

Using third-order cumulants to investigate spatial variation: a case study on the porosity of the Bunter Sandstone

R.M. Lark*

British Geological Survey, Keyworth, Nottinghamshire NG12 5GG, U.K.

1 Abstract

2 The multivariate cumulants characterize aspects of the spatial variability of a re-
3 gionalized variable. A centred multivariate Gaussian random variable, for example, has
4 zero third-order cumulants. In this paper it is shown how the third-order cumulants can
5 be used to test the plausibility of the assumption of multivariate normality for the porosity
6 of an important formation, the Bunter Sandstone in the North Sea. The results suggest
7 that the spatial variability of this variable deviates from multivariate normality, and that
8 this assumption may lead to misleading inferences about, for example, the uncertainty
9 attached to kriging predictions.

11 1. Introduction

12 Geostatistical analysis of spatially variable geological data allows us to quantify the
13 uncertainties in inferences made from partial samples by treating data as realizations of
14 a random field. In most cases the underlying model is multivariate Gaussian, and the
15 plausibility of this assumption is usually judged from the marginal distribution of obser-
16 vations (e.g. Webster and Oliver, 2007). Where necessary the data may be transformed,
17 for example to logarithms or, more generally, by the Box-Cox transformation. However, it
18 is recognized that the assumption of a Gaussian or trans-Gaussian (Gaussian after trans-
19 formation) distribution is not always safe, and, particularly, that it might not hold even
20 when it seems plausible for the marginal distribution of the data. Of particular concern
21 is the recognition that, under the multivariate Gaussian model, the first and second order

* *E-mail address:* mlarke@nerc.ac.uk (R.M. Lark).

22 moments entirely characterize the spatial distribution of a variable since all odd moments
23 larger than the first are zero and all even moments larger than the second can be written
24 in terms of it. However, it is known that the complex geometries that may be encountered
25 in geological data, the strongly-connected patterns of coarse-textured alluvium in former
26 braided streams are a *locus classicus*, might not be fully characterized by the first and
27 second moments, and more complex spatial distributions are necessary (e.g. Guardiano
28 and Srivastava, 1993).

29 It is therefore necessary to develop exploratory methods to examine the higher-
30 order behaviour of spatially variable data. Dimitrakopoulos et al. (2010) have shown
31 how higher order spatial cumulants of random variables can capture features of dense
32 training images that are not compatible with the assumption of an underlying multivariate-
33 Gaussian variable. The objective of the present paper is to show how such a cumulant
34 can be used in an inferential framework to test the strength of evidence against the null
35 hypothesis that, possibly relatively sparse, observations are drawn from a variable in which
36 these cumulants take values expected in the Gaussian case; and to identify exploratory
37 statistics that might be used to judge whether a Gaussian assumption is plausible. The
38 approach is illustrated using data on porosity of an important sedimentary formation under
39 the North Sea. A sound spatial stochastic model for this variable is necessary because the
40 pore-space in this unit may be important as a site for future carbon capture and storage
41 (Holloway, 2009).

42 2. Cumulants

43 A real-valued random variable, Z , with a probability density function $f_Z(z)$, has a
44 moment-generating function:

$$M(v) = E[\exp\{vZ\}] = \int_{-\infty}^{\infty} \exp\{vz\}f_Z(z)dz. \quad (1)$$

45 If $M(v)$ has a Taylor series expansion about the origin then it may be written as

$$M(v) = E[\exp\{vZ\}] = E\left[1 + vZ + \frac{v^2}{2!}Z^2 + \dots + \frac{v^r}{r!}Z^r + \dots\right]. \quad (2)$$

46 Note that the r th non-centred moment of Z ,

$$\mu'_r = \text{E}[Z^r],$$

47 is the coefficient of $\frac{v^r}{r!}$ in the r th term in this expansion, hence the name of the function.

48 Cumulants of the random variable may be defined in a similar and related way. The
49 cumulant generating function is

$$K(v) = \ln(\text{E}[\exp\{vZ\}]),$$

50 so we may write

$$1 + \mu'_1 \frac{v}{1!} + \mu'_2 \frac{v^2}{2!} + \dots + \mu'_r \frac{v^r}{r!} + \dots = \exp \left\{ \kappa_1 \frac{v}{1!} + \kappa_2 \frac{v^2}{2!} + \dots + \kappa_r \frac{v^r}{r!} + \dots \right\}, \quad (3)$$

51 where κ_r is the r th cumulant of Z .

52 The cumulants and moments of a distribution are related, for example (Kendall and
53 Stuart, 1977)

$$\begin{aligned} \mu'_1 &= \kappa_1, \\ \mu'_2 &= \kappa_1^2 + \kappa_2, \\ \mu'_3 &= \kappa_1^3 + 3\kappa_1\kappa_2 + \kappa_3. \end{aligned} \quad (4)$$

54 However, cumulants have certain properties which can make them more useful than
55 moments. In particular they generalize simply to the multivariate case (McCullagh and
56 Kolassa, 2009). Consider an n -variate random variate $\mathbf{Z} = \{Z_1, Z_2, \dots, Z_n\}$. One may
57 define entries in the mean vector of Z , the matrix of second non-centred moments and the
58 array of non-centred third moments as

$$\begin{aligned} \mathbf{E}_r &= \text{E}[Z_r] \\ \mathbf{E}_{rs} &= \text{E}[Z_r Z_s] \\ \mathbf{E}_{rst} &= \text{E}[Z_r Z_s Z_t] \end{aligned} \quad (5)$$

59 We denote linear combinations of the variables in \mathbf{Z} , and the powers of this term using
 60 Einstein's simplified convention for notation of multiple summations (Kuptsov, 2001):

$$v_r Z_r \equiv \sum_{r=1}^n v_r Z_r \quad (6)$$

61

$$v_r v_s Z_r Z_s \equiv \sum_{r=1}^n \sum_{s=1}^n v_r v_s Z_r Z_s = (v_r Z_r)^2 \quad (7)$$

62 where the term $v_r Z_r$ on the right is defined in Eq [6],

$$v_r v_s v_t Z_r Z_s Z_t \equiv \sum_{r=1}^n \sum_{s=1}^n \sum_{t=1}^n v_r v_s v_t Z_r Z_s Z_t = (v_r Z_r)^3 \text{ etc.} \quad (8)$$

63 Given this notation, the multivariate moment-generating function can be expanded as

$$M(v) = 1 + v_r E_r + \frac{v_r v_s E_{rs}}{2!} + \dots \quad (9)$$

64 and, similarly,

$$K(v) = \ln(M(v)) = \kappa^r \frac{v_r}{1!} + \kappa^{r,s} \frac{v_r v_s}{2!} \dots \quad (10)$$

65 As in the univariate case, the cumulants of increasing order, $\kappa^r, \kappa^{r,s}, \dots$ appear as coeffi-
 66 cients in the expansion. The moments and cumulants in the multivariate case are found
 67 to be related in a simple way, the moments of some order are given by the sum of products
 68 of cumulants over partitions of the superscripts so, for moments and cumulants of order
 69 up to three:

$$E_r = \kappa^r, \quad (11)$$

70

$$E_{rs} = \kappa^{r,s} + \kappa^r \kappa^s, \quad (12)$$

71 and

$$E_{rst} = \kappa^{r,s,t} + \kappa^{r,s} \kappa^t + \kappa^{r,t} \kappa^s + \kappa^{s,t} \kappa^r + \kappa^r \kappa^s \kappa^t. \quad (13)$$

72 The expressions above can be rearranged to express the cumulant of order k as functions
 73 of moments of order $m \leq k$ and cumulants of order $< k$:

$$\kappa^r = \kappa^r, \quad (14)$$

$$\kappa^{r,s} = \mathbf{E}_{rs} - \kappa^r \kappa^s, \quad (15)$$

75 and, rearranging Eq [13] and substituting Eq [15] for the second-order cumulants,

$$\begin{aligned} \kappa^{r,s,t} &= \mathbf{E}_{rst} - \kappa^{r,s} \kappa^t - \kappa^{r,t} \kappa^s - \kappa^{s,t} \kappa^r - \kappa^r \kappa^s \kappa^t, \\ &= \mathbf{E}_{rst} - [\mathbf{E}_{rs} - \kappa^r \kappa^s] \kappa^t - [\mathbf{E}_{rt} - \kappa^r \kappa^t] \kappa^s - [\mathbf{E}_{st} - \kappa^s \kappa^t] \kappa^r - \kappa^r \kappa^s \kappa^t, \\ &= \mathbf{E}_{rst} - \mathbf{E}_{rs} \kappa^t - \mathbf{E}_{rt} \kappa^s - \mathbf{E}_{st} \kappa^r + 2\kappa^r \kappa^s \kappa^t, \\ &= \mathbf{E}_{rst} - \mathbf{E}_{rs} \mathbf{E}_t - \mathbf{E}_{rt} \mathbf{E}_s - \mathbf{E}_{st} \mathbf{E}_r + 2\mathbf{E}_r \mathbf{E}_s \mathbf{E}_t. \end{aligned} \quad (16)$$

76 For zero mean \mathbf{Z} Eq[15] and Eq[16] simplify to

$$\kappa^{r,s} = \mathbf{E}_{rs} = \text{Cov}[Z_r, Z_s], \quad (17)$$

77 where $\text{Cov}[\cdot, \cdot]$ denotes the covariance of the terms in the brackets, and

$$\kappa^{r,s,t} = \mathbf{E}_{rst}, \quad (18)$$

78 i.e. the third cumulant is equal to the third moment. This is zero for multivariate Gaus-
79 sian \mathbf{Z} . In fact all multivariate cumulants of order $m > 2$ are zero for the Gaussian
80 case (Bilodeau and Brenner, 1999). This is demonstrated for the fourth cumulant in the
81 appendix.

