



Two-point declustering for weighting data pairs in experimental variogram calculations[☆]

Andrew Richmond

GeoRisk, London, UK

Received 12 May 2001; received in revised form 15 May 2001; accepted 16 May 2001

Abstract

The primary goal of calculating experimental variograms in kriging studies is to estimate the required variogram model. The destructuring effect of preferentially clustered data on the sample variogram is well-known but rarely considered in their experimental calculation and modelling due to a lack of suitable two-point declustering methods. This paper presents two methods of declustering that can be used to weight data pairs in experimental variogram calculations. Firstly, the traditional cell declustering method is extended to the two-point case, and secondly, data location clusters are identified and then explicitly used to determine declustering weights for data pairs. Experimental variograms calculated for a small, preferentially clustered and strongly positively skewed dataset show structural improvement when the two-point declustering weights are considered, resulting in significant changes to the estimated variogram model. Computer programs that implement the proposed two-point declustering and experimental variogram weighting techniques are presented. © 2002 Published by Elsevier Science Ltd.

Keywords: Variogram estimation; Kriging; Spatial structures; Destructuring effect; Geostatistics

1. Introduction

Kriging regards the continuous function $z(x)$ as a realisation of a stochastic process or continuous random function $Z(x)$, defined on the study area A , which satisfies the intrinsic hypothesis:

$$E\{Z(x) - Z(x+h)\} = 0 \quad \forall x, x+h \in A \quad (1)$$

and for any fixed displacement vector h (distance $|h|$ and direction θ),

$$\gamma(h) = \frac{1}{2} \text{var}[Z(x) - Z(x+h)] \quad \forall x, x+h \in A. \quad (2)$$

The two-point function γ of Eq. (2) is called the variogram, which is said to be isotropic if it depends only on the modulus of h , and anisotropic if it depends on both the distance h and the direction θ . In this paper, the two-point terminology is used as the two random variables, $Z(x)$ and $Z(x+h)$, in Eqs (1) and (2) relate to

the same attribute z at two different locations, x and $x+h$, rather than to two different attributes. In practice, the sample variogram $\hat{\gamma}$ is most commonly used to estimate the variogram γ , which is required for kriging.

For a sample design consisting of n data locations $\{x_\alpha, \alpha = 1, \dots, n\}$, there exist N_d data pairs (x_i, x_i+d) that are separated by distance d , where $h - \varepsilon \leq d < h + \varepsilon$. The distance d is used to provide distance classes of data location pairs. In practice, this is achieved by using distance and angle tolerances for the displacement vector h . The sample variogram is then defined as

$$\hat{\gamma}(d) = \frac{1}{2N_d} \sum_{i=1}^{N_d} [Z(x_i) - Z(x_i+d)]^2. \quad (3)$$

If $Z(x)$ is sampled at the n locations, providing specific observed values $\{z(x_\alpha), \alpha = 1, \dots, n\}$, which are then substituted for their corresponding random variables in Eq. (3), an estimate is produced. This estimate is calculated for a finite number of values of d , and the resulting experimental variogram is then modelled to provide a complete estimate of $\gamma(h)$.

[☆] Code available from server at <http://www.iamg.org/CGEditor/index.htm>

E-mail address: andrew@georisk.co.uk (A. Richmond).

An implicit goal of sampling geological phenomena is to obtain data from which the variogram γ can be estimated for kriging purposes. The requirement being to identify the type of variogram model (e.g. spherical or exponential) and the parameters of the variogram model components (e.g. nugget effect and sill). The modelling of γ is an iterative process that involves changing the values of variogram model components to find the perceived best fit of γ to $\hat{\gamma}$.

Reliability measures for the satisfactory estimation of the variogram are usually restricted to the number of data pairs for each distance d . This has resulted in several techniques for optimising the location of data points to maximise N_d (Russo, 1984; Warwick and Myers, 1987). However, as spatial data are correlated, the number of data pairs separated by a given distance d can be an unreliable measure of the precision of $\hat{\gamma}(d)$. To account for this spatial correlation, Morris (1991)

proposed the maximum equivalent uncorrelated pairs as measure of the estimation accuracy, and Zheng and Silliman (2000) derived a theoretical measure of the variance of the sample variogram $\hat{\gamma}$, both under the assumption of multivariate normality.

It is well known that the sample variogram in Eq. (3) may not be representative (accurate) if the data locations x_z are preferentially located in high- or low-valued areas. In practice, the destructuring of the sample variogram due to preferentially clustered data is: (1) simplistically accounted for by the a priori selection of one sample from each cluster with an inherent loss of information (Chiles and Delfiner, 1999) or (2) a more robust covariance measure is preferred (Isaaks and Srivastava, 1988).

Fig. 1A shows the GSLIB clustered dataset (Deutsch and Journel, 1997) overlain with a 5×5 unit grid with cell indices shown in Fig. 1B. Note that significant

