Delta)$  for these different thresholds and for windows  
 419 of different size, are plotted in Figure 9.

420 Figure 9 shows three things. First, the mean value of  $P_{PCLT}(\tau, \Delta)$  is larger than that  
 421 of  $P_{TG}(\tau, \Delta)$  for given  $\tau$  and  $\delta$ . That is to say, given that a value falls below a threshold,  
 422 there is a larger proportion of neighbouring values which do so for the PCLT process  
 423 than for the TG process. Second, the effect depends on the threshold, and increases as  
 424 the threshold becomes more extreme relative to the overall distribution. Third, the effect  
 425 depends on the window size. It is small for a large window, but it is also notable that the  
 426 difference is larger for the window width 4 m than the window width 2 m. This reflects  
 427 the spatial scale of the random function.

428 The  $P(\tau, \Delta)$  statistic was then estimated from the  $EC_a$  data for the same three  
 429 threshold values used in the simulations, and for  $\Delta = 4m$  given that this window showed  
 430 the largest differences between the two processes in the simulation. An independent ran-  
 431 dom subsample of 250 observations for which  $EC_a \leq \tau$  was obtained, the proportion of  $EC_a$   
 432 observations within a square window, width  $\Delta$  about each of these observations was com-  
 433 puted. The results are shown in Figure 10. The mean value of  $P_{TG}(\tau, \Delta)$  and  $P_{PCLT}(\tau, \Delta)$   
 434 from the simulations are plotted, and for each of these the 95% confidence interval for the  
 435 mean of a sample of 250 independent observations is also shown, based on the variances  
 436 of the values obtained by simulation. The estimates from the  $EC_a$  data are also plotted.  
 437 Note that for all three thresholds the values of  $P(\tau, \Delta)$  for the data are larger than the  
 438 upper limit of the confidence interval for the TG process. For  $\tau$  equal to the median and  
 439 the first quartile the values from the data are within the confidence interval for the PCLT  
 440 process, for the first octile the estimate is slightly smaller than the confidence interval for  
 441 the PCLT process, but closer to the expected value for the PCLT process than it is for  
 442 the TG process.



### 443 3. Discussion

444 The overall objective of this study was to identify a stochastic model for a soil  
445 property that varies according to some mode, and to base this identification as far as  
446 possible on knowledge of the underlying soil process and, at most, some simple descriptive  
447 statistics of the variable such as the empirical variogram and summary statistics. This  
448 was achieved in this study by employing general soil knowledge in a structured way. This  
449 is proposed as a framework for similar studies on soil variation in contrasting modes.

450 The PCLT model used here is a stochastic model of soil variability selected because  
451 it represents a particular model of soil variation. This places it in between the most  
452 common approach to stochastic modelling, where a Gaussian or TG model is selected for  
453 convenience, and approaches based on direct specification of the form of the covariance  
454 function from a mechanistic model of the process. The latter has been achieved only  
455 for a limited set of processes over a limited range of spatial scales, e.g. Whittle (1954,  
456 1962), Kolvos et al. (2004). Essentially the PCLT model is selected because it is in  
457 some sense an analogue of the soil process of interest. A similar approach has been used  
458 elsewhere. For example, Smith et al. (200) selected a ‘blur’ process based on convolution  
459 to model the space-time covariance of atmospheric pollutants, the convolution process  
460 was an analogue of pollutant dispersal. Similarly, Brochu and Marcotte (1993) selected a  
461 generalized Cauchy variogram model for observations on hydraulic head on the grounds  
462 that this process had physical analogies with a gravimetric field, which is mechanistically  
463 linked to the Cauchy model.

464 The use of stochastic geometric analogues of soil processes to generate stochastic  
465 models is attractive. It remains to be seen how wide a range of soil processes can be  
466 represented that way, and it is accepted that lateral textural variations in patterned ground  
467 are at once likely to be represented by simple geometric models and rather unrepresentative  
468 of soil variation in most conditions. None the less, the approach to the identification of  
469 models based on finding operators that are analogues for processes in the soil is likely to be  
470 more successful than the search for stochastic models based on strictly mechanistic models.  
471 It must also be noted that the stochastic geometric approach naturally reproduces non-  
472 Gaussian variation which must be characterized by moments of order higher than two,  
473 whereas the mechanistic approaches to spatial modelling are often explicitly based on

474 two-point statistics, the covariance function (e.g. Whittle, 1954; 1962).

475       The particular advantage of the stochastic geometric approach in this case study  
476 is how the inverse-distance PCLT model was better than the TG model in terms of the  
477 test statistic on the connectivity of values with small  $EC_a$ . If one wanted to generate  
478 conditional simulations of the soil in this environment as a basis for computing, for ex-  
479 ample, distributions of upscaled processes such as pollutant transport across a block of  
480 land, then the inverse-distance PCLT model would produce superior representations of the  
481 connectivity of material with large conductivities, and so of preferential flow pathways.

482       There is considerable scope for further development of this approach. Other dis-  
483 tance functions could be considered for this variable, and for others. In this study we  
484 looked for the simplest distance function that seemed to be compatible with soil knowl-  
485 edge, and there may be scope further to refine a framework for selecting a function. More  
486 specific soil knowledge could be used. For example, in the case study considered here,  
487 one could generate a simple conceptual 3D model of a polygon, with material with dif-  
488 ferent dielectrical properties, and compute the expected trend function from models of  
489 the EM properties of the soil. While the objective of this particular study was to restrict  
490 the use of direct observations on the target variable to simple descriptive statistics, one  
491 might also conduct specific surveys at fine scale on transects across polygons in order to  
492 identify plausible distance functions for further studies. It should also be noted that the  
493 homogeneous Poisson process, while a default spatial model, is not the only one available  
494 and might not be generally appropriate. While it was selected in this case on the basis  
495 of recent work on patterned ground (Cresto Aleina et al., 2012), it is likely that, at the  
496 limit, a more overdispersed spatial process would be more appropriate for this problem,  
497 with fewer close-spaced points than in the homogeneous Poisson case.

498       The model-fitting framework in this study made combined use of point estimates of  
499 the variogram, and a weighted least squares criterion for parameter estimation, subject  
500 to constraints identified from soil knowledge which imposed constraints on the modelled  
501 parameters based, in this case, on the coefficient of skewness. This remains a somewhat  
502 arbitrary procedure for parameter estimation. Ideally a likelihood-based estimator should  
503 be derived. This is unlikely to be straightforward, not least because the joint distribution  
504 function of any PCLT process is complex and requires geometrical functions for which

505 analytical expressions are not known. In other settings, when the likelihood function is  
506 expensive to evaluate, parameters may be estimated by an extension of the method of  
507 moments to include higher order moments than the usual first and second. An example of  
508 this is given by Iskander and Zoubir (1999), and it is suggested that a method of higher-  
509 order moments is most likely to be a tractable solution to fitting stochastic geometric  
510 models.