82 Dimitrakopoulos et al. (2010) describe the extension of multivariate cumulants to
83 the spatial random field $\mathbf{Z}(\mathbf{x})$. Consider the third-order cumulant. Given some location
84 \mathbf{x} we may define a set of three locations $\{\mathbf{x}, \mathbf{x} + \mathbf{h}_1, \mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2\}$ where \mathbf{h}_1 and \mathbf{h}_2 are
85 lag vectors such that $\mathbf{h}_1 = h_1 \mathbf{l}_1$ and $\mathbf{h}_2 = h_2 \mathbf{l}_2$ where h_1 and h_2 are scalar lag distances
86 and \mathbf{l}_1 and \mathbf{l}_2 are lag vectors of unit length. Note that this notation is somewhat different
87 to that of Dimitrakopoulos et al. (2010). Given such a configuration, and making the
88 ergodicity assumption that the distribution of $Z(\mathbf{x})$ is independent of \mathbf{x} , we may express

89 the third-order cumulant for the random field at these locations as a function of lag only:

$$\begin{aligned}
\kappa^3(\mathbf{h}_1, \mathbf{h}_2) &= \mathbb{E}[Z(\mathbf{x})Z(\mathbf{x} + \mathbf{h}_1)Z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2)] \\
&\quad - \mathbb{E}[Z(\mathbf{x})]\mathbb{E}[Z(\mathbf{x} + \mathbf{h}_1)Z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2)] \\
&\quad - \mathbb{E}[Z(\mathbf{x} + \mathbf{h}_1)]\mathbb{E}[Z(\mathbf{x})Z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2)] \\
&\quad - \mathbb{E}[Z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2)]\mathbb{E}[Z(\mathbf{x} + \mathbf{h}_1)Z(\mathbf{x})] \\
&\quad + 2\mathbb{E}[Z(\mathbf{x})]\mathbb{E}[Z(\mathbf{x} + \mathbf{h}_1)]\mathbb{E}[Z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2)], \tag{19}
\end{aligned}$$

90 given Equation (16) When $\mathbf{Z}(\mathbf{x})$ is a zero mean spatial field this simplifies to

$$\kappa^3(\mathbf{h}_1, \mathbf{h}_2) = \mathbb{E}[Z(\mathbf{x})Z(\mathbf{x} + \mathbf{h}_1)Z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2)]. \tag{20}$$

91 Note from the discussion above that, for a Gaussian random field, the cumulants
92 $\kappa^r(\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{r-1})$ for any lags and for $r > 2$ are zero. This does not depend on assump-
93 tions of ergodicity.

94 As proposed by Dimitrakopoulos et al. (2010) cumulants may be estimated for
95 specified lag combinations, such as $\mathbf{h}_1, \mathbf{h}_2$, by considering all sets of observations whose
96 locations are translations of the basic template $[\{0, 0\}, \mathbf{h}_1, \mathbf{h}_1 + \mathbf{h}_2]$. When observations
97 are not regularly spaced it is necessary, as with estimation of the empirical variogram,
98 to compute estimates for lag bins which allow for some variation or tolerance about a
99 central lag. Under the assumption of ergodicity (at least up to the order of the cumulant
100 of interest), the estimator for the third cumulant of a zero-mean random variable from a
101 set of observations at locations X is therefore

$$\widehat{\kappa^3}(\mathbf{h}_1, \mathbf{h}_2) = \frac{1}{N(\mathbf{h}_1, \mathbf{h}_2)} \sum_{\{\mathbf{x}, \mathbf{x} + \mathbf{h}_1, \mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2\} \in X} z(\mathbf{x})z(\mathbf{x} + \mathbf{h}_1)z(\mathbf{x} + \mathbf{h}_1 + \mathbf{h}_2), \tag{21}$$

102 where there are $N(\mathbf{h}_1, \mathbf{h}_2)$ sets of observations whose locations are translations of the basic
103 template $[\{0, 0\}, \mathbf{h}_1, \mathbf{h}_1 + \mathbf{h}_2]$.

104 3. Materials and Methods

105 3.1. Data on the Bunter Sandstone porosity.

106 The data used in this study are all from the Bunter Sandstone formation. The Bunter
107 Sandstone is a sheet-sand complex comprising mainly fine-grained but locally medium- or
108 coarse-grained material (Cameron et al., 1992). It was deposited as fluvial channel sands in
109 arid conditions in the lower Triassic. The Bunter Sandstone is a significant formation in the
110 North Sea and corresponds to the Sherwood Sandstone group onshore. It is an important
111 gas reservoir in the North Sea and is potentially important for carbon capture and storage
112 (Holloway, 2009; Senior, 2010). For this reason the porosity of the Bunter Sandstone is of
113 interest. The porosity of this material is affected by various factors including the structure
114 of the original sediments, the depositional overburden, cementation of the material and
115 subsequent diagenetic transformation (Bifani, 1986).

116 The data are derived from analysis of cores extracted from 32 wells across the North
117 Sea. The cores were of variable length, and were sampled by extracting plug samples of one
118 inch diameter, the diameter of the plug being in the vertical direction. The recorded depth
119 of the plug was at its centre. The samples were not collected at absolutely regular intervals,
120 the mean spacing down-core between successive samples was 0.6 m. Where coherent plugs
121 could not be extracted a comparable volume of chipped material was removed. Each
122 sampled specimen was washed to remove all hydrocarbons and oven-dried to a constant
123 weight before porosity was determined by helium porosimetry. These are the best data
124 available on the porosity of the Bunter Sandstone, but it is acknowledged that there may
125 be some observational errors due to dissolution of halite cements during washing of the
126 samples (Ketter, 1991). The analyses reported in this paper are limited to porosity data
127 from plugs in water-filled sections of the cores, excluding results from gas-filled material.
128 A total of 1282 measurements from the 32 cores were available.

129 *3.2. Calculations.*

130 *3.2.1. Exploratory analysis and linear mixed model.* The number of wells is too small to
131 allow spatial modelling of the lateral variability of porosity in this formation. For this

132 reason a linear mixed model of the following form was fitted for exploratory purposes

$$Z(i, x) = \mu + K_i + \eta(i, x), \quad (22)$$

133 where $Z(i, x)$ is a random variable: the porosity at depth x within the i th well. Note that
 134 we define locations within wells by scalar depths, effectively the data within any well are
 135 in one dimension. The mean porosity over all depths and wells is μ , K_i is a random effect
 136 drawn from a random variable with mean zero and variance σ_B^2 ; it represents the difference
 137 between the mean porosity for the i th well and the overall mean porosity. The term $\eta(i, x)$
 138 is also a random effect of mean zero and variance, σ_W^2 . This random effect accounts for
 139 the within-well variability. The covariance of the values of η at any two depths in the
 140 same borehole is

$$\begin{aligned} \text{Cov} [\eta(i, x), \eta(i, x')] &= \sigma_W^2, \quad x = x' \\ &= (1 - \xi) \sigma_W^2 R(|x - x'|; \boldsymbol{\psi}), \quad x \neq x' \end{aligned} \quad (23)$$

141 where $R(\cdot; \boldsymbol{\psi})$ is a correlation function with parameters in $\boldsymbol{\psi}$ and $\xi \in [0, 1]$ is the nugget
 142 ratio, the proportion of the variance of η which is not correlated at spatial scales resolved
 143 by the sampling. This may include measurement error. Because the argument of the
 144 correlation function is the distance between two locations within a borehole rather than
 145 two absolute positions, the correlation structure is said to be second-order stationary
 146 (Journel and Huijbregts, 1977). Various correlation functions may be considered, provided
 147 that they guarantee a positive definite correlation matrix for η at any set of unique sites.
 148 One such function is the exponential:

$$R(|x - x'|; [r]) = \exp\{-|x - x'|/r\}, \quad (24)$$

149 with r , a distance parameter the only element in $\boldsymbol{\psi}$. An alternative is the spherical
 150 function:

$$\begin{aligned} R(|x - x'|; [a]) &= 0 \quad a > |x - x'|, \\ &= 1 - \frac{3|x - x'|}{2a} + \frac{1}{2} \left(\frac{|x - x'|}{a} \right)^3 \quad a \leq |x - x'|, \end{aligned} \quad (25)$$

151 for which a is the distance parameters, the range of the covariance function. Under these
152 correlation models the term η at two locations in a borehole are expected to be more
153 similar the closer they are in space.