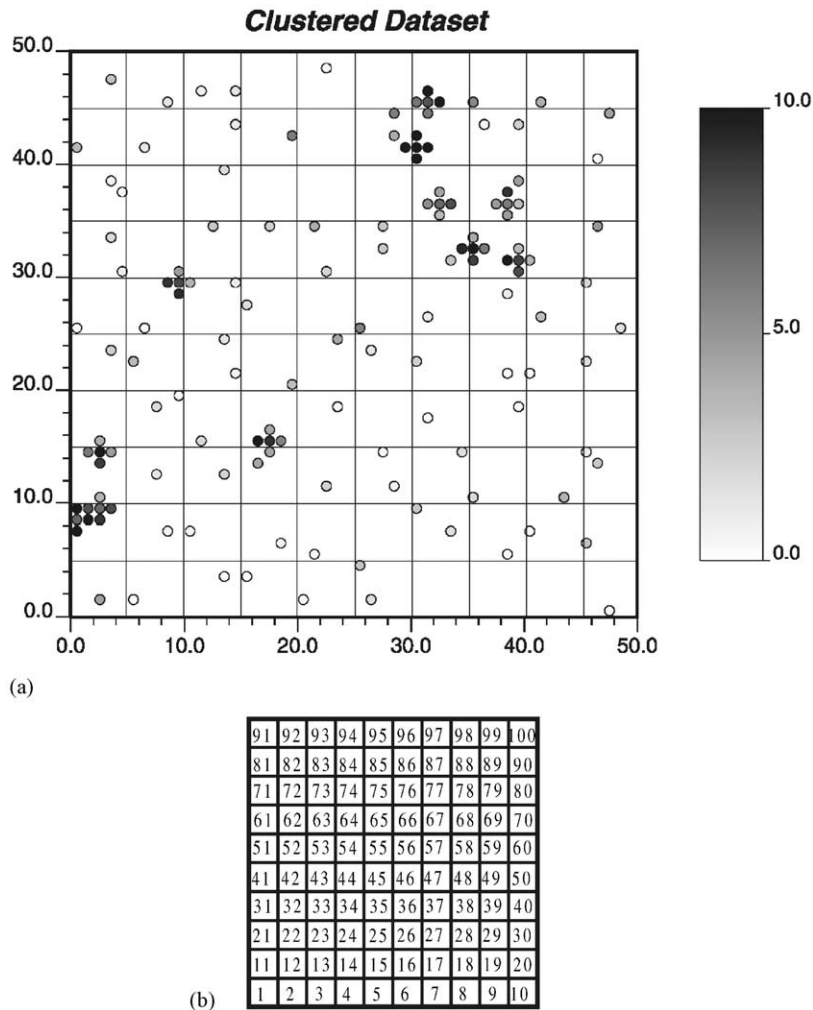


Fig. 1. (A) Clustered dataset with 5×5 unit grid; and (B) the grid indices.

clustering is present in high-grade areas, for example, the clusters present in cells 11 and 78. The clustering of data significantly destructures the experimental variograms calculated from this dataset, shown later. To understand how this distortion is introduced into experimental variograms, it is necessary to consider the location of data pairs used for experimental variogram calculations.

The data pair vectors obtained from the GSLIB dataset for various lags during experimental variogram calculations are shown in Fig. 2. Note that in all three plots, the vector densities vary throughout the study area, as expected from clustered data. For example, in the first lag there are significantly more data pair vectors in the top right and bottom left higher valued parts of the study area, shown in the top diagram of Fig. 2. However, from these diagrams it is difficult to quantify the influence of various data locations and their values on the experimental variogram based on the degree of local clustering.

Fig. 3 shows the number of times a data location was paired during experimental variogram calculations for the three lags shown in Fig. 2. In this diagram, note that the data pairing magnitudes are not related to the degree of local clustering but on the data clustering at distance d . For example, in the cluster centred on $x = 32.5$ and $y = 36.5$ for the 10th lag, three data locations have been paired once, one data location is paired twice, and the remaining data location is not paired. Moreover, the data location shown at the origin of the lag vector ($x = 38.5$ and $y = 5.5$) in the three diagrams is paired 0, 2, and 5 times for lags one, five, and ten, respectively. Thus, weighting the squared difference of the data pair values in experimental variogram calculations by the product of the univariate data location declustering weights for each distance d (Rivoirard, 2000) is not appropriate.

To account for the influence of preferentially clustered data on experimental variograms, and, consequently on the estimated variogram, it is first proposed to extend the univariate cell declustering method (Journel, 1983) to the two-point case, and secondly to explicitly use data location clusters to determine two-point declustering weights. The two-point declustering results from both methods are then used to weight the squared difference of the data pair values during the calculation of experimental variograms, removing the distortion associated with data clustering. Effects on the estimated variogram when using the proposed methods are investigated.

2. Two-point declustering

Methods for obtaining declustered first-order moments, such as the mean and variance of $Z(x)$, have been discussed by several authors. Traditional declustering

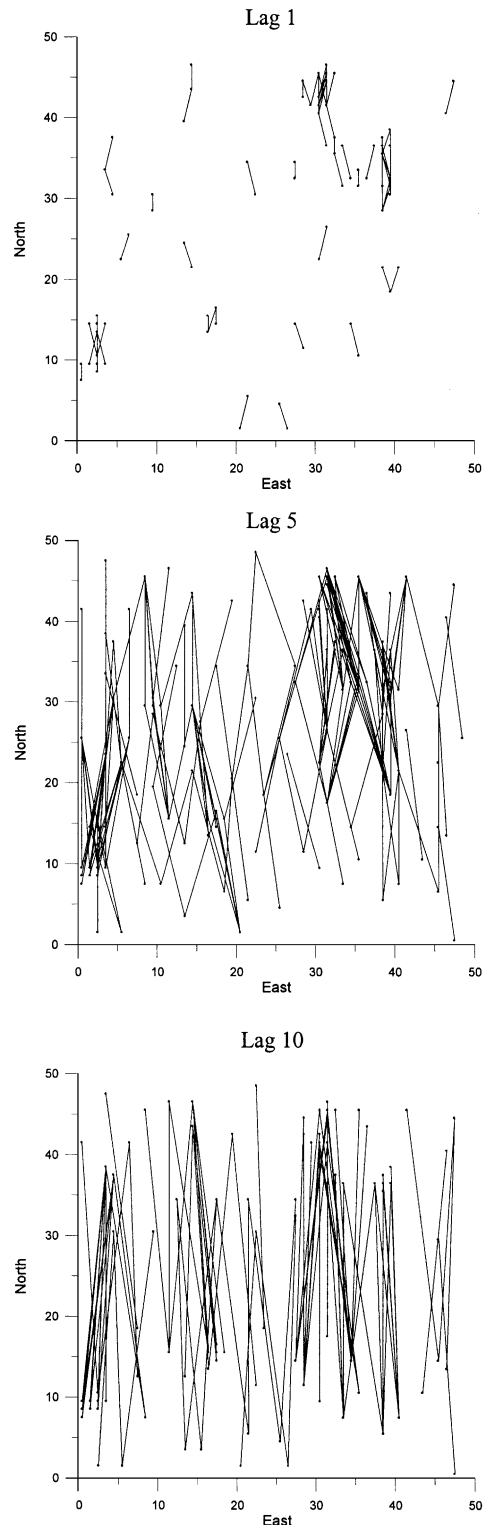


Fig. 2. Data pairs for various lags. (vector azimuth = 000° ; vector tolerance = 20° ; lag = 3 units; bandwidth = 5 units).

