

Spatial analysis of compositional data: A historical review

doi:10.1016/j.gexplo.2015.12.010 ☆

Vera Pawlowsky-Glahn

Universitat de Girona

Juan José Egozcue

Universitat Politècnica de Catalunya

Abstract

Like the statistical analysis of compositional data in general, spatial analysis of compositional data requires specific tools. An historical overview of their development is presented in three steps: (a) the recognition of the problem, known as spurious spatial covariance, (b) first attempts to use the logratio approach, and (c) the application of the principle of working in coordinates using isometric logratio representations. Also mentioned are the use of matrix-valued variation-variograms as a tool to model crossvariograms, and the simplicial approach to indicator kriging, that solves inconsistencies in the standard approach to indicator kriging.

Keywords: compositional data analysis, Geostatistics, simplex, variation-variogram, simplicial indicator kriging

2010 MSC: 62-07, 86A32

*Vera Pawlowsky-Glahn

Email address: vera.pawlowsky@udg.edu (Juan José Egozcue)

¹Dept. Computer Science, Applied Mathematics, and Statistics, Girona, Spain

²Dep. Applied Mathematics III, Barcelona, Spain

1. Introduction

According to Chilès and Delfiner (2012), the term *Geostatistics* was introduced by Matheron (1962) to designate his own methodology for ore reserve estimation. Since then, Geostatistics expanded amazingly, as the methodology finds application in many fields, not only in geo- and environmental sciences. Independently, in the 1980's, J. Aitchison started developing *compositional data analysis* (CoDa) (Aitchison and Shen, 1980; Aitchison, 1982, 1986) introducing what nowadays is known as *the log-ratio approach*. Although most type of data to which Geostatistics is applied are compositional, like ore grade, chemical or mineralogical composition of rocks, contaminants in air or water, it was not recognized until 1984 that spurious spatial correlation might be at work (Pawlowsky, 1984). We summarise in what follows the steps that have been undertaken since then to solve the problems derived from the compositional character of some spatially dependent data. We limit our contribution to the historical development, omitting most formal derivations which can be found in the references cited.

2. Spurious spatial covariance

The problem of spurious spatial covariance of regionalized compositions, or *r-compositions* for short, was first stated in Pawlowsky (1984). The results are illustrative, and are therefore briefly exposed.

According to our present understanding, a random vector, \mathbf{Z} , with D strictly positive components representing parts of a whole, is a composition if it carries only relative information (Pawlowsky-Glahn et al., 2015c). Note that the term *relative information* is equivalent to *information lies in the ratios between components*, not in the absolute values. The same definition holds for a spatially distributed random vector, $\mathbf{Z}(x)$, at any point x of a spatial domain \mathcal{R} .

27 In 1984, r-compositions were still understood as random vectors subject to
 28 a constant sum constraint, or *closed r-compositions*. We know now that compo-
 29 sitions in general, and r-compositions in particular, are equivalence classes, and
 30 that a closed composition is just a representation. This means, that the results
 31 obtained under this assumption hold for any representation of the equivalence
 32 classes.

33 For the understanding of spurious spatial covariance or correlation, it is
 34 mathematically easier to work with a closed representation. Therefore, in what
 35 follows, we work with a *closed r-composition*, i.e. with a spatially distributed
 36 random vector, $\mathbf{Z}(x)$, with D strictly positive parts or components, that is
 37 subject to a constant sum constraint for all $x \in \mathcal{R}$,

$$\sum_{i=1}^D Z_i(x) = \kappa, \quad (1)$$

38 with κ a given positive constant depending on the units of the random vector.
 39 The constant κ is usually 1 (parts per unit), 100 (percentages), or 10^6 (parts
 40 per million).