511       There is scope for further work on the comparison of realizations of the PCLT and  
512 TG processes with respect to multiple point statistics and for weighing the evidence that  
513 one model rather than the other best represents particular data. We used a relatively  
514 simple statistic in this paper, given that our data are not-quite regularly sampled and  
515 so do not constitute an image. However, it would be interesting to see how statistics  
516 developed for images (e.g. De Iaco and Maggio, 2011) might be used to evaluate alternative  
517 stochastic models. That said, the statistic which we used in this paper was not a general  
518 measure of spatial structure but rather was focussed on a particular problem of direct  
519 interest (i.e. the connectedness of areas likely to have larger hydraulic conductivities).  
520 This is arguably more relevant than a generalized measure. It would be interesting to  
521 develop methods to quantify the spatial structure of random fields as this affects particular  
522 processes. For example, one might compare the outcomes of a process model for the  
523 dispersal of contaminant plumes when it is run with input data on conductivity or similar  
524 model parameters which are realizations of contrasting random processes.

525       Any PCLT model could be used in conventional spatial prediction by kriging since  
526 the variogram or, equivalently, the covariance function can be specified. However, since  
527 the PCLT covariance function is not available in closed form, it would generally be more  
528 efficient to use a standard variogram function for kriging; and since kriging uses only the  
529 two-point statistics of a variable there is unlikely to be any benefit in using the PCLT model  
530 rather than a standard spatial model for this purpose. The value of the PCLT model is not  
531 to provide an alternative form of the covariance function, but rather for spatial prediction  
532 of non-Gaussian variables whose multivariate distribution is not entirely characterized by  
533 the covariance function. Spatial prediction in such cases may be done by codes  
534 such as SNESIM (Strebelle, 2002) or the direct sampling (DS) algorithm of Mariethoz et  
535 al. (2010) which allow one to obtain conditional distributions at unsampled sites from

536 multiple realizations of a non-Gaussian process. These procedures require training images  
537 of the variables of interest, and the availability of sufficient training images of adequate  
538 quality is a potential limitation on the use of multiple point geostatistical methods in  
539 soil science. For this reason Pyrcz et al. (2008) developed a library of training images  
540 for a particular geological setting (fluvial and deepwater reservoirs) by a combination of  
541 stochastic and object-based simulation methods. If an appropriate PCLT process could be  
542 identified for a particular soil variable, then it might be used similarly to generate training  
543 images, either for a library or as required for a multiple point conditional simulation. It is  
544 easy to generate multiple training images from a PCLT model. This would be particularly  
545 advantageous for the DS algorithm, because it has been noted (e.g. Meerschman et al.,  
546 2013) that multiple realization generated by the DS algorithm sometimes all include exact  
547 copies of significant patches of the (single) training image. This could be avoided by  
548 modifying the DS algorithm to sample multiple training images in random order, when  
549 these can readily be generated.

#### 550 **4. Conclusions**

551 We have shown how a structured use of soil knowledge allows us to fit an appropriate  
552 stochastic geometric model to data on a soil property in a particular environment. Further-  
553 more, we have shown that this model appears to capture features of the spatial variation of  
554 our target variable better than the standard Gaussian model, even after transformation of  
555 the data to marginal normality. There is more work to be done in the development of this  
556 approach, and exploring its practical implications but we believe this case study shows that  
557 there is considerable potential. In particular, realizations of PCLT processes may be bet-  
558 ter than standard TG simulations for predicting outcomes of non-linear processes such as  
559 contaminant transport, and for quantifying the uncertainty of such predictions. If PCLT  
560 models succeed in capturing the multiple point behaviour of soil variables, then PCLT  
561 simulation can be used to provide an inexhaustible supply of training images for existing  
562 multiple point prediction code. This removes one major limitation on the application of  
563 this emerging geostatistical methodology.

564 **Acknowledgements**

565 This paper is published with the permission of the Director of the British Geological  
566 Survey (Natural Environment Research Council). The second author was funded by the  
567 Fund for Scientific Research-Flanders (FWO-Vlaanderen). We are grateful to Professor  
568 Alex McBratney and to two anonymous referees for helpful comments on this paper.

**References**

- Brochu, Y., Marcotte, D. 2003. A simple approach to account for radial flow and boundary conditions when kriging hydraulic head fields for confined aquifers. *Mathematical Geology*, 35, 111–139.
- Brown, P.E., Roberts, G.O., Kåresen, K.F., Tonellato, S. 2000. Blur-generated non-separable space-time models. *Journal of the Royal Statistical Society B*, 62, 847–860.
- Catt, J.A. 1979. Soils and Quaternary Geology in Britain. *Journal of Soil Science* 30, 607–642.
- Chilès, J.P., Guillen, A., 1984. Variogrammes et krigeages pour la gravimétrie et le magnétisme. *Sciences de la Terre – Série Informatique Géologique*, 20, 455–468.
- Cockx, L., Ghysels, G., Van Meirvenne, M., Heyse, I. 2006. Prospecting frost-wedge pseudomorphs and their polygonal network using the electromagnetic induction sensor EM38DD. *Permafrost and Periglacial Processes*, 17, 163–168.
- Cressie, N. 1985. Fitting variogram models by weighted least squares. *Mathematical Geology*, 17, 563–586.
- Cresto Aleina, F., Brovkin, V., Muster, S., Boike, J., Kutzbach, L., Zuyev, S. 2012. Poisson-Voronoi diagrams and the polygonal tundra. *Geophysical Research Abstracts*, 14, EGU2012-1963-1.
- De Iaco, S., Maggio, S. 2011. Validation techniques for geological patterns simulations based on variogram and multiple point statistics. *Mathematical Geoscience* 43, 483–500.