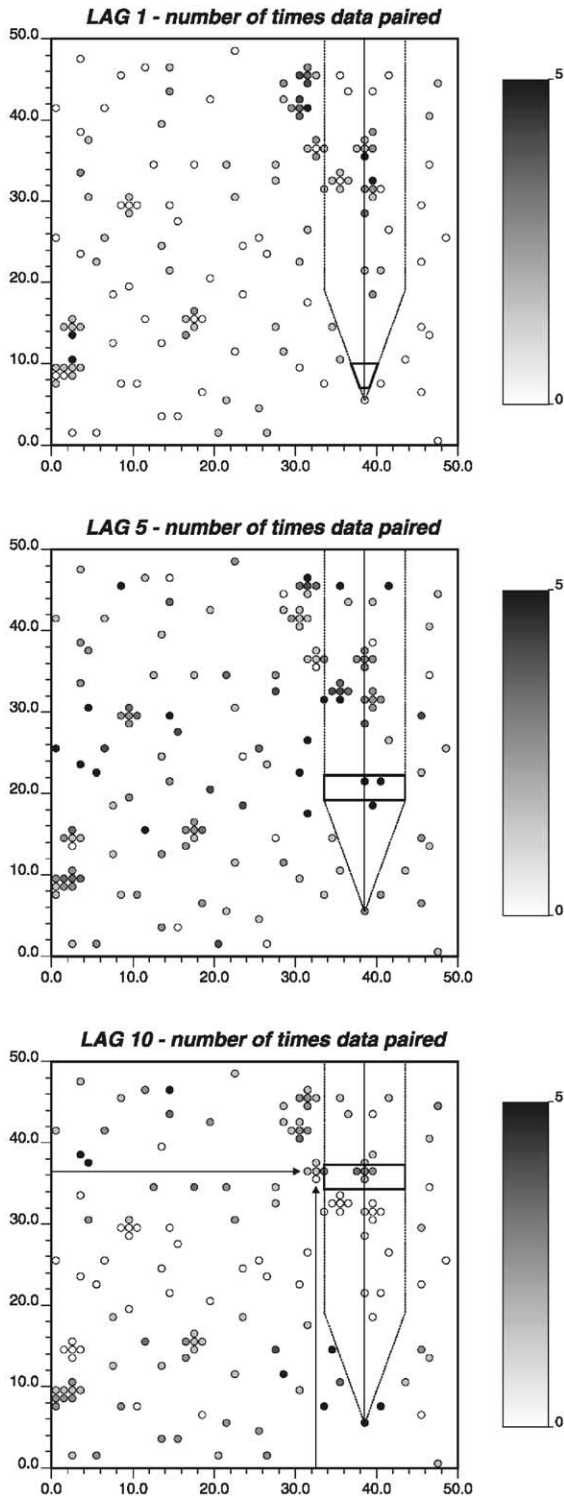


Fig. 3. Number of times data locations are paired for various lags. (vector azimuth = 000°; vector tolerance = 20°; lag = 3 units; bandwidth = 5 units).

techniques such as polygons of influence David (1977) and cell declustering (Journel, 1983) account for spatial clustering under the implicit assumption of spatial independence. Declustering using a global kriging approach (Isaaks and Srivastava, 1989) or redundancy co-efficients obtained from spatial correlation matrices (Bourgault, 1997) explicitly account for spatial correlations with an average covariance function. However, it is well known that the importance of the clustering effect on the estimate of the cumulative distribution function increases when the dependency between spatial locations increases (Cressie, 1993). Bogaert (1999) proposed a least-squares approach to weight data values to account explicitly for the presence of spatial dependence that varied with the magnitude of data values.

Techniques for declustering two-point statistics, such as the experimental variogram, have rarely been developed in the literature. Rivoirard (2000) suggested weighting data pairs in experimental variogram calculations using various weighting schemes, including the product of the univariate sample declustering weights. This section presents two methods for two-point declustering.

2.1. Two-point cell declustering

Weighted cell declustering for univariate statistics involves overlaying a regular grid on the area considered, and for each occupied cell, assign a weight that is inversely proportional to the number of data locations present within that cell, to each data location within that cell (Journel, 1983; Deutsch, 1989). For the two-point case, the procedure for determining the declustering weights is

1. overlay area A with a grid of C regular cells, c_i , $i = 1, \dots, C$;
2. for each pair of data locations $\{x_{2i} \in c_i, x_{2j} \in c_j\}$ separated by distance d count the number of vectors v_{ij} that originate in cell c_i , and terminate in cell c_j ; and
3. for distance d , the weight for a pair of data locations $\{x_{2i} \in c_i, x_{2j} \in c_j\}$ is $1/v_{ij} = w_{2ij}$.