41 Following Matheron (1965), Geostatistics can be used with regionalized vari-
 42 ables satisfying stationarity conditions. Second order stationarity requires re-
 43 gionalized variables to have a constant mean and the autocovariance only de-
 44 pending on the lag between pairs of variables $\mathbf{Z}(x_i)$ and $\mathbf{Z}(x_j)$; a less stringent
 45 condition is the *intrinsic hypothesis*, which assumes that the first order differ-
 46 ences are second order stationary. Under these kind of assumptions, Geostatis-
 47 tics builds on modelling the mean and the spatial autocovariance, or related
 48 parameters, like the variogram. The following development handles the com-
 49 ponents of the closed r-composition $\mathbf{Z}(x) = (Z_1(x), Z_2(x), \dots, Z_D(x))$ at two
 50 spatial locations, say x and $x + h$ in \mathcal{R} , where h denotes the lag between them.

51 From Eq. (1), for any lag h it holds

$$\sum_{i=1}^D (Z_i(x) - Z_i(x+h)) = \sum_{i=1}^D Z_i(x) - \sum_{i=1}^D Z_i(x+h) = \kappa - \kappa = 0. \quad (2)$$

Hence, multiplying both sides of Eq. (2) by $(Z_j(x) - Z_j(x+h))$,

$$\sum_{i=1}^D (Z_i(x) - Z_i(x+h)) (Z_j(x) - Z_j(x+h)) = 0,$$

52 for any $j = 1, 2, \dots, D$. Taking expectations,

$$\sum_{i=1}^D \text{cov} [(Z_i(x) - Z_i(x+h)), (Z_j(x) - Z_j(x+h))] = 0. \quad (3)$$

53 Given that a variance is always positive, Eq. (3) can be rewritten for any
54 $j = 1, 2, \dots, D$, as

$$\begin{aligned} & \text{var} [(Z_j(x) - Z_j(x+h)) (Z_j(x) - Z_j(x+h))] \\ &= - \sum_{i \neq j} \text{cov} [(Z_i(x) - Z_i(x+h)), (Z_j(x) - Z_j(x+h))] . \end{aligned} \quad (4)$$

55 Note that Eq. (4) depends only on the fact that $\mathbf{Z}(x)$ is the closed representation
56 of an r-composition, and not on the type of spatial dependence of its components.
57 Equation (4) implies that non-stochastic factors determine the value of cross-
58 covariances. They cannot be all null simultaneously, as the variance is, by
59 definition, always positive. Also, if the closed r-composition was generated by
60 closure of independent random variables, a dependence will appear, which is
61 spurious, as it is not generated by the phenomenon itself (Pawlowsky, 1984).
62 This result is well known for compositional data in general as the *closure problem*
63 (Chayes, 1960). It has many implications in standard multivariate analysis
64 which can be directly extended to r-compositions.

65 For a closed intrinsic r-composition $\mathbf{Z}(x)$, Eq. (4) can be written in terms of
 66 variograms, $\gamma_j(h)$, and crossvariograms, $\gamma_{ij}(h)$,

$$\gamma_j(h) = - \sum_{i \neq j} \gamma_{ij}(h), \quad j = 1, 2, \dots, D. \quad (5)$$

67 for any lag h . As stated in Pawlowsky (1984), the obvious conclusion is the
 68 need of non-zero cross-variograms for r-compositions, some of which have to
 69 be negative—as the variogram is, by definition, positive. It is clear that the
 70 only case in which cross-variograms could be all null or all positive is that the
 71 variogram is null, i.e. the r-composition is constant. The fact that variograms
 72 and cross-variograms of r-compositions are subject to non-stochastic controls
 73 leads to the conclusion that, when based on raw data, they are spurious.