154 The variance parameters of the linear mixed model in Eq. [22] — the variances σ_B^2
155 and σ_W^2 , the nugget ratio ξ and the terms in ψ — are best estimated by residual maximum
156 likelihood (REML) (Verbeke and Mohlenbergs, 2000). This entails the assumption that
157 the random effects can plausibly be regarded as realizations of a normal random field. In
158 the context of this study we examined the plausibility of this assumption (which we know
159 cannot be strictly true because porosity is bounded in the interval [0,100]), by examining
160 the marginal distribution of the residuals from an ordinary least squares fit of the LMM.
161 Exploratory statistics were computed for the residuals, including the robust measure of
162 skewness, the octile skew, proposed by Brys et al. (2003). Because porosity is a proportion,
163 as noted above, we repeated this exploratory analysis after a logistic transformation of the
164 porosities. Finally, the parameter of a Box-Cox transformation was estimated by maximum
165 likelihood by means of the BOXCox procedure in the MASS package for the R platform
166 (Venables and Ripley, 2002) and exploratory analysis was undertaken on residuals after
167 this transform. Results are presented below, but the following procedures may be followed
168 on the basis either that the residuals appear to have a reasonably normal distribution or
169 that this is plausible after an appropriate transformation.

170 The parameters of the linear mixed model were then estimated by REML. The lme
171 procedure in the NLME library for R (Pinheiro et al., 2013; R Development Core Team,
172 2010) was used, and spherical and exponential correlation functions for η were considered.
173 The variance parameters for η were tested by cross-validation. Each residual from the
174 well mean was removed from the data set in turn and predicted by ordinary kriging from
175 the remaining values in the same well. This was done using the XVOK2D algorithm in the
176 GSLIB library (Deutsch and Journel, 1997). For each observation, $\eta(i, x)$ this provides a
177 kriging estimate, $\tilde{\eta}(i, x)$, and the prediction error variance (kriging variance) $\sigma_K^2(i, x)$. A

178 useful diagnostic (Lark, 2009) is the standardized squared prediction error, with mean one
 179 and median 0.455 for normal kriging errors when the variance parameters are correct:

$$\theta(i, x) = \frac{\{\eta(i, x) - \tilde{\eta}(i, x)\}^2}{\sigma_K^2(i, x)}. \quad (26)$$

180 The linear mixed modelling framework was used to test the hypothesis that porosity
 181 depends on depth down the well. Neither exploratory plots of the data nor these models
 182 provided any evidence for a trend in porosity with depth, and so I proceeded with the
 183 model in Equation [22] where the mean porosity is constant within any well.

184 *3.2.2 Estimating κ^3 for particular templates.* For a zero-mean ergodic random variable
 185 $\eta(i, x)$ on a set of one-dimensional wells, K , the third-order cumulant, defined for a random
 186 field in Eq. [20], is defined for scalar lag distances h_1 and h_2 by

$$\begin{aligned} \kappa_\eta^3(h_1, h_2) = & \quad (27) \\ \mathbb{E}[\eta(i, x_1)\eta(i, x_2)\eta(i, x_3); |x_2 - x_1| = h_1, |x_3 - x_2| = h_2, (x_2 - x_1)(x_3 - x_2) > 0]_{i \in K}. \end{aligned}$$

187 Note that under this definition the locations are in order x_1, x_2, x_3 up or down the well,
 188 and the cumulant is symmetric in the sense that $\kappa_\eta^3(h_1, h_2) = \kappa_\eta^3(h_2, h_1)$.

189 In practice, when sampling is not on a regular array, it is necessary to allow some
 190 tolerance in the definition of the lag distances (Dimitrakopoulos et al., 2010). In this study
 191 we define a scalar-lag class \tilde{h} as the interval $[h - \tau, h + \tau]$ where τ is the tolerance. We
 192 define the indicator variable

$$\begin{aligned} I(i, x_1, x_2, x_3; \tilde{h}_1, \tilde{h}_2) &= 1 \quad i \in K, |x_2 - x_1| \in \tilde{h}_1, |x_3 - x_2| \in \tilde{h}_2, (x_2 - x_1)(x_3 - x_2) > 0 \\ &= 0 \quad \text{otherwise.} \end{aligned} \quad (28)$$

193 We then define the estimate $\widehat{\kappa}_\eta^3(h_1, h_2)$ by

$$\begin{aligned} \widehat{\kappa}_\eta^3(h_1, h_2) = & \quad (29) \\ \frac{1}{N_{h_1, h_2}} \sum_{i \in K} I(i, x_1, x_2, x_3; \tilde{h}_1, \tilde{h}_2) \{z(i, x_1) - \bar{z}_i\} \{z(i, x_2) - \bar{z}_i\} \{z(i, x_3) - \bar{z}_i\}, \end{aligned}$$

194 where $z(i, x_1)$ is the observed value of the variable at depth x_1 in the i th well, and \bar{z}_i is
 195 the average value of the variable over all observations in the i th well. The summation is

196 over all sets of three observations within all wells in the set K and N_{h_1, h_2} is the sum of
 197 the indicator over all these observations.

198 In this study the estimate $\widehat{\kappa}_\eta^3(h_1, h_2)$ was computed for lag distances 25 cm, 50 cm, \dots , 500 cm
 199 with lag tolerance $\tau = 12.5$ cm.

200 *3.2.3. Testing $\widehat{\kappa}^3$ against a null hypothesis of normality.* As noted above the expected
 201 value of the third cumulant for a multivariate normal random variable is zero. Values of
 202 $\widehat{\kappa}^3$ for some h_1, h_2 provide evidence against this null hypothesis, but this evidence must be
 203 assessed accounting for the sample variance of the estimates. This is complicated by the
 204 lack of independence of the observations from which the estimate is obtained, so a Monte
 205 Carlo simulation procedure was developed.

206 Under the null hypothesis of multivariate normality the variability of the data is
 207 entirely accounted for by the variances and associated parameters of the random effects
 208 in the linear mixed model, Eq [22]. The Monte Carlo procedure requires that we can
 209 generate realizations of the random term η from the linear mixed model. We denote the
 210 set of values of this random variable by the $N \times 1$ vector $\boldsymbol{\eta}$ which corresponds to the full
 211 set of N observations. The covariance matrix of the random variate $\boldsymbol{\eta}$ is denoted by \mathbf{V}
 212 where

$$\mathbf{V} = \xi \sigma_W^2 \mathbf{I} + (1 - \xi) \sigma_W^2 \mathbf{R}, \quad (30)$$

213 where \mathbf{I} is a $N \times N$ identity matrix and \mathbf{R} is an $N \times N$ correlation matrix such that the
 214 entry $\mathbf{R}\{k, l\}$ for the l th observation $\eta(i, d)$ and the k th $\eta(j, d')$ is:

$$\begin{aligned} \mathbf{R}\{k, l\} &= 0, \quad \forall i \neq j \\ &= R(|d - d'|; \boldsymbol{\psi}), \quad \forall i = j, \end{aligned} \quad (31)$$

215 where R is a correlation function with parameters in $\boldsymbol{\psi}$. In this study the correlation
 216 function fitted by REML, and the estimated parameters were used. Once \mathbf{V} has been
 217 computed it is possible to find its Cholesky factorization:

$$\mathbf{V} = \mathbf{L}\mathbf{L}^*, \quad (32)$$

218 where \mathbf{L} is a lower-triangular matrix with real and positive diagonal elements and \mathbf{L}^* is its
 219 conjugate transpose. This factorization is guaranteed to exist because the matrix \mathbf{R} , as a
 220 covariance matrix computed from an authorized correlation function, is positive-definite
 221 and symmetric with real values. It is then possible to generate a realization of $\boldsymbol{\eta}$ by
 222 computing

$$\boldsymbol{\eta} = \mathbf{L}\mathbf{g}, \quad (33)$$

223 where the elements of \mathbf{g} are independent values with a standard normal distribution.

224 In this study the IMSL subroutine CHFAC was used to compute the Cholesky fac-
 225 torization. One may then substitute the elements of $\boldsymbol{\eta}$ for the values of z in Eq. [29] to
 226 compute $\widehat{\kappa}_{\boldsymbol{\eta}}^3(h_1, h_2)$ for the same lag distances for which this was computed for the original
 227 data. It is immaterial that the between-well random effect is not simulated here since
 228 the mean value for each well is subtracted from each observation in Eq. [29]. Since $\boldsymbol{\eta}$
 229 is simulated for the same locations as the data, the value of $\widehat{\kappa}_{\boldsymbol{\eta}}^3(h_1, h_2)$ for some lag dis-
 230 tances computed from the simulated data can be regarded as a realization of the sampling
 231 distribution of our observed statistic under the null hypothesis of a multivariate normal
 232 distribution. Note also that the sample error of each well mean, which contributes to the
 233 error of the estimation of $\widehat{\kappa}_{\boldsymbol{\eta}}^3$ which is estimated on the assumption of zero mean, also
 234 appears in the simulation procedure and so is included in the Monte Carlo approximation
 235 to the sampling distribution of $\widehat{\kappa}_{\boldsymbol{\eta}}^3$. In this study 100 000 realizations of $\boldsymbol{\eta}$ were generated
 236 and used to compute the sampling distribution of $\widehat{\kappa}_{\boldsymbol{\eta}}^3(h_1, h_2)$ for the specified lags under
 237 the null hypothesis.