The GSLIB clustered dataset was used to calculate two-point declustering weights using the proposed methodology and the 5×5 unit grid shown in Fig. 1. The cell size was chosen as it minimised the declustered mean of the data (Deutsch and Journel, 1997). The plots in Fig. 4 show the number of times v_{ij} , the data locations fall in the various cell pairs. In these diagrams, positive values of v_{ij} are mostly 1, but occasionally exceed 10. For example, in the 10th lag, $v_{ij} = 14$ for the cell pair (11, 71). Therefore, in this example, data pairs originating and terminating in essentially the same area, due to data clustering, can influence two-point statistics by

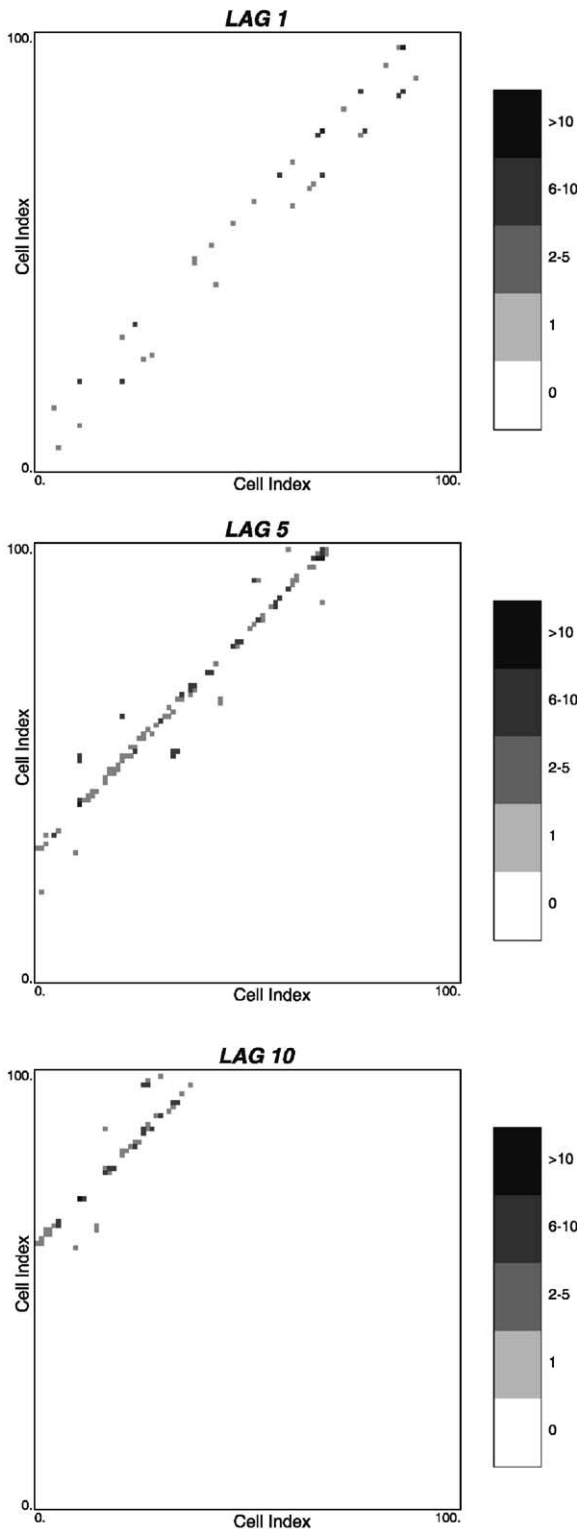


Fig. 4. Number of vectors v_{ij} that belong to various cell pairs for same lags shown in Fig. 2.

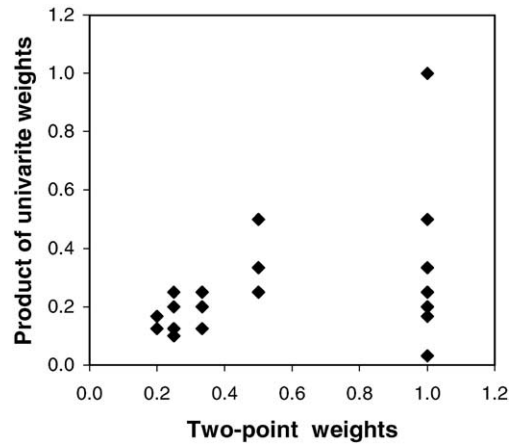


Fig. 5. Scatter plots of two-point cell declustering weights $1/v_{ij}$ versus product of univariate cell declustering weights for 5th lag. (Note: there is significant overplotting of values).

more than 10 times than those data pairs located in sparsely sampled areas. When using two-point cell declustering, the idea is to weight the data pairs inversely proportional to the number of times the data locations belong to the same cells.

The two-point cell declustering weights are sensitive to the distance d as v_{ij} varies with d , shown in Fig. 4, and the cell grid configuration. By varying the grid origin and cell sizes, individual clusters of data locations either will fall entirely within a single cell (ideal), or will be split over multiple cells (not ideal).

Note that, in general, the two-point cell declustering weights are not the product of the univariate cell declustering weights for the corresponding set of cells. This is shown in Fig. 5, a scatter plot of the two-point cell declustering weights versus the product of the corresponding univariate cell declustering weights. However, for any lag in which all the data locations in a cell c_i were paired with all data locations in a cell c_j , then the two-point cell declustering weights will equal the product of the univariate cell declustering weights for c_i and c_j .

For a separation vector $h = 0$, a data location can only be paired with itself, i.e. $x_x = x_x$ and $c_i = c_j$. Consequently, the weights applied to these data pairs are simply inversely proportional to the number of data locations present within that cell. Thus, univariate cell declustering (Journel, 1983) is simply a special case of the proposed two-point cell declustering method.