74 Under the assumption that the sample space is the whole real space endowed
 75 with the standard Euclidean space structure and geometry, or a subset with the
 76 induced structure and geometry, for $\mathbf{Z}(x)$ satisfying the second order stationary
 77 hypothesis, the following equalities hold:

$$\begin{aligned} \sum_{i=1}^D Z_i(x) &= \kappa, \\ \sum_{i=1}^D \mathbb{E}(Z_i(x)) &= \sum_{i=1}^D m_i = \kappa, \\ \sum_{i=1}^D (Z_i(x) - m_i) &= 0, \end{aligned} \quad (6)$$

78 with $\mathbb{E}(Z_i(x)) = m_i$, the expected value of $Z_i(x)$, $i = 1, 2, \dots, D$. Multiplying
 79 both sites of Eq. (6) by $(Z_j(x) - m_j)$ and taking expectations, it holds

$$\sum_{i=1}^D \text{cov} [(Z_i(x) - m_i), (Z_j(x) - m_j)] = 0, \quad j = 1, 2, \dots, D, \quad (7)$$

80 and therefore, for any lag h ,

$$C_j(h) = - \sum_{i \neq j} C_{ij}(h), \quad j = 1, 2, \dots, D, \quad (8)$$

81 where $C_j(h)$ stands for the auto-covariance of component j , and $C_{ij}(h)$ for the
82 cross-covariance of components i and j . Consequently, also the cross-covariances
83 cannot be all null, and some of them have necessarily to be negative. Being
84 subject to non-stochastic controls, they are spurious.

85 As summarized in Pawlowsky-Glahn and Burger (1992), the problems de-
86 rived from the nature of spatially distributed compositional data, when the raw
87 data are analysed, are

- 88 1. The mathematical necessity of at least one non-zero cross-covariance.
- 89 2. The bias towards negative cross-covariances.
- 90 3. The singularity of the cross-covariance matrix for any lag h .
- 91 4. The distorted description and interpretation of the spatial dependence
92 between the compositional variables under study.

93 Nowadays we know that the problem of spurious spatial covariance or correla-
94 tion is generated by the fact that compositional data are analysed as *real data*,
95 with the usual Euclidean geometry. In fact, most statistical methods have been
96 developed for real data without constraints under the implicit assumption that
97 the Euclidean geometry holds. This means that the difference between obser-
98 vations is measured as an absolute difference, that the sum and its opposite
99 make sense. This holds even with constraints, i.e. restricting the support of the
100 sample to a subset of real space without changing the geometry.

101 **3. The beginning — 1986: The additive log-ratio approach**

The initial approach (Pawlowsky, 1986; Pawlowsky-Glahn and Olea, 2004) was to use the additive log-ratio (alr) transformation (Aitchison, 1982, 1986). The r-composition is transformed into log-ratios as

$$\mathbf{W}(x) = \left(\ln \frac{Z_1}{Z_D}, \ln \frac{Z_2}{Z_D}, \dots, \ln \frac{Z_{D-1}}{Z_D} \right),$$

102 thus obtaining a regionalized vector of $D - 1$ components which can be treated
103 using cokriging. As we are aware nowadays, this was done under the implicit
104 assumption that the Euclidean geometry holds for alr transformed vectors.
105 Under this assumption the alr-transformation leads to BLU (Best Linear Un-
106 biased) estimates (Pawlowsky-Glahn and Egozcue, 2002). Nevertheless, soon
107 problems appeared, like the fact that cokriging seemed to lead to worse re-
108 sults than kriging, a fact that stands in contradiction with theoretical results
109 (Pawlowsky-Glahn and Olea, 2004, p. 160-161). The reasons for these problems
110 could not be explained in a consistent way until the algebraic-geometric struc-
111 ture of the sample space of compositional data was recognized (Aitchison et al.,
112 2002; Billheimer et al., 2001; Pawlowsky-Glahn and Egozcue, 2001) and the alr
113 was understood within this framework. Essentially, the problem was the com-
114 putation of variances and covariances using the alr coordinates, which at that
115 moment was not clear.