238 Two approaches were used to examine the extent to which the empirical cumulants
 239 of the data are consistent or otherwise with a null hypothesis of normality. The first was
 240 to find the maximum absolute value of the estimated cumulants over all lag distances,

$$\widehat{\kappa}_{\boldsymbol{\eta}, \max}^3 = \max \left\{ \left| \widehat{\kappa}_{\boldsymbol{\eta}}^3(h_1, h_2) \right| ; h_1 = 25, 50, \dots, 500\text{cm}; h_2 = 25, 50, \dots, 500\text{cm} \right\}. \quad (34)$$

241 This statistic was evaluated for the empirical residuals from the well means, and
 242 then for each of a set of 100 000 realizations of $\boldsymbol{\eta}$, generated as described above. Since the

243 expected value of the cumulant under the null hypothesis of a multivariate Gaussian ran-
 244 dom variable is zero a large value of $\widehat{\kappa}_{\eta, \max}^3$ provides evidence against this null hypothesis.
 245 The strength of evidence is measured by a p -value which can be approximated by ordering
 246 the values of $\widehat{\kappa}_{\eta, \max}^3$ from the simulations and computing the proportion of these which
 247 exceed the observed value.

248 The second approach was to test the separate cumulants for each lag pair h_1, h_2 . For
 249 some observed lag pair at which the observed cumulant is $\widehat{\kappa}_{\eta}^3(h_1, h_2)$ the p -value for the null
 250 hypothesis of a zero cumulant is computed by finding the proportion of the 100 000 realiza-
 251 tions of $\boldsymbol{\eta}$ for which the cumulant fall outwith the interval $\left[-\left|\widehat{\kappa}_{\eta}^3(h_1, h_2)\right|, +\left|\widehat{\kappa}_{\eta}^3(h_1, h_2)\right|\right]$.
 252 These p -values were inspected for a set of lag combinations, excluding those with fewer
 253 than 600 supporting triplets of observations. This is a multiple hypothesis test, in which
 254 we examine a family of null hypotheses which are not mutually independent. For that
 255 reason it is necessary to control the family-wise error rate (FWER), α_{rmFW} , which is the
 256 probability of one or more of the family of null hypotheses' being rejected although all
 257 of them are true. The simplest way to control the family-wise error rate for a set of m
 258 hypotheses is to reject only those for which $p < \alpha_{FW}/m$. This is the Bonferroni control
 259 of FWER, and is valid for non-independent hypotheses (Snedecor and Cochran, 1980).
 260 However, it is relatively lacking in power. An alternative, also valid for non-independent
 261 hypotheses, is the procedure due to Holm (1979). In Holm's procedure one orders the
 262 null hypotheses H_1, H_2, \dots, H_m in order of ascending p -value, p_1, p_2, \dots, p_m . One then
 263 evaluates for successive $k = 1, 2, \dots, m$ whether

$$p_k > \frac{\alpha_{FW}}{m + 1 - k}.$$

264 Let k_r be the smallest value of k for which this expression is true. One may then reject,
 265 with FWER α_{FW} , the null hypotheses $H_1, H_2, \dots, H_{k_r-1}$. This procedure was followed to
 266 find the subset of lag pairs for which the null hypothesis that the cumulant is zero could
 267 be rejected.

268 *3.2.4 Exploring the implications of a non-zero cumulant.* In order to gain insight into

269 the nature of the variability of a variable with non-zero third order cumulants for lag-
270 pairs h_1, h_2 I examined 3-D plots of the triplets of observations $\{z(d), z(d + h_1), z(d + h_2)\}$
271 using the SCATTERPLOT3D package in R. This is comparable to the examination of two-
272 dimensional scatterplots of $\{z(\mathbf{x}), z(\mathbf{x} + \mathbf{h})\}$ which is sometimes advocated as an exploratory
273 technique in geostatistics (Goovaerts, 1997).

274 4. Results

275 Table 1 presents summary statistics for residuals for porosity from the well mean,
276 and the same residuals for data after logistic or Box-Cox transformation. Note that there
277 is little appreciable effect of the Box-Cox transformation, and the 95% confidence interval
278 of the Box-Cox parameter included the value 1, under which the transform is equivalent
279 to adding a constant to the variable and has no effect on the shape of the distribution.
280 The residuals after a logistic transform are more skewed than in the other two cases. All
281 of these exploratory statistics suggest that an assumption of normality of the residuals
282 with no transformation seemed plausible. Figure 1 shows the histogram of these residuals
283 and their empirical Quantile-Quantile plot which should lie on the bisector.

284 Table 2 shows the results of the REML estimation of the variance parameters for
285 the linear mixed model for porosity set out in Eq[22]. Figure 2 shows the histogram of
286 cross-validation errors for the selected model (exponential) and the Q-Q plot. These show
287 that the errors are close to normal in their distribution. The mean and median standard
288 square cross validation errors are in Table 2. Note that the mean is close to 1.0, but the
289 median is rather smaller than is expected.

290 It was found that the numbers of triplets of observations from which to estimate the
291 cumulant for particular lag pairs N_{h_1, h_2} varied. For most pairs of lags there were between
292 600 and 1600 triplets, so those lags supported by fewer observations were discarded. Figure
293 3 shows the estimated values $\widehat{\kappa}_\eta^3(h_1, h_2)$ which are plotted only in the lower half of the
294 plot (where $h_1 > h_2$). The dots in the upper half of the plot indicate the lags at which

295 the number of supporting triplets of observations was fewer than 600.

296 The largest absolute value of the third cumulant over the lags considered was 57.8
297 for lag-pair {50 cm, 250 cm}. Table 3 shows the percentiles of the maximum absolute value
298 of the third cumulant over 100 000 realizations of the Gaussian model, and also percentiles
299 of the third cumulant for lags {50 cm, 250 cm}. Figure 4 shows the approximate density
300 functions for (a) the maximum absolute value of the third cumulant over all lags and
301 (b) the third cumulant for lags {50 cm, 250 cm} from the 100 000 realizations. The
302 density was obtained by the KERNELDENSITY procedure in GenStat (Goedhart, 2009).
303 This, and the percentiles in Table 3, indicate that the cumulant is distributed more or less
304 symmetrically about zero under the null hypothesis of a multivariate Gaussian distribution.
305 The percentiles of the maximum absolute value of the third cumulant over all lags in Table
306 3 shows that the approximate p -value for the evidence provided by the absolute maximum
307 third cumulant for these data against a null hypothesis of normality is less than 0.01, but
308 larger than 0.001.

309 In the upper half of Figure 3 are plotted those cumulants which were significantly
310 different from zero as judged by the p -values computed for each lag pair from the 100 000
311 realizations, with FWER controlled at 0.05. There are six lag pairs at which the cumulants
312 are significantly non-zero. Note that the significant cumulants are negative for smaller lags
313 — {50 cm, 250 cm}, {50 cm, 225 cm} and {100 cm, 225 cm} — and positive for the longer
314 lags, {50 cm, 425 cm}, {150 cm, 275 cm} and {275 cm, 500 cm}.

315 Three-dimensional scatter-plots were examined for data triplets (residuals from the
316 well mean) with the smallest (most negative) and largest (most positive) cumulant, cor-
317 responding to lags {50 cm, 250 cm} and {150 cm, 275 cm} respectively. I do not attempt
318 to reproduce them here but the effects that they show can be illustrated by two two-
319 dimensional plots of residuals for two locations, x_1 and $x_2 = x_1 + h_1$ with, respectively
320 $\eta(x_3) > 0$ and $\eta(x_3) \leq 0$ where $x_3 = x_2 + h_2$. These plots are shown in Figure 5, along
321 with the correlations between the variables on the plots. Note that the ‘positive quad-

322 rants' of the plot, where $\eta(x_1)\eta(x_2) > 0$, have been given a grey background. It is apparent
323 that the correlation between $\eta(x_1)$ and $\eta(x_2)$ differ between the cases where $\eta(x_3) > 0$ and
324 $\eta(x_3) \leq 0$, and these differences are significant in each case with $p < 0.001$. This difference
325 in correlation is an expression of the non-zero third cumulant of $\boldsymbol{\eta}$ which has been found
326 for these lag pairs, since it means that distribution of observations between the positive
327 and negative quadrants of these plots is different for the case where $\eta(x_3) > 0$ and where
328 $\eta(x_3) \leq 0$. Furthermore, this difference in correlation is inconsistent with the assumption
329 of second-order stationarity under which the correlation between $\eta(x_1)$ and $\eta(x_2)$ should
330 depend only on h_1 .

331 5. Discussion

332 In the work above five general results were obtained from the exploratory analysis
333 of the porosity data.

- 334 1. Summary statistics and histograms on the marginal distribution of the data, includ-
335 ing after transformation. (Figure 1, Table 1).
- 336 2. A plot of the third cumulant of the centred data for a range of lags (Figure 3).
- 337 3. P -values for tests of the null hypothesis of an underlying multivariate Gaussian
338 process based on the third cumulants.
- 339 4. Scatter plots of data triplets and associated correlations (Figure 4).
- 340 5. Results from the cross validation of the fitted linear mixed model (Figure 2, Table
341 2).