- Deutsch, C.V., Journel, A.G. 1997. GSLIB: Geostatistical software library and user's guide. 2nd Edition. Oxford University Press, New York.
- Dudewicz, E.J., Mishra, S.N. 1988. Modern mathematical statistics. John Wiley & Sons, New York.
- Goovaerts, P. 1997. Geostatistics for natural resources evaluation. Oxford University Press, New York.
- Gómez-Hernández, J.J., Wen, X-H. 1998. To be or not to be multi-Gaussian? A reflection on stochastic hydrology. *Advances in Water Resources*, 21, 47–61.
- Guardiano, F., Srivastava, R.M., 1993. Multivariate geostatistics: beyond bivariate moments. In: Soares, A. (Ed.), *Geostatistics-Troia*, Vol. 1. Kluwer, Dordrecht, pp. 133–144.
- Harding, S.A., Murray, D.A., Webster, R. 2010. MVARIOGRAM procedure in (ed.) R.W. Payne, *GenStat Release 13 Reference Manual, Part 3 Procedure Library PL21*. VSN International, Hemel Hempstead.
- Kachanoski, R.G., Hendrickx, J.M.H., de Jong, E. 2002. Electromagnetic induction. In: (eds) J.H. Dane, G.C. Topp. *Methods of Soil Analysis, Part 1, Physical Methods*, Third Edition. Soil Science Society of America. (2002), pp.497–501.
- Iskander, R., Zoubier, A.M. 1999. Estimation of the parameters of the  $K$ -distribution using higher order and fractional moments. *IEEE Transactions on Aerospace and Electronic Systems*, 35, 1453–1456.
- Kolvos, A., Christakos, G., Hristopulos, D.T., Serre, M.L. 2004. Methods for generating non-separable spatiotemporal covariance models with potential environmental applications. *Advances in Water Resources* 27, 815–830.
- Lark, R.M. 2009. A stochastic-geometric model of soil variation. *European Journal of Soil Science*, 60, 706–719.
- Lark, R.M. 2012a. Towards soil geostatistics. *Spatial Statistics* 1, 92–99.

- Lark, R.M. 2012b. A stochastic geometric model for continuous local trends in soil variation. *Geoderma*, 189–190, 661–670.
- Mariethoz, G., Renard, P., Straubhaar, J. 2010. The Direct Sampling method to perform multiple-point geostatistical simulations. *Water Resources Research* 46, W11536.
- Matheron, G. 1962. *Traité de Géostatistique Appliqué, Tome 1. Memoires du Bureau de Recherches Géologiques et Minières*, Paris.
- McBratney, A.B., Webster, R. 1986. Choosing functions for semi-variograms of soil properties and fitting them to sampling estimates. *Journal of Soil Science* 37, 617–639.
- Meerschman, E., Van Meirvenne, M., De Smedt, T., Islam, M.M., Meeuws, F., Van De Vijver, E., Ghysels, G. 2011. Imaging a polygonal network of ice-wedge casts with an electromagnetic induction sensor. *Soil Science Society of America Journal* 75, 2095–2100.
- Meerschman, E., Pirot, G., Mariethoz, G., Straubhaar, J., Van Meirvenne, M., Renard, P. 2013. A practical guide to performing multiple-point statistical simulations with the Direct Sampling algorithm. *Computers & Geosciences*, 52, 307–324.
- Okabe, A., Boots, B., Sugihara, K., Chiu, S.K. 2000. *Spatial Tessellations: Concepts and Applications of Voronoi Diagrams*. 2nd Edition John Wiley & Sons, Chichester.
- Payne, R.W. (ed) 2010. *GenStat Release 13 Reference Manual, Part 2 Directives*. VSN International, Hemel Hempstead.
- Pyrzcz, M.J., Boisvert, J.B., Deutsch, C.V. 2008. A library of training images for fluvial and deepwater reservoirs and associated code. *Computers & Geosciences* 43, 542–560.
- R Development Core Team. 2012. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org/>.
- Saey, T., Van Meirvenne, M., Vermeersch, H., Ameloot, N., Cockx, L. 2009. A pedo-transfer function to evaluate the soil profile textural heterogeneity using proximally sensed apparent electrical conductivity. *Geoderma*, 150, 389–395.

- Strebelle, S. 2002. Conditional simulation of complex geological structures using multiple-point statistics. *Mathematical Geology*, 34, 1–21.
- Venables, W.N., Ripley, B.D. 2002. *Modern Applied Statistics with S*. Fourth Edition. Springer, New York.
- Visual Numerics, 2006. IMSL Fortran Numerical Library Version 6.0. Visual Numerics, Houston, Texas.
- Walters, J.C. 1994. Ice-wedge casts and relict polygonal patterned ground in North-East Iowa, USA. *Permafrost and Periglacial Processes* 5, 269–281.
- Webster, R. & Oliver, M.A. 2007. *Geostatistics for Environmental Scientists*. 2nd Edition. John Wiley & Sons, Chichester.
- Whittle, P. 1954. On stationary processes in the plane. *Biometrika*, 41, 434–449.
- Whittle, P. 1962. Topographic correlations, power-law covariance functions and diffusion. *Biometrika*, 49, 305–314.



**Table 1.** Summary statistics of the raw data on EC<sub>a</sub>.

Statistic	mS m <sup>-1</sup>
Average	31.37
Median	31.13
Standard deviation	2.2
Skewness	0.36
Quartile 1	29.9
Quartile 3	32.76
Octile 1	29.03
Octile 7	34.08

**Table 2.** Summary statistics of the data on  $EC_a$  after the Box-Cox transformation and for the transformed data after standardization. Variogram parameters for the standardized data are also given.

Statistic	Transformed data	Transformed and standardized data
Average	1.508	0
Median	1.507	-0.056
Standard deviation	0.01	1
Skewness	0	0
Quartile 1	1.501	-0.646
Quartile 3	1.514	0.668
Octile 1	1.497	-1.085
Octile 7	1.52	1.216
$\zeta^*$	-0.57	
Variogram parameters <sup>†</sup>		
$c_0$		0.12
$c_1$		0.84
a		1.91
$\kappa$		1.49

\* Maximum likelihood estimate of the parameter of the Box-Cox transform, see Eq.(1)

† Powered (stable) exponential model, see Eq.(2).

## Figure captions

1.  $EC_a$  data, coordinates are in metres relative to a local datum.
2. Histogram of  $EC_a$  data.
3. Empirical variogram of transformed and standardized  $EC_a$  data with a fitted model.
4. Values of the coefficient of skewness for an inverse-distance PCLT process with different values of the parameters  $\lambda$  and  $\alpha$ . The two contours bound the region where we regard the variable as mildly positively skewed.
5. Profile plot of the weighted sum of squares for the fit of the inverse-distance PCLT variogram function against  $\lambda$ , with  $\alpha$  fixed at 2.5 m.
6. Empirical variogram of the untransformed  $EC_a$  data with the fitted inverse-distance PCLT variogram.
7. Marginal distribution function of the standardized inverse-distance PCLT random function with  $\alpha=2.5$  m and  $\lambda=0.07$  m<sup>-2</sup> (line). The points are from the empirical cumulative distribution function of the standardized  $EC_a$  data.
8. Realizations of (a) the inverse-distance PCLT random function and (b) the TG random function (back transformed to original units) on a 0.25-m square grid.
9. Plot of the difference between the mean of  $P_{PCLT}(\tau, \Delta)$  and that of  $P_{TG}(\tau, \Delta)$  for different window widths ( $\Delta$ ) and with  $\tau$  set to the median, first quartile and first octile of the  $EC_a$  data. Mean for 10 000 realizations of each random function.
10.  $P(\tau, \Delta)$  with  $\Delta = 4$  m plotted against  $\tau$  set to the median, first quartile or first octile. The solid disc, ●, is the mean value from 10 000 realizations of the PCLT random function, the solid square, ■, is the mean value from 10 000 realizations of the TG random function. The horizontal bars show the 95% confidence interval for the mean of based on 250 independently and randomly selected locations that mean the conditioning criteria. The crosses, × show the mean values for 250 independently and randomly selected sites in the  $EC_a$  data set.