116 The covariance structure of compositional data can be described by the so-
117 called variation matrix (Aitchison, 1982, 1986). This matrix contains the vari-
118 ances of each possible log-ratio of pairs of compositional parts. It was shown that
119 the variation matrix completely describes the covariance structure of the compo-
120 sition, independently of which transformation is used to analyse the data. These
121 facts inspired the introduction of the spatial structure of r-compositions, first
122 defined in Pawlowsky (1986) and summarised in Pawlowsky-Glahn and Burger

123 (1992) and (Pawlowsky-Glahn and Olea, 2004, p. 29):

DEFINITION 3.1 [*Spatial covariance structure*] *The spatial covariance structure of a D -part r -composition is defined as the set of functions of the lag h*

$$\sigma_{ij.k\ell}(h) = \text{Cov} \left(\ln \frac{Z_i(x)}{Z_k(x)}, \ln \frac{Z_j(x+h)}{Z_\ell(x+h)} \right), \quad i, j, k, \ell \in 1, 2, \dots, D, x \in \mathcal{D}, .$$

124 At a first glance, the geostatistical analysis of $\mathbf{W}(x)$ can be performed as a
125 cokriging. This means that variograms and cross-variograms have to be fitted
126 to their empirical versions. However, the spatial covariance structure allows the
127 modelling of each component of $\sigma_{ij.k\ell}(h)$ by a simple variogram, thus avoiding
128 modelling of cross-variograms. A matrix transformation can transform the spa-
129 tial covariance structure into the cross-variograms required for a cokriging of
130 $\mathbf{W}(x)$.

131 As stated in Pawlowsky-Glahn and Burger (1992), the most difficult part—
132 compared to a spatial analysis of several variables—is that, in addition to the
133 usual difficulties, problems have to be reformulated in terms of logratios, and
134 interpretation and description of spatial dependencies have to be made in the
135 same terms.

136 **4. The breakthrough 2000...**

137 Around the year 2000, compositional data analysis attains a further maturity
138 level. The achievements can be summarized in two main points: (1) The sim-
139 plex, as sample space of compositional data, is endowed with a Euclidean space
140 structure, called Aitchison geometry (Pawlowsky-Glahn and Egozcue, 2001);
141 and (2) compositional data are no longer conceived as vectors constrained to
142 a constant sum but as equivalence classes of proportional vectors with positive
143 components (Barceló-Vidal et al., 2001). These new points of view influenced

144 the way of identifying and analysing r-compositions and they are briefly de-
145 scribed in the following sections.

146 Subsequent developments (Tolosana-Delgado, 2006), based on the sample
147 space approach and the *Principle of Working in Coordinates* (Mateu-Figueras et al.,
148 2011; Pawlowsky-Glahn, 2003), proved the potential for the log-ratio approach
149 within the Aitchison geometry of the simplex, setting the foundations for a
150 rigorous theory. Based on the principles of scale invariance, subcompositional
151 dominance, and permutation invariance, the operations of perturbation, power-
152 ing, and the inner product associated to the distance introduced by Aitchison
153 (1982, 1986, 1997), provide, as mentioned, the simplex with a Euclidean space
154 structure (Billheimer et al., 2001; Pawlowsky-Glahn and Egozcue, 2001), differ-
155 ent, but nonetheless isometric to the Euclidean space structure of real space.
156 The Euclidean space structure of the simplex was termed *Aitchison geometry*
157 in Pawlowsky-Glahn and Egozcue (2001). It opened up the door to a deeper
158 understanding of the nature of compositional data, of the available methods to
159 analyse them, and of the problems linked to different approaches. In particu-
160 lar, the advantage of using isometric log-ratio transformations was recognised.
161 Within this family of transformations, those known as balances (Egozcue et al.,
162 2003; Egozcue and Pawlowsky-Glahn, 2005) have shown a high potential based
163 on their interpretability, and can be used for spatial analysis of compositional
164 data.