342 The significance tests on the cumulants — item (3) in the list above — allow us to re-
343 ject the null hypothesis of an underlying multivariate Gaussian random variable. Whether
344 this is, of itself, of direct practical relevance is open to debate. Webster and Oliver (2007)
345 suggest that significance tests for conformity to distributions are not particularly valuable
346 for the purpose of assessing the plausibility of distributional assumptions. We know in

347 most cases that a variable is not strictly normally distributed, and, particularly with large
348 data sets, we do not expect the null hypothesis of normality to be accepted. For example,
349 with the data in this paper, we know that they cannot have a Gaussian distribution at the
350 limit since porosity is bounded on the interval $[0, 1]$. However, the exploratory statistics of
351 these data indicated that they are close to symmetrically distributed with a bell-shaped
352 histogram, and that neither the logistic nor the Box-Cox transformation improved this.
353 Following the guidelines of Webster and Oliver (2007) one would normally proceed on the
354 basis that a normality assumption is plausible.

355 How is this approach extended to the consideration of multivariate normality? The
356 plot of the cumulants (Figure 3) may indicate possible systematic deviations from the
357 expected value (zero), e.g. clustering of small (large negative) or large positive values at
358 particular lags, and the significance test indicates whether or not the general pattern is
359 compatible with sampling error from an underlying Gaussian process. The cumulant plot
360 also leads us to the particular data triplet plots which merit further investigation. These
361 triplet plots are visualizable projections of the data which allow us to see the particular
362 deviation from normality which the corresponding cumulant represents. In this case we
363 can identify notable differences between the correlation of $\eta(x_1)$ with $\eta(x_2)$ conditional
364 on the value of $\eta(x_3)$. This is not consistent with an assumption of stationarity in the
365 covariance. This is consistent with the cross-validation results, presented in Table 2. Note
366 that the median squared standard prediction error is rather less than the expected value
367 of 0.455. This may be due to the non-stationarity of the underlying variable, as found
368 by Lark (2009) in the comparison of kriging results from stationary and non-stationary
369 variance models. It is also possible that outlying data values could influence both the
370 squared standard prediction errors and estimates of the cumulants. The exploratory data
371 analysis did not indicate any marginal outliers in the data, but spatial outliers, values
372 unusual in their local context, may be present. One possible area for future work is to
373 develop robust estimators of the cumulants, but it would be necessary to find estimators

374 that do not import distributional assumptions through the use of particular consistency
375 corrections (Lark, 2000) while remaining reasonably efficient.

376 In short, the analysis of the third cumulants of the variable provides us with a basis
377 for identifying particular plots of the data which allow us to examine its deviation from a
378 stationary normal process directly, and to interpret other results such as those from the
379 cross-validation. I would agree with Webster and Oliver (2007) that, in general, we should
380 not base decisions about the validity of distributional assumptions on tests of conformity
381 to the particular distribution. Further work is required to develop exploratory statistics
382 based on the cumulants, which allow us to make an informed pragmatic judgement about
383 the plausibility of the distributional assumption. We require, for example, rules of thumb
384 such as that enunciated by Webster and Oliver (2007) that some transformation of data
385 is required if the coefficient of skewness exceeds 0.5. Such rules of thumb might be based
386 on plots of the cumulant such as Figure 3, and must be based on experience of a range of
387 data sets and the robustness of the Gaussian assumption when predicting or simulating
388 the measured variable.

389 Note that in the case study there was no evidence for any trend in porosity with
390 depth, and so it was assumed that the mean porosity in any well was constant. If a trend
391 was found then this would be subtracted from the observations before computation of the
392 cumulants, and the Monte Carlo procedure to approximate the sample distribution of the
393 cumulant under the null hypothesis would have to be extended to include the contribution
394 of the uncertainty in the estimation of the trend just as the reported procedure accounted
395 for the uncertainty in the estimation of the well means.

396 Given the sparsity of wells, and the distances between them, the current study
397 was limited to cumulants in one dimension, attention was also focussed on the third
398 cumulants. Any third cumulant in one dimension is defined for a lag pair, and so can
399 easily be displayed in 2-D plots. The extension of this method to higher-order cumulants,
400 to two or more dimensions, or both would make it harder to use visualization in the

401 analysis of data. However, the general principles used in this paper, for the estimation
402 of empirical cumulants and the use of multiple hypothesis testing methods to find lag-
403 combinations at which the data provide evidence against a multivariate Gaussian model,
404 could be extended to sets of more than two lag combinations in a straightforward way, and
405 so to higher-order cumulants and more than one dimension. Plots for visual interpretation
406 could then be generated as appropriate projections, in the same spirit of the triplet plots
407 used in this paper. Those considerations aside, the one-dimensional case illustrated here
408 remains of considerable relevance since many porosity or conductivity fields in geology can
409 only be examined intensively down-core. This is because of the relative sparsity of cores,
410 particularly offshore, and the fact that they are often widely spaced which limits the scope
411 to examine lateral variability.

412 These results give reason for concern about the suitability of prediction error vari-
413 ances and other measures of uncertainty based on the multivariate Gaussian model of
414 porosity in the Bunter Sandstone. It should also be recalled that regionalized variables with
415 non-Gaussian distributions may have more complex geometrical structure than Gaussian
416 variables, particularly with respect to the connectivity of extreme values (e.g. Guardiano
417 and Srivastava, 1993). This means that simulations of porosity fields from multivariate
418 Gaussian random variables, even if these well-reproduce the marginal statistics of porosity,
419 may fail to represent all aspects of the spatial structure of the variable (such as the vol-
420 umes of regions of continuous large or small porosity) which may be relevant to questions
421 of fluid flow or potential gas storage in the field.

422 One way to deal with this may be by copula methods (e.g. Haslauer et al, 2012),
423 although the development of appropriate spatial copula models other than the Gaussian
424 which can be fitted to sizeable data sets is at an early stage. An alternative is to use the
425 methods of multiple point geostatistical modelling,(e.g. Strebelle, 2001), but these require
426 large data sets for training. One solution would be to find a non-Gaussian stochastic model
427 which reproduces the cumulants of interest. A possible general form of the model would be

428 one in which a well is divided into intervals by randomly located boundaries (occurring as
 429 a Poisson process, so that the boundaries have an exponential distribution). The resulting
 430 segments of the well could be regarded as distinct geological facies. In the simplest such
 431 model all observations within any one of the segments thus-formed take a value drawn
 432 from a centred Gaussian random variable, Y . It is known (Lark, 2010) that this random
 433 field is not multivariate Gaussian (although its marginal distribution is). However, one
 434 can see that its third cumulant is zero since:

$$\kappa^3(\mathbf{h}_1, \mathbf{h}_2) = p_1(\mathbf{h}_1, \mathbf{h}_2)E[Y^3] + p_2(\mathbf{h}_1, \mathbf{h}_2)E[Y]E[Y^2] + p_3(\mathbf{h}_1, \mathbf{h}_2)E[Y]^3, \quad (35)$$

435 where $p_1(\mathbf{h}_1, \mathbf{h}_2)$ is the probability that all three locations in the template fall in different
 436 segments, $p_2(\mathbf{h}_1, \mathbf{h}_2)$ is the probability that two sites fall in one segment and one in another
 437 and $p_3(\mathbf{h}_1, \mathbf{h}_2)$ is the probability that all three locations fall into the same segment. These
 438 probabilities need not be evaluated since it is clear, from the fact that the variable is
 439 centred and Gaussian, so $E[Y] = E[Y^3] = 0$, that all three terms are zero. In a more
 440 complex version of this model one might postulate, for example, a correlation between
 441 the thickness of the segment and its expected porosity. In some preliminary simulations
 442 it was found that the resulting random variable may have a marginal distribution which
 443 appears Gaussian when the correlation between segment thickness and mean porosity is
 444 not too strong, but that the third cumulants were systematically smaller than zero for
 445 pairs of short lags (Figure 6). This is not offered as an alternative model for the Bunter
 446 Sandstone porosity, but simply as an indicator that the kind of spatial variation that has
 447 been found in reality might be reproduced by an appropriate stochastic model. This is a
 448 topic for further work, and should account for known general properties of the geological
 449 units. For example, while one might postulate relationships between grain size and facies
 450 thickness in depositional environments, porosity is also affected by overburden, diagenetic
 451 transformations of the sandstone and other processes which may be spatially dependent
 452 but are not obviously reproducible by a stochastic geometry.

453 **6. Conclusions**

454 It has been shown how the third cumulant of a spatial variable observed in linear
455 data sets (wells) can be used in an inferential context to test the null hypothesis that the
456 underlying distribution of the variable is multivariate Gaussian, and to guide exploratory
457 analysis to test the plausibility of this distributional assumption. This approach was
458 applied to data on porosity of the Bunter Sandstone and showed that there were features
459 of its distribution which appear incompatible with the assumption of stationarity and
460 multivariate Gaussian variation. This has potential implications for the use of standard
461 geostatistical methods to characterize the uncertainty that attends inferences about this
462 variable. This might require that multiple point geostatistics are used for this variable.
463 Alternatively some non-Gaussian random variable might be postulated as a model, and an
464 example of one which has some common features with the data is discussed. In practice
465 it might be possible to develop such a model for porosity of the Bunter Sandstone; such
466 a model should take account of our understanding of the depositional and diagenetic
467 processes that control this variable.