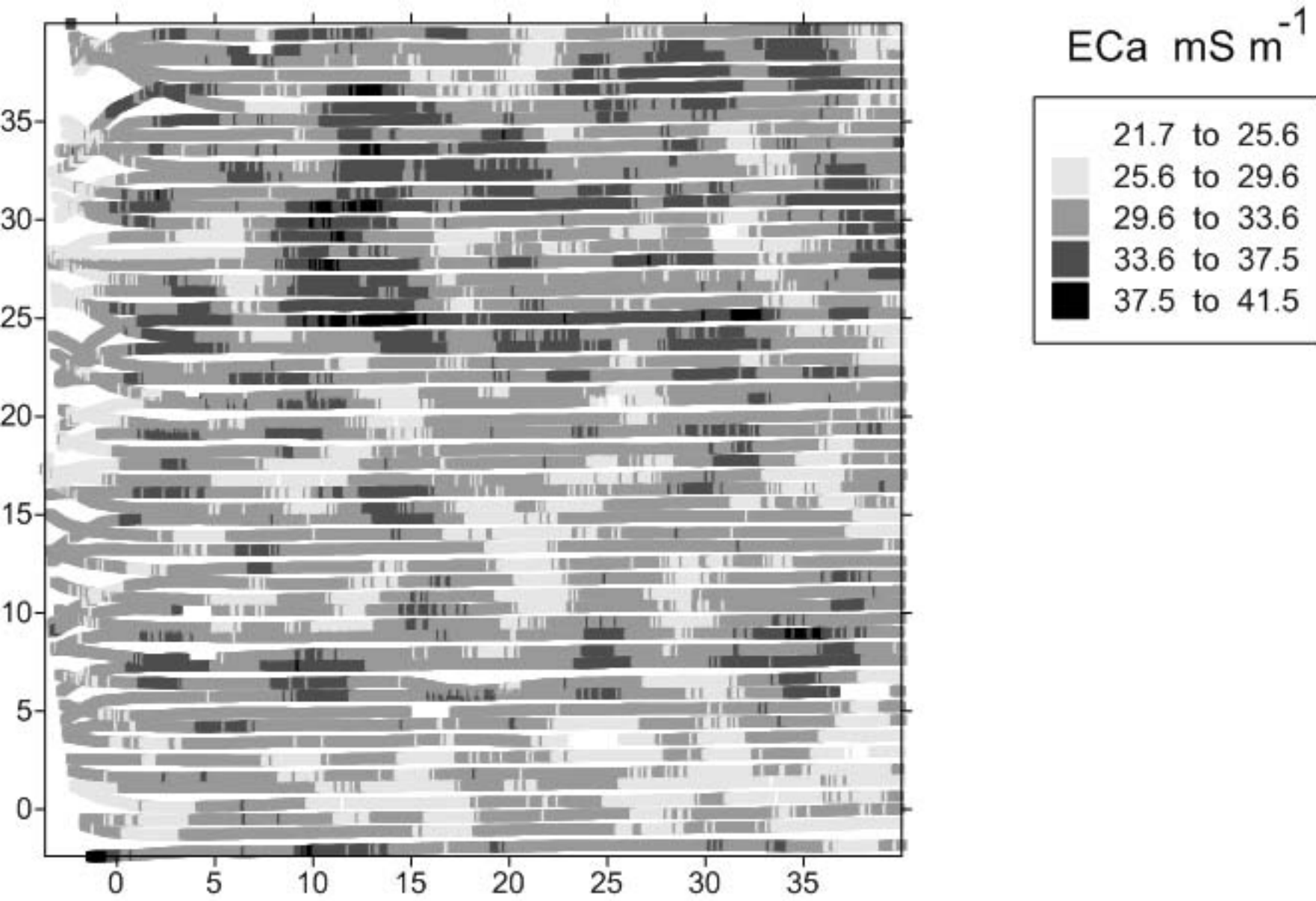


Figure 2