165 4.1. *Compositions are representatives of equivalence classes*

166 In Aitchison (1986), the so called *principle of scale invariance of composi-*
167 *tions* was formulated. It states that the analysis of a composition must remain
168 invariant when the composition is multiplied by any positive constant. This was
169 the motivation for preconising the use of log-contrasts as the main tool in the
170 analysis. Log-contrast are combinations of the logarithms of the parts such that,

171 when all parts of the composition are multiplied by a positive constant, the value
172 of the combination remains unaltered. Also, vectors of positive components are
173 reduced to constant sum by using the closure operation. These concepts were
174 clear from the beginning of compositional data analysis, but there was a lack of
175 mathematical formulation reflected in the wording of compositional data anal-
176 ysis. For instance, when referring to the closure problem as the only origin of
177 pitfalls in compositional data analysis.

178 The progress consists in thinking that all vectors having proportional positive
179 components are equivalent and convey the same compositional information. A
180 composition is then an equivalence class which can be represented by choosing an
181 arbitrary element of the class. Equivalence classes can be represented in many
182 ways and each choice defines a potential sample space, whether constraint to
183 a constant sum or not (see explanations in Pawłowsky-Glahn et al., 2015c, ch.
184 2). When compositional data are represented as data subject to a constant sum
185 constraint, their sample space is a simplex, and the simplex is nothing else but a
186 choice of one out of all the possible sample spaces of compositions. This choice is
187 not only convenient because it is the usual choice in practice, but also because
188 it is mathematically easy to define a meaningful and interpretable Euclidean
189 vector space structure in the simplex (Pawłowsky-Glahn and Egozcue, 2001).

190 Other representations of compositions are possible. For instance, when com-
191 positions of air pollutants are expressed in $\mu\text{g}/\text{m}^3$ or solutes are given in Mol per
192 liter, concentrations do not add to a constant and they are not represented in
193 the simplex. Simply, the representative of the equivalence class has been taken
194 in another way, but still the ratios of the parts are the relevant information. In
195 these kind of representations, perturbation is also easily interpretable. Other
196 possibilities are less intuitive, for instance, when compositions are represented
197 in an orthant of a hypersphere (e.g. Wang et al., 2007).

198 It is remarkable that this interpretation of compositions as equivalence classes
199 only arose in 2000. This may be the reason why, in the decade from 1980 to
200 1990, concentrations in units like mol/litre, concentration of a single element,
201 or just removing a large component, were considered to be non-compositional
202 and, consequently, free of the difficulties of analysing compositional data.

203 The influence of these new concepts in geostatistics is reflected in the iden-
204 tification of what is an r-composition, independently of whether the collected
205 data are closed to a constant or not.

206 4.2. Aitchison geometry of the simplex and consequences

207 The simplex endowed with perturbation (the compositional sum), powering
208 (compositional multiplication by real numbers) and Aitchison distance, con-
209 stitute a $(D - 1)$ -dimensional Euclidean vector space (Billheimer et al., 2001;
210 Pawlowsky-Glahn and Egozcue, 2001). The Euclidean space structure of the
211 simplex was termed *Aitchison geometry* in Pawlowsky-Glahn and Egozcue (2001).
212 The value of this mathematical result is supported by the fact that perturbation
213 is an interpretable operation in most compositional scenarios. In fact, perturba-
214 tion can be interpreted as filtering in geochemistry or particle size analysis; or
215 as the Bayes formula for probabilities (for details, see Pawlowsky-Glahn et al.,
216 2015c, ch. 2).

217 The Aitchison geometry points out that orthonormal basis of the space ex-
218 ist, and that the corresponding (Cartesian) coordinates can efficiently repre-
219 sent compositions; orthogonal projections are possible; the concepts of linear
220 combination, linear dependence, Euclidean distances, and all the typical geo-
221 metrical elements are available. All these tools are readily used once composi-
222 tions are represented by their coordinates with respect to a basis of the space,
223 as perturbation is the sum in coordinates, powering is scaling in coordinates,
224 and the Aitchison distance is the standard Euclidean distance between coor-

225 dinates. This constitutes the core of the *Principle of Working in Coordinates*
226 (Mateu-Figueras et al., 2011).