Acknowledgements

I am grateful to colleagues at BGS (Sam Holloway, Sarah Hannis) for guidance about the Bunter Sandstone and its properties, and to IHS who are owners of a subset of the porosity data for permission to use them. This paper is published with the permission of the Executive Director of the British Geological Survey (NERC).

References

- Bifani, R. 1986. Esmond Gas Complex. In: Brooks, J., Goff, J.C. & Van Hoorn, B. (eds) Habitat of Paleozoic Gas in NW Europe. Geological Society, London, Special Publications, 23, 209–221.
- Bilodeau, M., Brenner, D. 1999. Theory of Multivariate Statistics. Springer, New York.
- Brys, G., Hubert, M., Struyf, A. 2003. A comparison of some new measures of skewness. In: Developments in Robust Statistics (eds R. Dutter, P. Filzmoser, U. Gather and P.J. Rousseeuw), pp. 98–113. Physica-Verlag, Heidelberg.
- Cameron, T.D.J., Crosby, A., Balson, P.S., Jeffrey, D.H., Lott, G.K., Bulat, J., Harrison, D.J. 1992. United Kingdom offshore report: the geology of the southern North Sea. London, HMSO for the British Geological Survey.
- Deutsch, C.V., Journel, A.G. 1997. GSLIB: Geostatistical Software Library and User's Guide. Oxford University Press, New York.
- Dimitrakopoulos, R., Mustapha, H., Gloaguen, E. 2010. High-order statistics of spatial random fields: exploring spatial cumulants for modeling complex non-Gaussian and non-linear phenomena. *Mathematical Geoscience*, 42, 65–99.
- Goedhart, P.W. 2009. KERNELDENSITY procedure in (ed.) R.W. Payne, GenStat Release 12 Reference Manual, Part 3 Procedure Library PL20. VSN International, Hemel Hempstead.
- Goovaerts, P. 1997. Geostatistics for Natural Resources Evaluation. Oxford University Press, New York.
- Guardiano, F., Srivastava, R.M., 1993. Multivariate geostatistics: beyond bivariate moments. In: Soares, A. (Ed.), Geostatistics-Troia, Vol. 1. Kluwer, Dordrecht, pp. 133–144.

- Haslauer, C.P., Guthke, P., Bárdossy, A., Sudicky, E.A. 2012. Effects of non-Gaussian copula-based hydraulic conductivity fields on macrodispersion. *Water Resources Research*, 48, W07507.
- Holloway, S. 2009 Storage capacity and containment issues for carbon dioxide capture and geological storage on the UK continental shelf. *Proceedings of the Institution of Mechanical Engineers. Part A, Journal of Power and Energy*, 223 (A3), 239–248.
- Holm, S. 1979. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6, 65–70.
- Jansen, P.H.M., Stoica, P. 1988. On the expectation of the product of four matrix-valued Gaussian random variables. *IEEE Transactions on Automation and Control*, 33, 867–870.
- Journel, A.G., Huijbregts, C.J. 1978. *Mining Geostatistics*. Academic Press, London.
- Ketter, F.J. 1991. The Esmond, Forbes and Gordon Fields, Blocks 43/8a, 43/13a, 43/15a, 43/20a, UK North Sea. In: Abbotts, I.L. (ed.) *United Kingdom Oil and Gas Fields, 25 Years Commemorative Volume*. Geological Society, London, *Memoirs*, 14, 425–432.
- Kuptsov, L.P. 2001. Einstein's Rule. In: Hazewinkel, M. (Ed.) *Encyclopaedia of Mathematics*, Springer, Berlin.
- Lark, R.M. 2000. A comparison of some robust estimators of the variogram for use in soil survey. *European Journal of Soil Science*, 51, 137–157
- Lark, R.M. 2009. Kriging a soil variable with a simple non-stationary variance model. *Journal of Agricultural, Biological and Environmental Statistics*, 14, 301–321.
- Lark, R.M. 2010. Two contrasting spatial processes with a common variogram: inference about spatial models from higher-order statistics. *European Journal of Soil Science*, 61, 479–492.

- Lark, R.M, Bellamy, P.H. Rawlins, B.G. 2006. Spatio-temporal variability of some metal concentrations in the soil of eastern England, and implications for soil monitoring. *Geoderma* 133, 363–379.
- McCullagh, P., Kolassa, J. 2009. Cumulants. *Scholarpedia* 4, 4699.
- Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D., The R Development Core Team, 2013. nlme: Linear and Nonlinear Mixed Effects Models. R package version 3.1-110.
- R Development Core Team 2010. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org/>.
- Senior, B. 2010. CO₂ Storage in the UK — Industry Potential. Department Department for Energy and Climate Change (DECC) Report URN 10D/512, DECC, London.
- 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.
- Verbeke, G., Molenberghs, G. 2000. *Linear Mixed Models for Longitudinal Data*. Springer, New York.
- Visual Numerics, 2006. *IMSL Fortran Numerical Library Version 6.0*. Visual Numerics, Houston, Texas.
- Webster, R. & Oliver, M.A. 2007. *Geostatistics for Environmental Scientists*. 2nd Edition John Wiley & Sons, Chichester.

Appendix. The fourth cumulant of a multivariate-Gaussian random variable is zero.

Using the notation from section 2, and considering the zero-mean case for brevity of notation, the fourth multivariate cumulant can be written as

$$\kappa^{r,s,t,u} = E_{rstu} - \{E_{rs}E_{tu} + E_{rt}E_{su} + E_{ru}E_{st}\}, \quad (36)$$

see McCullagh and Kolassa, 2009. Now, for multivariate-Gaussian $\mathbf{Z} \equiv [Z_1, Z_2, Z_3, Z_4]$ the expected product $Z_1 Z_2 Z_3 Z_4$ is

$$\begin{aligned} E[Z_1 Z_2 Z_3 Z_4] &= \text{Cov}[Z_1, Z_2] \text{Cov}[Z_3, Z_4] \\ &+ \text{Cov}[Z_1, Z_3] \text{Cov}[Z_2, Z_4] \\ &+ \text{Cov}[Z_1, Z_4] \text{Cov}[Z_2, Z_3], \end{aligned} \quad (37)$$

because of the disappearance of odd-order moments, see, for example, Jansen and Stoica (1988). Note that the term in braces on the RHS of Eq. [36] is equivalent to the RHS of Eq. [37], from which it follows immediately that $\kappa^{r,s,t,u} = 0$.

Table 1. Summary statistics of residuals from mean well porosity using the original data, data after a logistic transformation and data after a Box-Cox transformation.

	Original data	After Box-Cox* transformation	After logistic transformation
Mean	0.00	0.00	0.00
Median	0.28	0.26	0.09
Skewness	-0.04	-0.09	-1.36
Standard deviation	6.61	5.25	0.64
Quartile 1	-4.44	-3.51	-0.69
Quartile 3	4.27	3.32	0.59
Octile skewness	-0.07	-0.08	-0.20

*The maximum-likelihood estimate of the Box-Cox transformation parameter was 0.92 with 95% confidence interval [0.83,1.01].

Table 2. Results from REML estimation of random effects parameters, and cross-validation.

Model	log-Likelihood	AIC	
Exponential	-4166.1	8342.3	
Spherical	-4176.8	8363.7	

Selected model			
Model	Random effects parameters		
	σ_B^2	σ_W^2	ξ
Exponential	23.59	45.23	18×10^{-9}

Cross-validation results	
Mean error	0.004
Mean standardized squared error	1.06
Median standardized squared error	0.32

Table 3. Quantiles of (a) Maximum value of the third cumulant over all lags, $\widehat{\kappa^3}_{\eta, \max}$; and (b) Value of the third cumulant for lag pair {50 cm, 250 cm}, $\widehat{\kappa^3}_{\eta}(50 \text{ cm}, 250 \text{ cm})$; computed from 100 000 realizations of the random model for $\boldsymbol{\eta}$.

$\widehat{\kappa^3}_{\eta, \max}$	
Quantile	Value
0.5	29.5
0.9	40
0.95	43.9
0.99	53.3
0.999	68.7

$\widehat{\kappa^3}_{\eta}(50 \text{ cm}, 250 \text{ cm})$	
Quantile	Value
0.001	-26.3
0.01	-19.4
0.05	-13.6
0.1	-10.5
0.5	0.0
0.9	10.5
0.95	13.7
0.99	19.8
0.999	27.1

Figure Captions

1. (a) Histogram of residuals from well mean porosity and (b) Gaussian Q-Q plot with bisector.
2. (a) Histogram of residuals cross-validation kriging errors and (b) Gaussian Q-Q plot with bisector.
3. Map of estimates, $\widehat{\kappa}_\eta^3(h_1, h_2)$ (below the diagonal). Symbols appear above the diagonal where the estimate was judged significantly different from zero. Small grey circles indicate where the estimate is supported by fewer than 600 triplets.
4. Estimated density functions for (top) the maximum absolute value of the third cumulant over all lag pairs under a null hypothesis of a multivariate Gaussian distribution and (bottom) the third cumulant for lag pair {50 cm, 250 cm}.
5. Scatter plots of data triplets for observations at x_1 and $x_2 = x_1 + h_1$ for (left) $\eta(x_3) < 0$ and (right) $\eta(x_3) > 0$, $x_3 = x_2 + h_2$. Top row, $h_1 = 50$ cm, $h_2 = 250$ cm; bottom row, $h_1 = 150$ cm, $h_2 = 275$ cm.
6. Map of estimates $\widehat{\kappa}_\eta^3(h_1, h_2)$ for a simulated random variable in which wells are divided randomly into segments and segment porosity is weakly correlated with segment thickness.