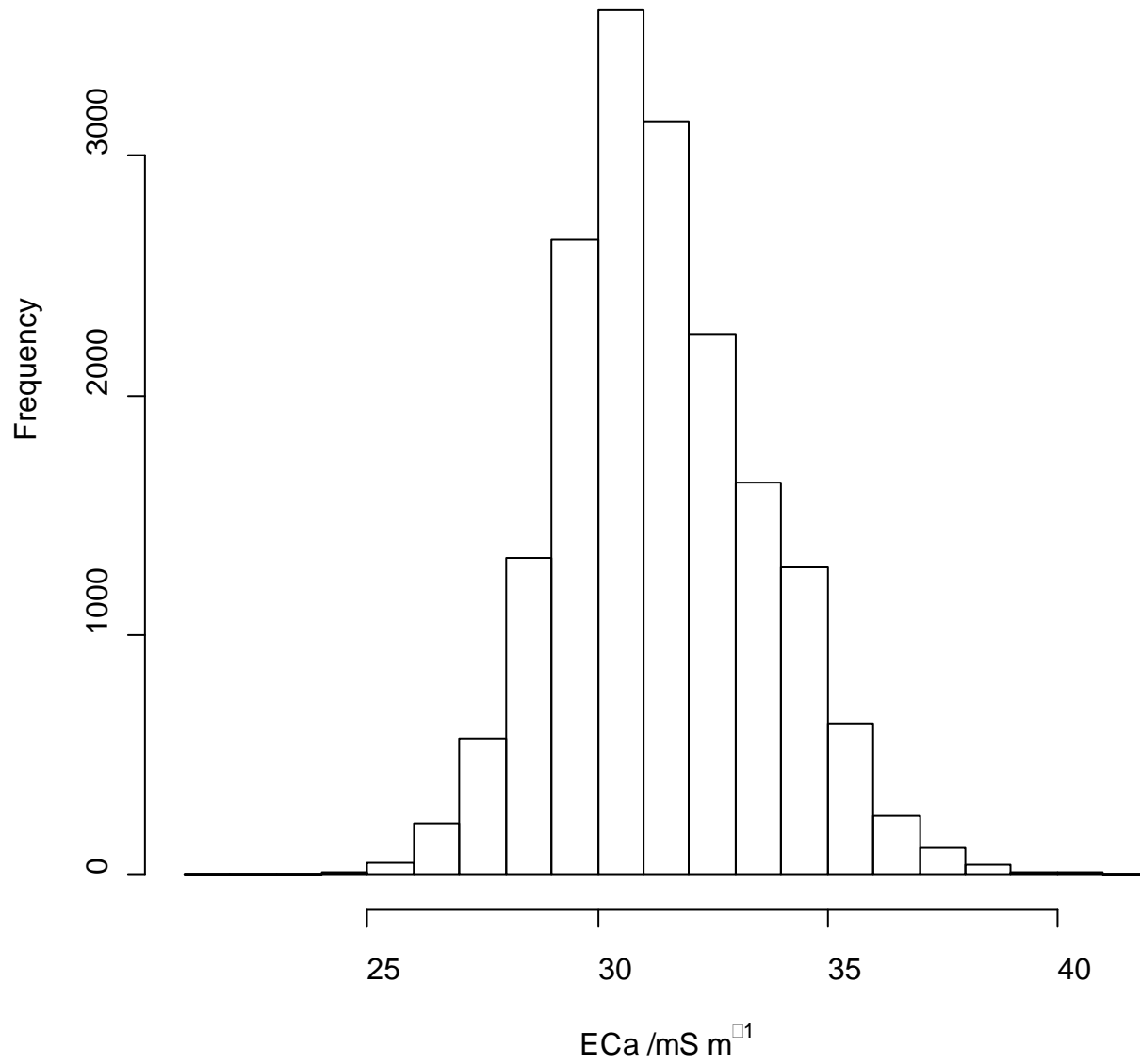
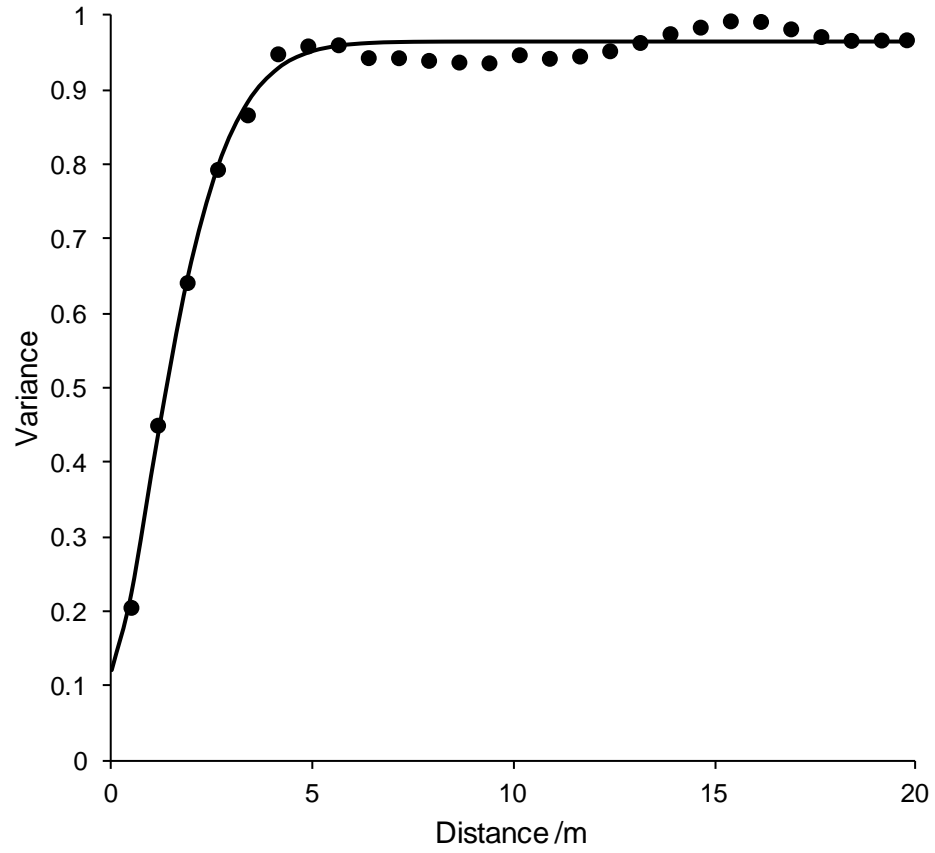