2.2. Two-point declustering with data location clusters

An alternative method of two-point declustering that is not sensitive to the cell configuration explicitly considers data location clusters. The idea is

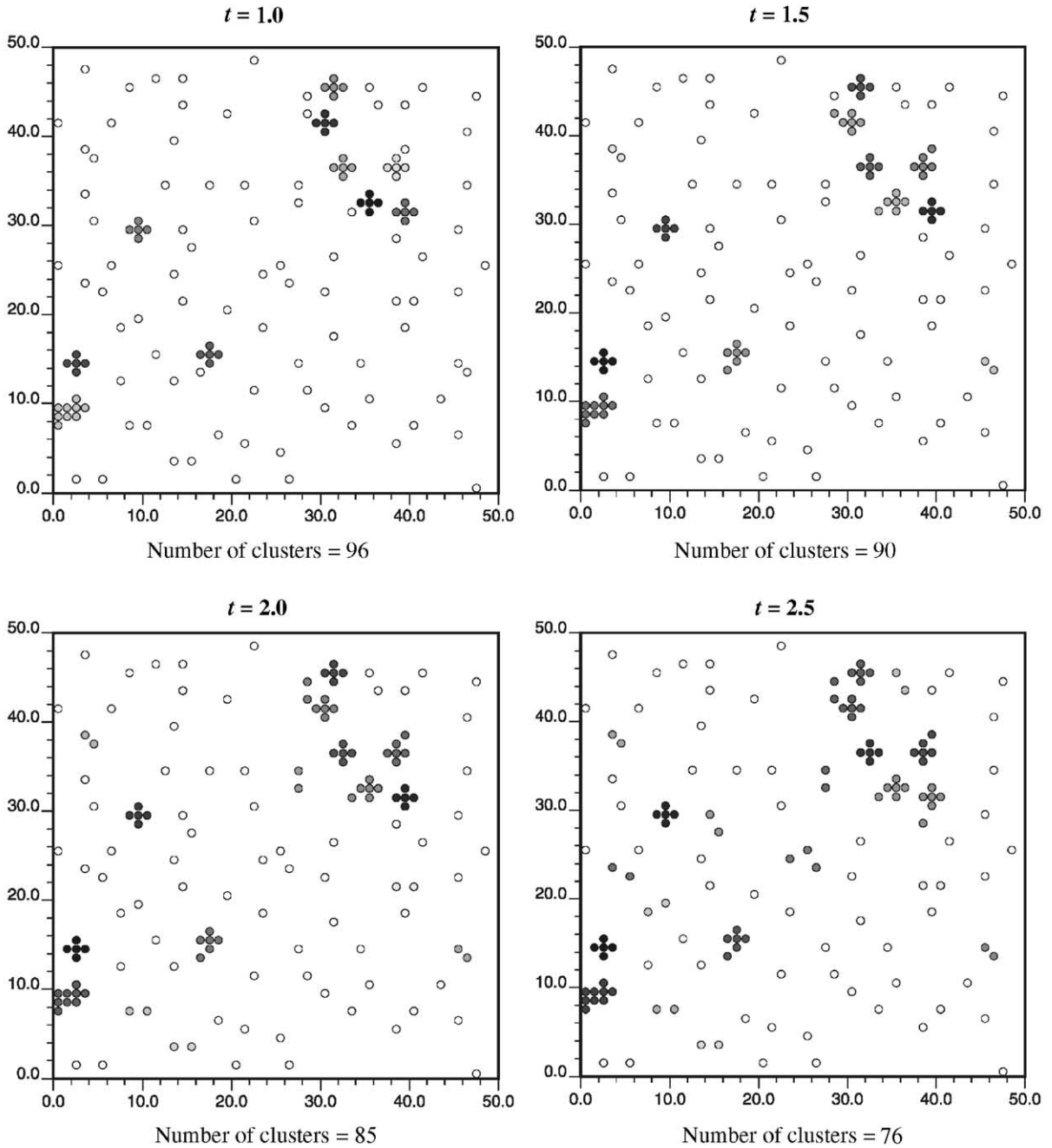


Fig. 6. Clusters for various distances t . (Note: only clusters with two or more data locations are shaded, and shading intensity is unique to cluster).

1. identify the number of data location clusters G within A , g_i , $i = 1, \dots, G$;
 2. for each pair of data locations $\{x_z \in g_i, x_{z'} \in g_{i'}\}$ separated by distance d count the number of vectors $v_{ii'}$ that originate in cluster g_i , and terminate in cluster $g_{i'}$; and
 3. for distance d , the weight for a pair of data locations $\{x_z \in g_i, x_{z'} \in g_{i'}\}$ is $1/v_{ii'} = w_{zz'}$.
- There is considerable literature on methods of determining clusters for spatially correlated data (e.g. Allard and Guillot, 2000). In this study, if two data

locations are less than a specified distance t apart they are considered to belong to the same cluster, i.e. if $|x_\alpha - x_{\alpha'}| \leq t$ then $x_\alpha, x_{\alpha'} \in g_j$. Fig. 6 shows data location clusters identified for various t values. Note that the gross dimension of the cluster is not equal to t , but by the number and relative positions of the data locations that jointly satisfy the constraint $|x_\alpha - x_{\alpha'}| \leq t$. In Fig. 6, as t increases the number of clusters decreases. This is due to originally ungrouped data being considered clustered, or several shorter-scale clusters merging. The former (dashed ellipse) and latter (dashed square) scenarios are identified in the top and bottom pairs of diagrams in Fig. 6, respectively. Using this clustering algorithm, a practical distance t should minimise the merging of shorter-scale clusters.

The GSLIB clustered dataset was used to calculate two-point declustering weights using the proposed methodology and $t = 1.5$ units. The diagrams in Fig. 7 show the number of times v_{ij} the data locations fall in the various cluster pairs. In Fig. 7, for the first lag, there are several data pairs sourced entirely from a single cluster, shown as values greater than or equal to one on the 45° diagonal in the top diagram. As with the corresponding two-point cell declustering plots (Fig. 4), the magnitude of v_{ij} varies considerably from 1 to >10 . When employing this two-point declustering technique, the idea is to weight the data pairs inversely proportional to the number of times the data locations belong to the same clusters.

3. Weighting data pairs in experimental variogram calculations

The two-point declustering weights from the previous section can be used to weight the sample variogram in Eq. (3), i.e.

$$\hat{\gamma}(d) = \frac{1}{2W_d} \sum_{\alpha=1}^{N_d} [z(x_\alpha) - z(x_{\alpha'})]^2 w_{\alpha\alpha'}, \quad (4)$$

where $W_d = \sum_{\alpha=1}^{N_d} w_{\alpha\alpha'}$, the sum of two-point weights for distance d ; and $|x_\alpha - x_{\alpha'}| \in d$.