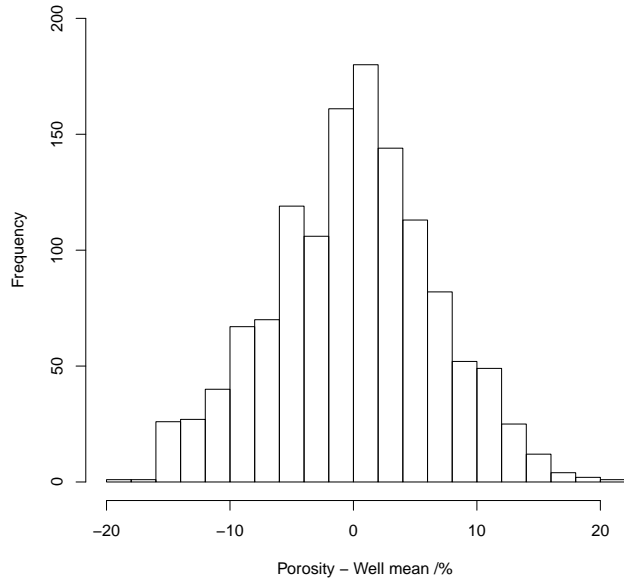


Figure 1(a)

Normal Q-Q Plot

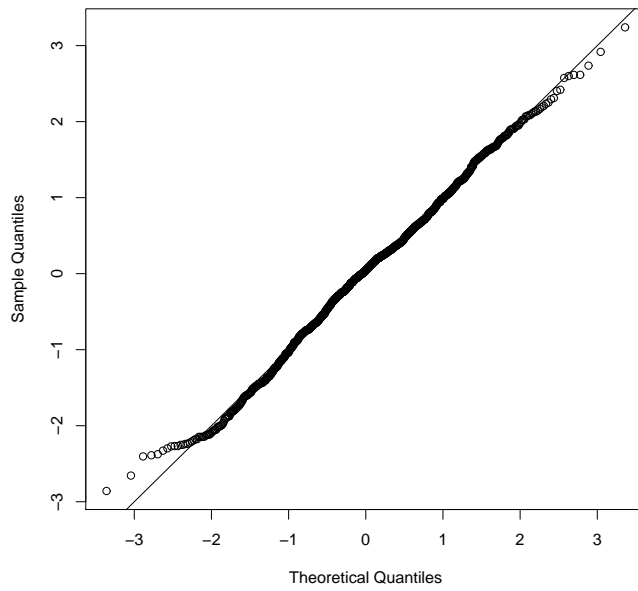


Figure 1(b)

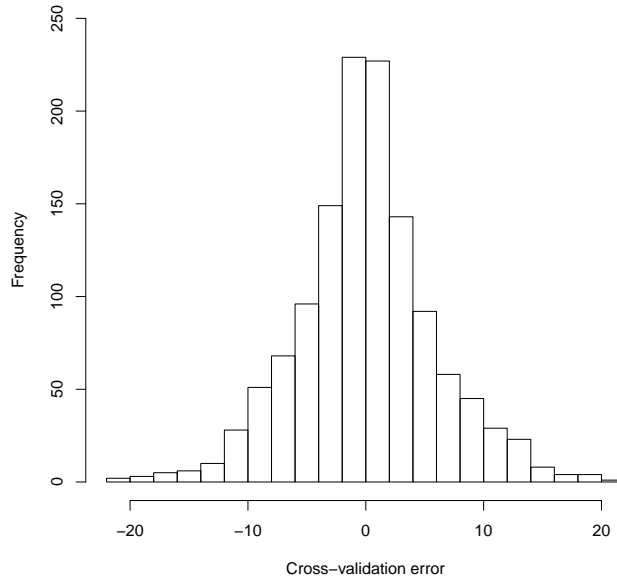


Figure 2(a)

Normal Q-Q Plot

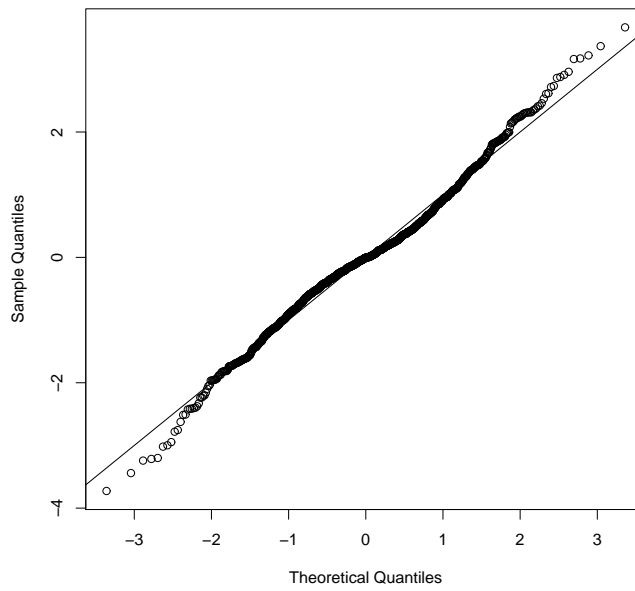


Figure 2(b)