Figure 3



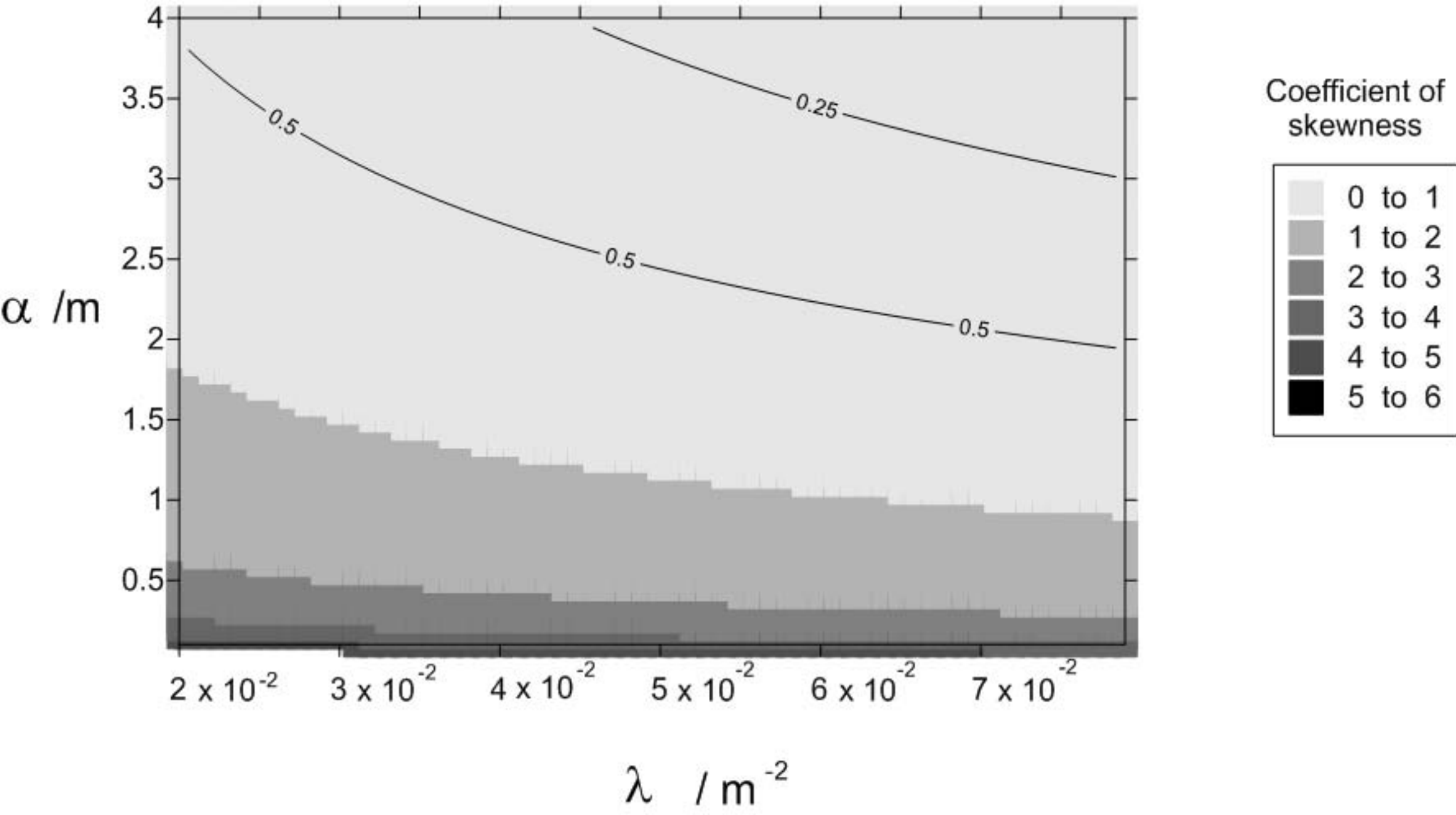


Figure 5

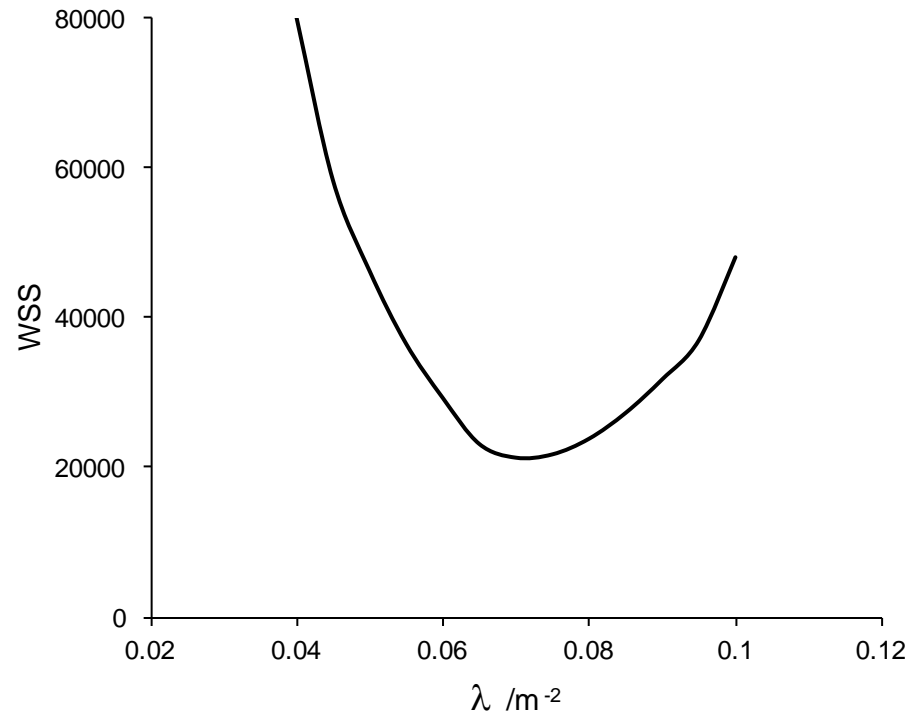




Figure 6

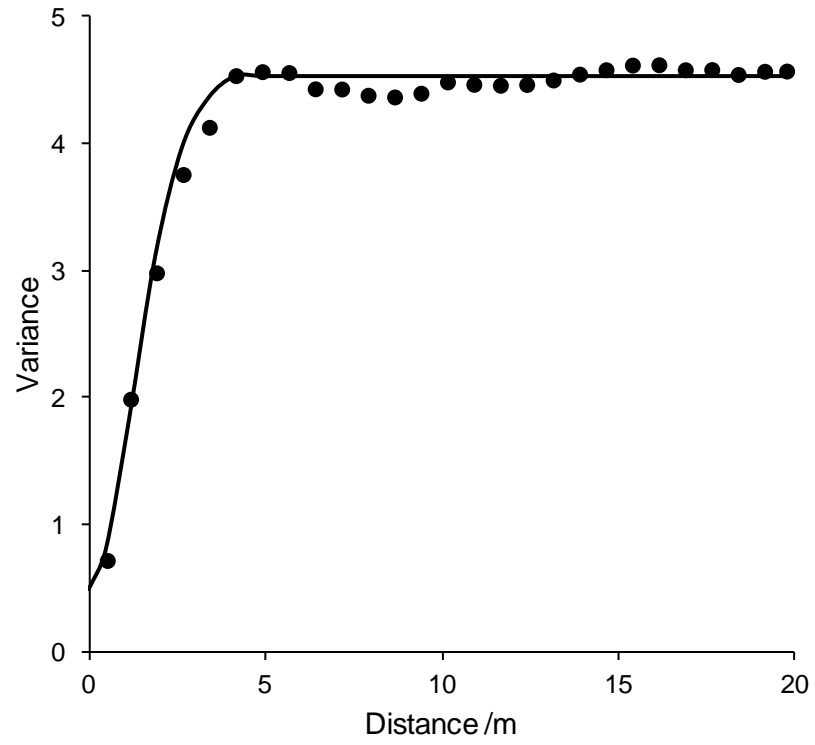