Computer programs *bicell* and *biclus* were used to calculate experimental variograms for the GSLIB dataset using the proposed methodology. For these calculations, the data values were initially transformed to their natural logarithms prior to using Eqs. (3) and (4). Log experimental variogram values were calculated for comparison purposes using the reference dataset from which the GSLIB clustered dataset were drawn (Deutsch and Journel, 1997).

Fig. 8 shows the unweighted and weighted log experimental variogram values in the north–south and east–west directions. In Fig. 8, the unweighted log experimental variogram values for the north–south

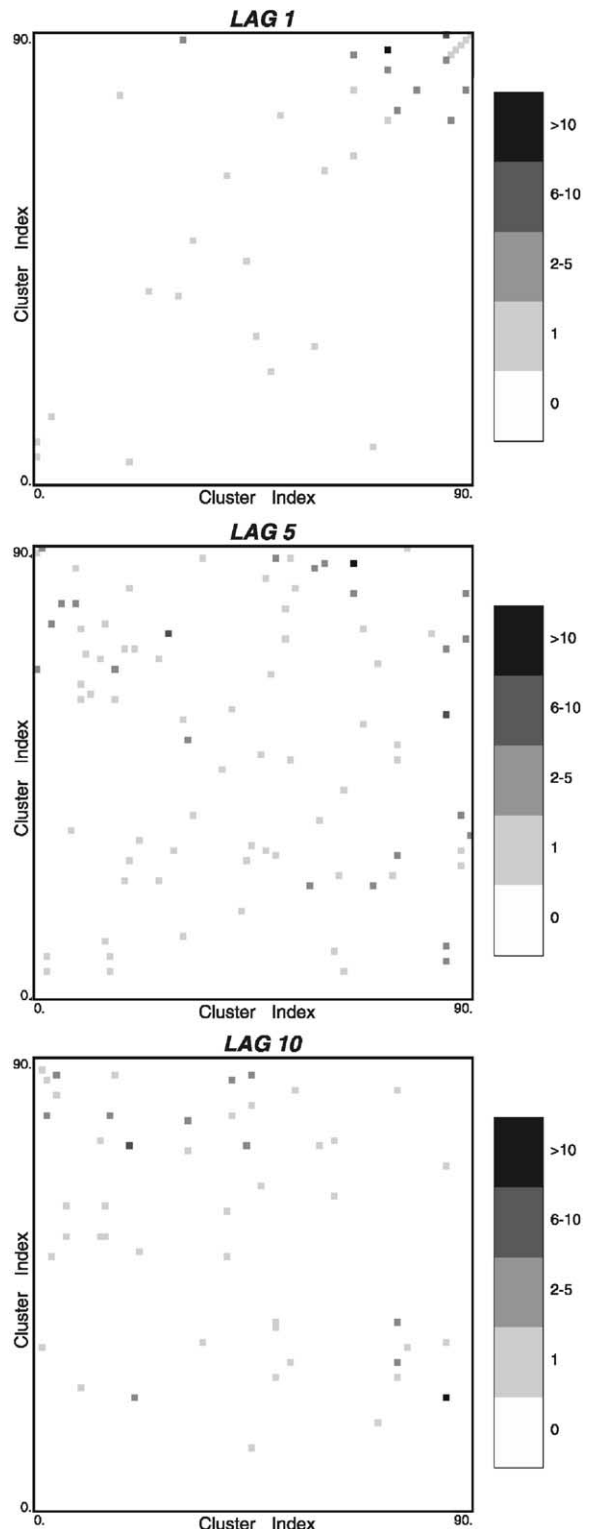


Fig. 7. Number of vectors v_{ij} that belong to cluster pairs for various lag vectors. (Note: the cluster index is assigned randomly; vector azimuth = 000°; vector tolerance = 20°; lag = 3 units; bandwidth = 5 units; cluster distance = 1.5 units).