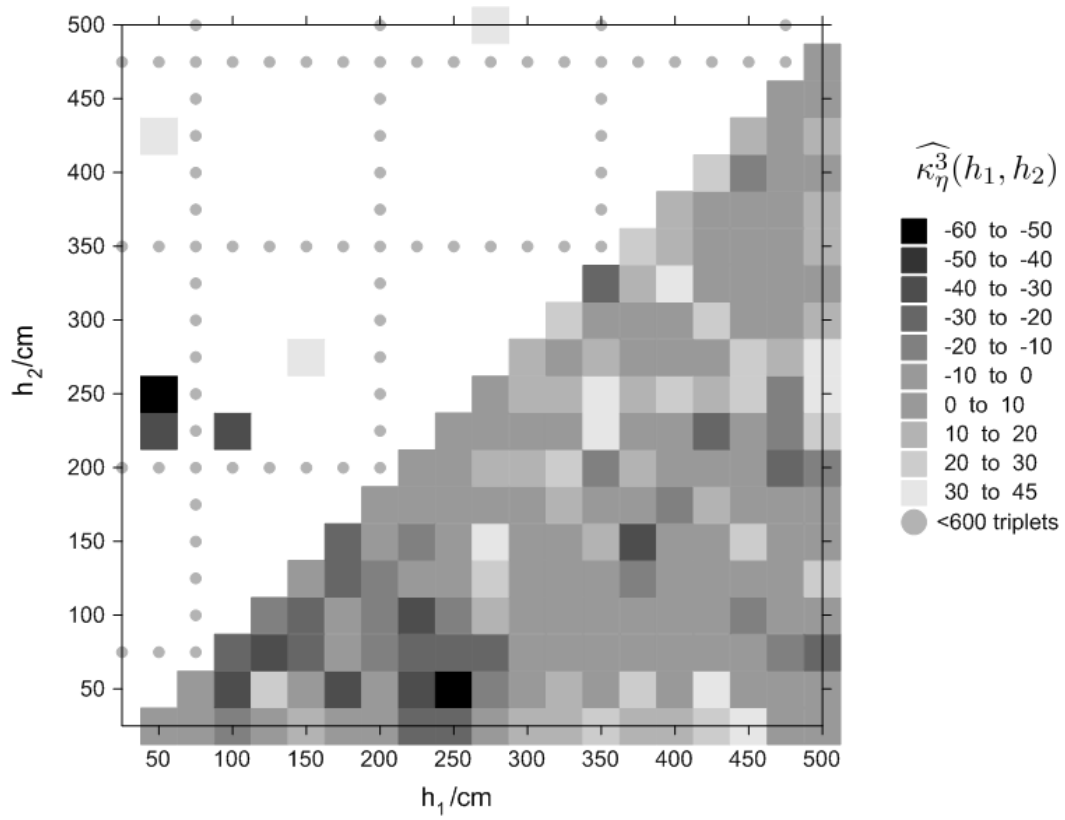


Figure 3

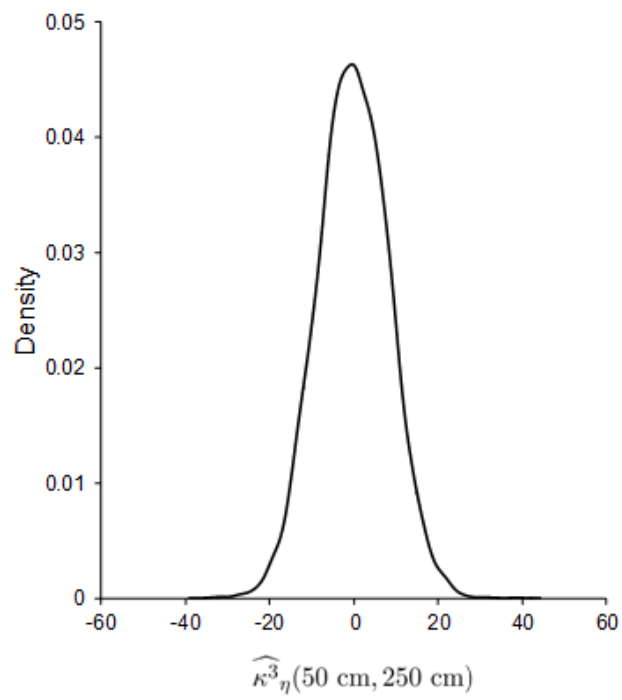
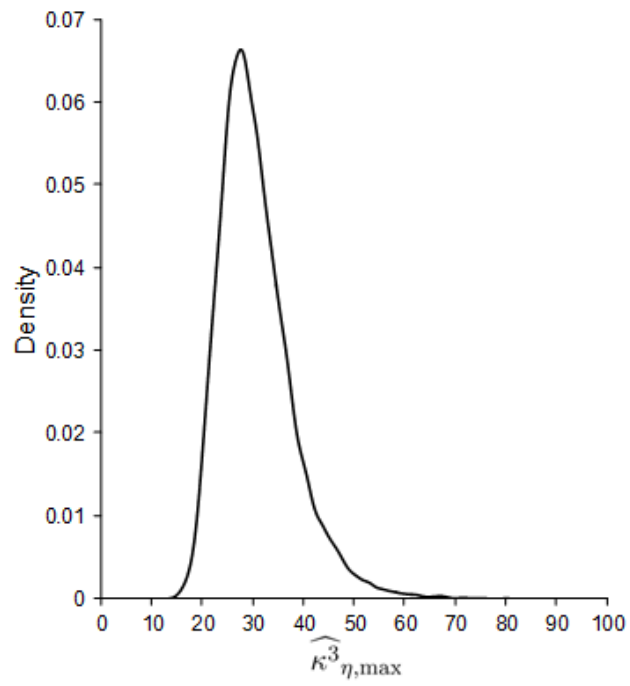


Figure 4

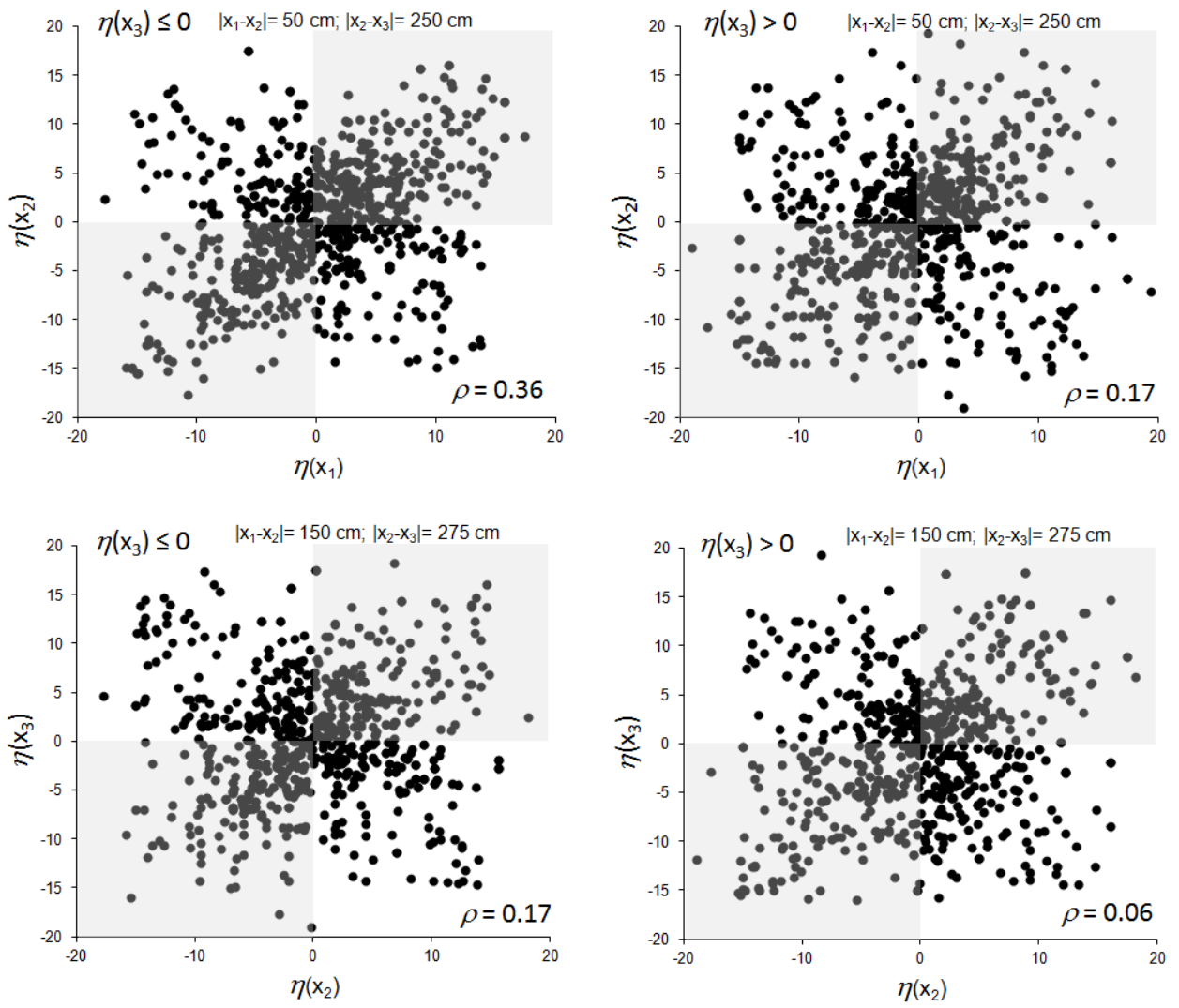


Figure 4

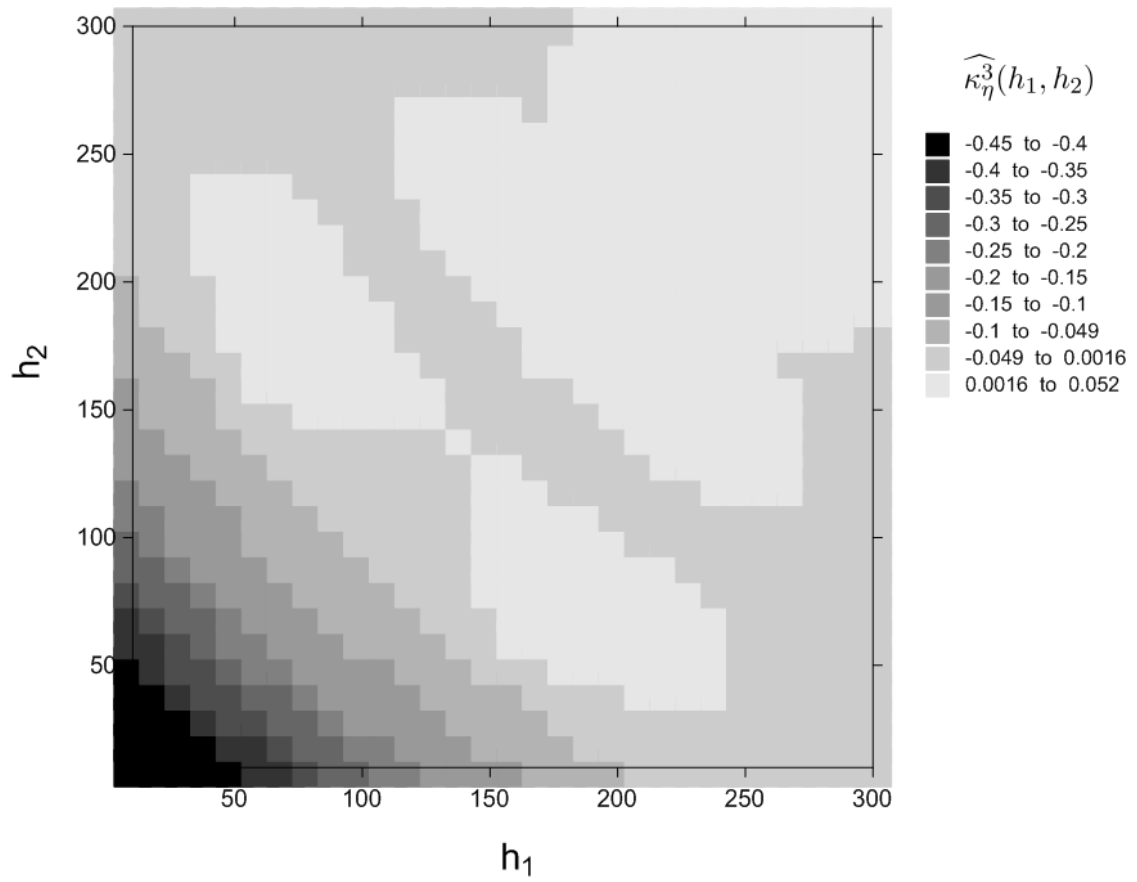


Figure 6