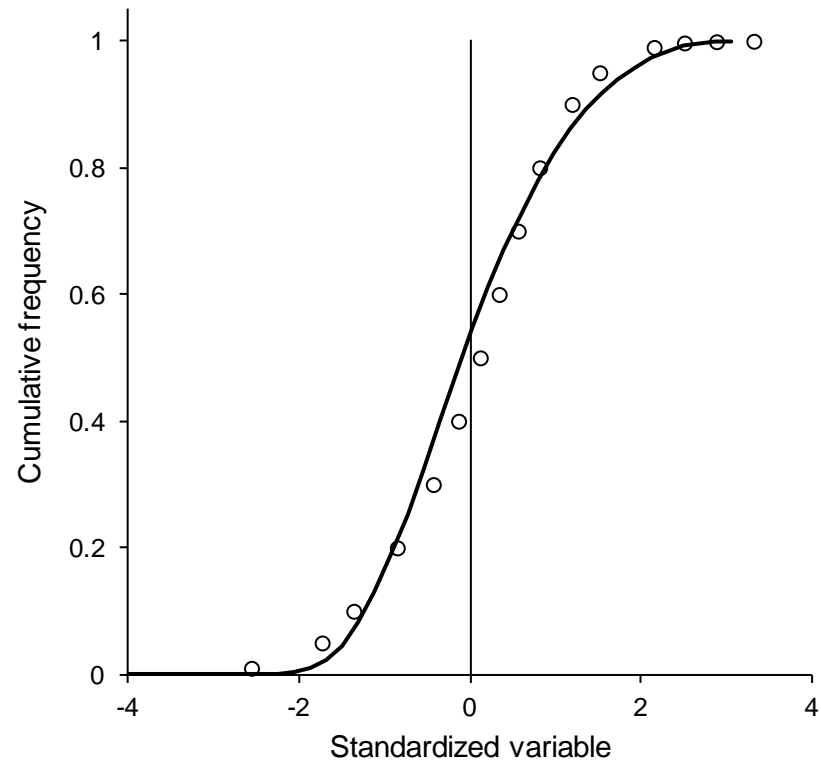
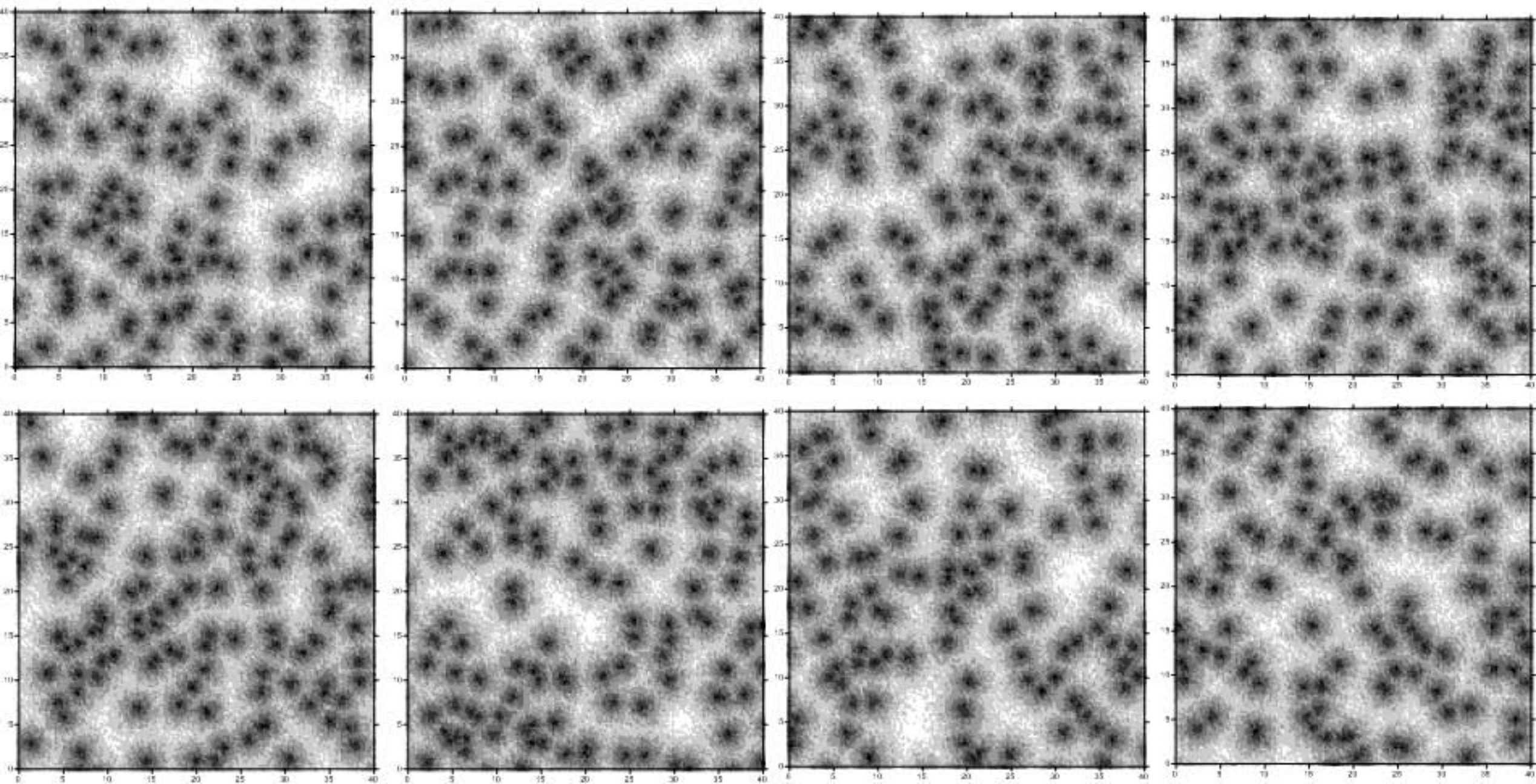
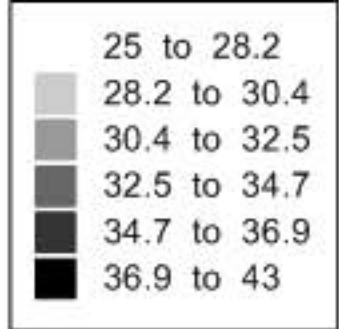