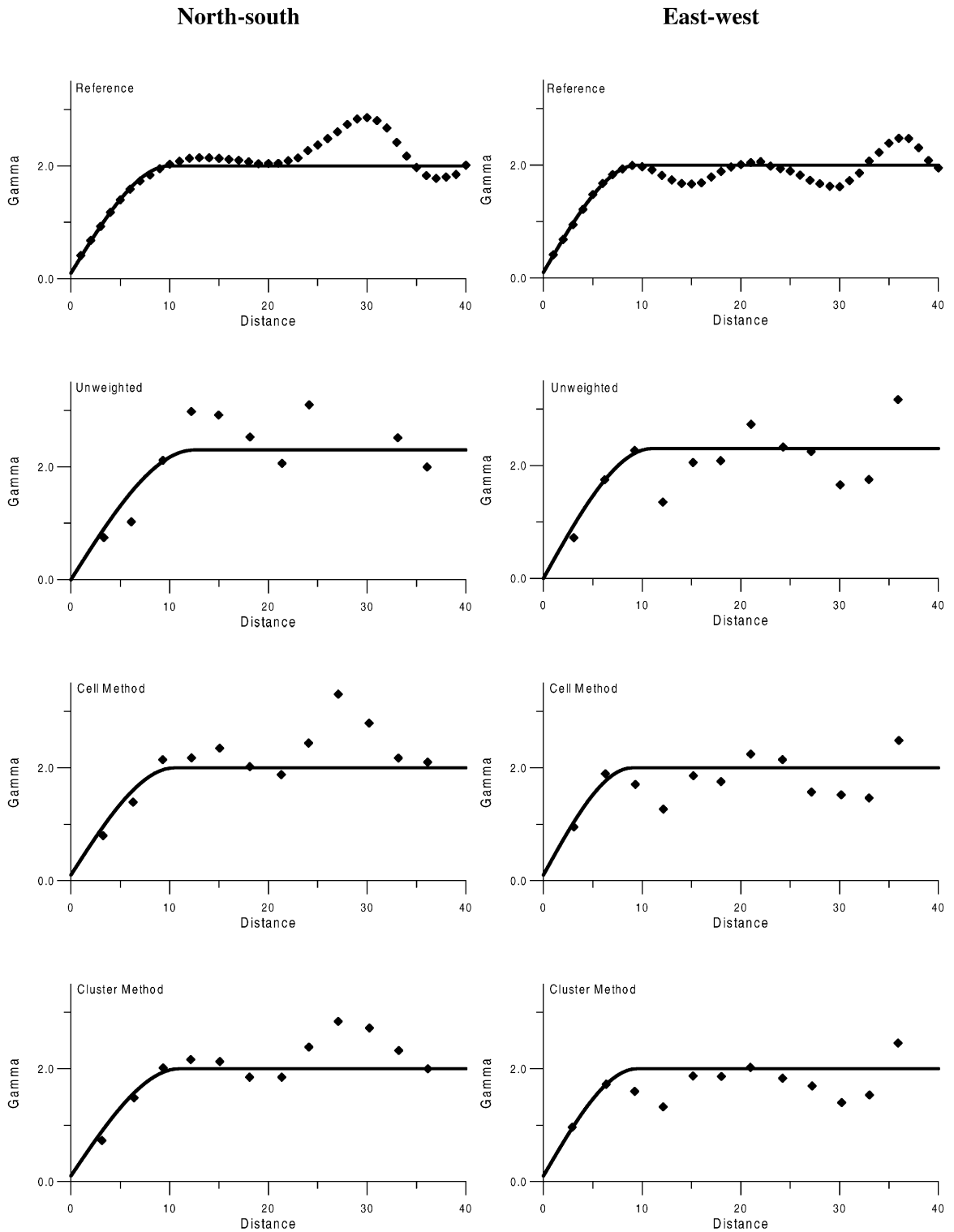


Fig. 8. Log experimental variograms and the modelled variogram. (Vector tolerance = 20°; lag = 3 units; bandwidth = 5 units; cell size = 5 × 5 units; cluster distance = 1.5 units.)

Table 1
Estimated log variogram model paramters (sph. = spherical)

Method	North–south				East–west			
	Model type	C_0	C_1	a_1	Model type	C_0	C_1	a_1
Reference	Sph.	0.1	1.9	10.0	Sph.	0.1	1.9	9.5
Unweighted	Sph.	0.0	2.3	12.5	Sph.	0.0	2.3	11.0
Weighted (cell)	Sph.	0.1	1.9	10.5	Sph.	0.1	1.9	9.0
Weighted (cluster)	Sph.	0.1	1.9	11.0	Sph.	0.1	1.9	9.5

```

Parameters for BICELL
*****

START OF PARAMETERS:
cluster.dat          \ data file
1  2  0             \  x, y, z columns
3                   \  variable column
2                   \  vtyp (1=semi; 2=log-variogram)
3.0                 \  lag - unit separation distance
0.0                 \  ltol- lag tolerance
12                  \  nlags - number of lags
      0.0  20.0  5.0 \  azi,atol,aband
      0.0  20.0  5.0 \  dip,dtol,dband
0.0  0.0  0.0       \  x,y,z cell origins
5.0  5.0  1.0       \  x,y,z cell sizes
10  10  1           \  x,y,z cell numbers
bicell.out          \  output file
1                   \  debug level (0,1)
bicell.dbg          \  debug file

```

Fig. 9. Example parameter file for *bicell* program.

direction fluctuate markedly about the reference log experimental variogram values, and the spatial structure is unclear. Isaaks and Srivastava (1988) suggest centring, as for covariances, and standardising, as for correlograms, to stabilise and provide interpretable spatial structures. However, uncertainty of the quantities used for centering and normalisation may also introduce a bias into these measures (Rivoirard, 2000). Experimental log-variogram values employing both of the weighting methods do not differ significantly from the reference log experimental variogram values, and clearly provide more interpretable structures. Similar relationships can be noted for the east–west log experimental variograms in Fig. 8.

All experimental log variograms were then modelled, shown as solid lines in Fig. 8. The estimated variogram model parameters are included in Table 1. This table indicates that the variogram estimated from the unweighted experimental variogram has a lower nugget effect, higher spatial variance, and longer range than the variogram estimated from the reference experimental

variogram. For both proposed data pair weighting techniques the estimated variogram model parameters closely match the parameters of the variogram model estimated from the reference data.

4. Computer programs

FORTTRAN programs *bicell* and *biclus* implement both two-point declustering algorithms to weight data pairs in experimental variogram calculations. The programs are modelled after and use the GSLIB algorithm to pair data for the various lag distances d (Deutsch and Journel, 1997).

The input parameters for two-point cell declustering (*bicell*) are shown in Fig. 9 and documented as follows:

- *dataft*: input data file.
- *ixl, iyl, izl*: columns for the x, y, z co-ordinates.
- *ivr*: variable column.

```

Parameters for BICLUS
*****

START OF PARAMETERS:
cluster.dat          \ data file
1  2  0             \  x, y, z columns
3                  \  variable column
2                  \  vtyp (1=semi; 2=log-variogram)
3.0                \  lag - unit separation distance
0.0                \  ltol- lag tolerance
12                 \  nlags - number of lags
      0.0   20.0   5.0 \  azi, atol, aband
      0.0   20.0   5.0 \  dip, dtol, dband
1.5                \  cluster distance t
biclus.out         \  output file
1                 \  debug level (0,1)
biclus.dbg        \  debug file

```

Fig. 10. Example parameter file for *biclus* program.

- *vtyp*: variogram type (1=variogram; 2=log variogram).
- *lag*: unit lag separation distance.
- *ltol*: lag tolerance.
- *nlags*: number of lags.
- *azi*, *atol*, *aband*: azimuth, half window azimuth tolerance, and azimuth bandwidth.
- *dip*, *dtol*, *dband*: dip, half window dip tolerance, and dip bandwidth.
- *xo*, *yo*, *zo*: *x*, *y*, *z* cell grid origin.
- *xsiz*, *ysiz*, *zsiz*: *x*, *y*, *z* cell sizes.
- *xn*, *yn*, *zn*: *x*, *y*, *z* number of cells.
- *outft*: output file.
- *idbg*: debug level.
- *dbgft*: debug file.

The input parameters for two-point declustering using clusters (*biclus*) are shown in Fig. 10 and documented as follows:

- *dataft*: input data file.
- *ixl*, *iyl*, *izl*: columns for the *x*, *y*, *z* co-ordinates.
- *ivr*: variable column.
- *vtyp*: variogram type (1=variogram; 2=log variogram).
- *lag*: unit lag separation distance.
- *ltol*: lag tolerance.
- *nlags*: number of lags.
- *azi*, *atol*, *aband*: azimuth, half window azimuth tolerance, and azimuth bandwidth.
- *dip*, *dtol*, *dband*: dip, half window dip tolerance, and dip bandwidth.
- *htol*: cluster distance *t*.
- *outft*: output file.
- *idbg*: debug level.
- *dbgft*: debug file.

The output files from *bicell* and *biclus* contain *nlags* lines for both the unweighted and weighted variograms, including the lag number, the average distance between data pairs for the lag; the variogram value; the lag mean; the lag variance; and the number of data pairs. Information output to the debug files can be used to generate some of the diagrams shown in this paper.

The computer programs are memory intensive since an array of size $nlags \times C \times C$ or $nlags \times G \times G$ is required. For large study areas, it may be necessary to utilise an alternative algorithm that reduces the computer programming requirements to two-dimensional arrays. This could be achieved by considering the vector centres of data pairs as the clustered attribute and using a univariate declustering method to determine the data pair weights.

5. Conclusions

This study provides evidence of the destructuring effect of preferentially clustered data on the experimental variogram when the data are strongly positively skewed, the number of samples is small, and the clustering is significant. However, unlike univariate statistics in which bias results from local clustering, the influence of data values on the variogram is related to clustering at distance *d* from its location. Consequently, two-point declustering weights vary with *d*.

Two-point declustering to weight data pairs in experimental variogram calculations was proposed and involved using either a cell declustering method or by considering explicitly the data location clusters. The former is sensitive to the grid configuration, and the latter to the manner in which data locations are grouped into clusters. The two-point weights for the cell declustering method are not the products of the

univariate cell declustering weights. However, univariate cell declustering weights can be obtained by considering a null vector in the two-point approach. Other methods of determining clusters should be considered when employing the second declustering method (e.g. Allard and Guillot, 2000).

In the example, destructuring of unweighted log experimental variograms due to clustering was significant and resulted in poor estimation of the log variogram model. Log experimental variograms showed significant structural improvement when the proposed two-point declustering weights were considered. In addition, the estimated log variogram model using this approach closely matched the log variogram model estimated from the reference data.

When the number of data is large and preferential clustering is present, the proposed approach to weighting data pairs in experimental variogram calculations will provide more accurate experimental variogram values. However, it is not clear whether these more accurate experimental variograms will result in significant improvements to the estimated variogram model.

References

- Allard, D., Guillot, G., 2000. Clustering geostatistical data. In: Kleingeld, W.J., Krige, D.C. (Eds.), Proceedings of the 6th International Geostatistics Congress, Cape Town, South Africa, 15pp.
- Bogaert, P., 1999. On the optimal estimation of the cumulative distribution function in the presence of spatial dependence. *Mathematical Geology* 31 (2), 213–239.
- Bourgault, G., 1997. Spatial declustering weights. *Mathematical Geology* 29 (2), 277–290.
- Chiles, J.P., Delfiner, P., 1999. *Geostatistics: Modeling Spatial Uncertainty*. Wiley, New York, 672pp.
- Cressie, N., 1993. *Statistics for Spatial Data*. Wiley, New York, 900pp.
- David, M., 1977. *Geostatistical Ore Reserve Estimation*. Elsevier, Amsterdam, 364pp.
- Deutsch, C.V., 1989. Declus: a fortran 77 program to determine optimal spatial declustering weights. *Computers & Geosciences* 15 (3), 325–332.
- Deutsch, C.V., Journel, A.G., 1997. *GSLIB: Geostatistical Software Library and User's Guide*, 2nd edition. Oxford University Press, New York, 369pp.
- Isaaks, E.H., Srivastava, R.M., 1988. Spatial continuity measures for probabilistic and deterministic geostatistics. *Mathematical Geology* 20 (4), 313–341.
- Isaaks, E.H., Srivastava, R.M., 1989. *An introduction to applied geostatistics*. Oxford University Press, New York, 561pp.
- Journel, A.G., 1983. Non-parametric estimation of spatial distributions. *Mathematical Geology* 15 (3), 445–468.
- Morris, M.D., 1991. On counting the number of data pairs for semivariogram estimation. *Mathematical Geology* 23 (7), 929–943.
- Rivoirard, J., 2000. Weighted variograms. In: Kleingeld, W.J., Krige, D.C. (Eds.), Proceedings of the 6th International Geostatistics Congress, Cape Town, South Africa, 11pp.
- Russo, D., 1984. Design of an optional sampling network for estimating the variogram. *Soil Science Society of America Journal* 48, 708–716.
- Warwick, A.W., Myers, D.E., 1987. Optimization of sampling locations for variogram calculations. *Water Resources Research* 23 (3), 496–500.
- Zheng, L., Silliman, S.E., 2000. Estimating the theoretical semivariogram from finite numbers of measurements. *Water Resources Research* 36 (1), 361–366.