Figure 7





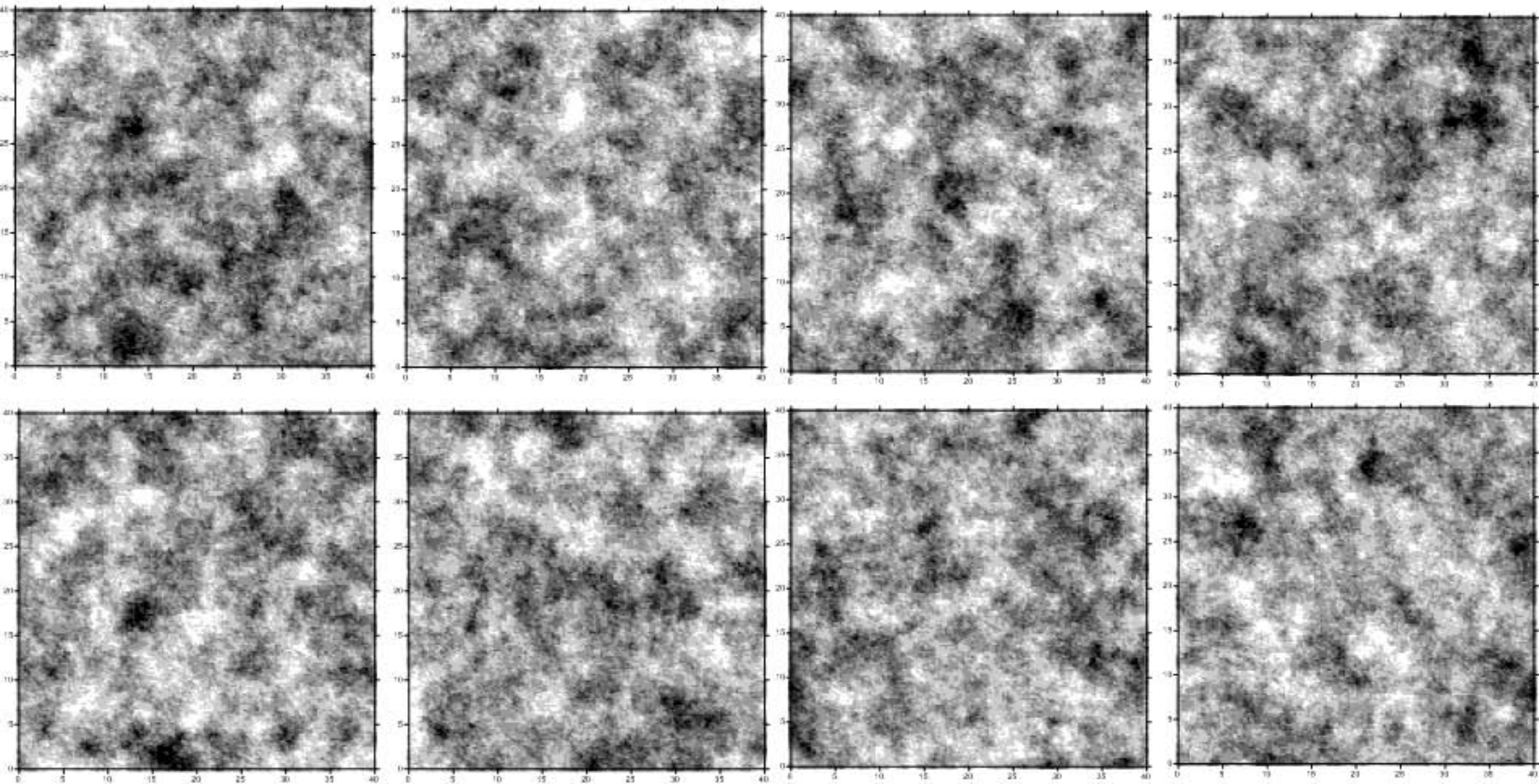
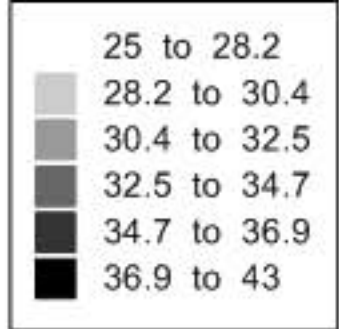


Figure 9

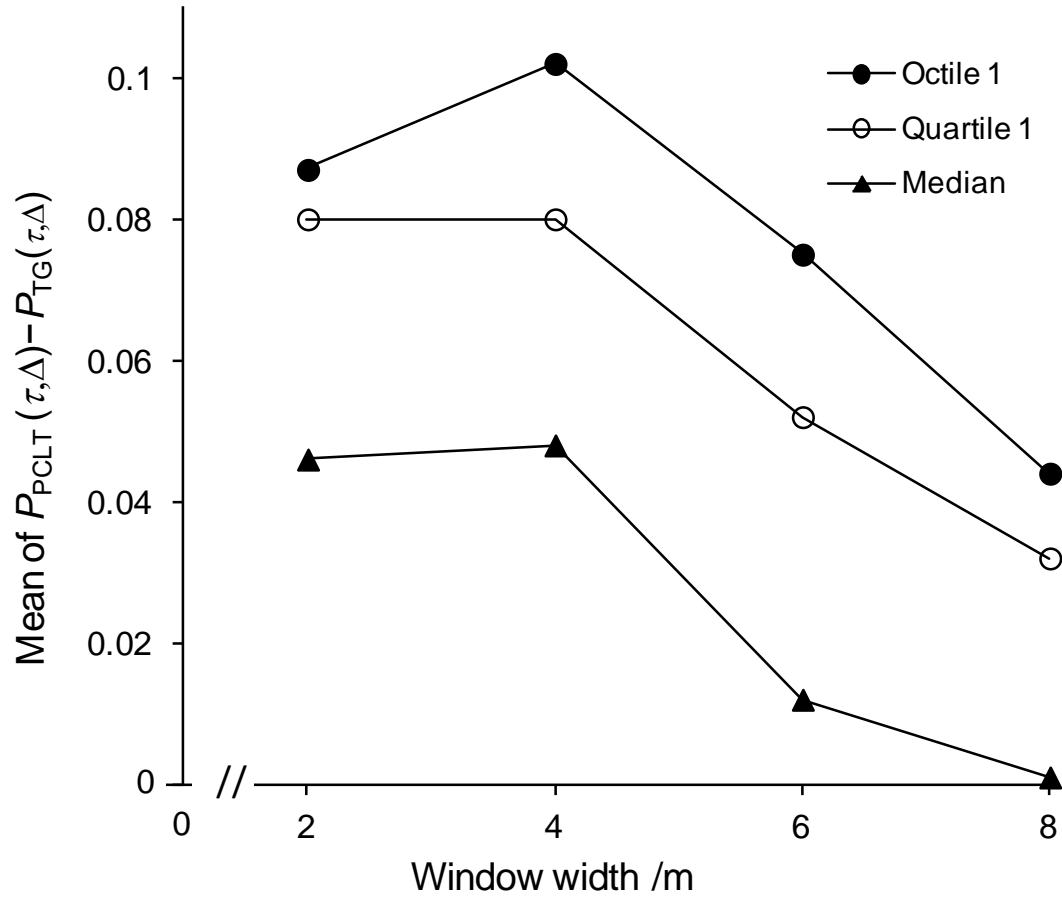


Figure 10