An important step ahead is the construction of orthonormal (Cartesian) coordinates in the simplex. The function assigning orthonormal (Cartesian) coordinates to a composition has been named isometric log-ratio transformation (ilr) (Egozcue et al., 2003). The ilr-transformation is not unique, as there are infinitely many basis of the space. As a consequence, it was clear that the alr-transformation is an assignation of coordinates with respect to an oblique basis (Egozcue et al., 2003), while the centered log-ratio transformation (clr) (Aitchison, 1986)

$$\text{clr}(\mathbf{Z}(x)) = \left(\ln \frac{Z_1(x)}{g(\mathbf{Z}(x))}, \ln \frac{Z_2(x)}{g(\mathbf{Z}(x))}, \dots, \ln \frac{Z_D(x)}{g(\mathbf{Z}(x))} \right),$$

where $g(\mathbf{Z}(x))$ is the geometric mean of the $\mathbf{Z}(x)$ components, gives coordinates with respect to a generating system of the simplex. The clr-transformation was not used for geostatistical analysis, as its covariance matrix is always singular. It is, nevertheless, extremely useful for computation in compositional data analysis. For example, the ilr-coordinates are readily obtained through a clr-transformation as

$$\text{ilr}(\mathbf{Z}(x)) = V \text{clr}(\mathbf{Z}(x)),$$

227 where V , called contrast matrix (Egozcue et al., 2011; Pawlowsky-Glahn et al.,
228 2015c), is a $(D, D - 1)$ -matrix satisfying the property that $V'V$ is the identity
229 matrix. An easy way of building coordinates, called balances, was introduced
230 in (Egozcue et al., 2003; Egozcue and Pawlowsky-Glahn, 2005). The procedure,
231 called sequential binary partition (SBP), provides such contrast matrices, and
232 the resulting ilr-coordinates are called *balances* (Egozcue and Pawlowsky-Glahn,
233 2005).

234 It is remarkable that ilr, alr, and clr transformations are different assign-
235 nations of coordinates to a composition more than different transformations
236 leading to different approaches. An important point is that, within the Aitchi-
237 son geometry of the simplex, the predictors used in all classes of kriging are
238 linear, as they are linear combinations of coordinates. However, in the case
239 of alr-coordinates distances and covariances should be handled very carefully,
240 paying special attention to the fact that they are representations in an oblique
241 coordinate system. This explains the problems detected when using cokriging
242 on alr-coordinates (Pawłowsky-Glahn and Olea, 2004, p. 108), where this fact
243 was not taken into account.

244 *4.3. Cokriging of regionalized compositions*

245 Initially, the problems for cokriging of r-compositions appeared to be centred
246 on the modelling of cross-variograms of log-ratio transformed data, although it
247 was known that a simple matrix transformation leads from the matrix-valued
248 variation-variogram, the matrix of variograms of all possible simple log-ratios, to
249 any log-ratio representation (Pawłowsky, 1986; Pawłowsky-Glahn and Burger,
250 1992; Pawłowsky-Glahn and Olea, 2004; Tolosana-Delgado, 2006; Tolosana-Delgado et al.,
251 2011). Later, Tolosana-Delgado and Boogaart (2013) recognised the potential
252 of the matrix-valued variation-variogram, specially to model cross-variograms
253 using first a Linear Model of Coregionalisation for the matrix-valued variation-
254 variogram, followed by a matrix transformation to obtain the corresponding
255 variograms and cross-variograms for the coordinates chosen by the scientist to
256 represent the available data. Note that the matrix-valued variation-variogram is
257 a matrix with all its entries simple variograms and no cross-variogram. Standard
258 cokriging can then be applied to obtain the desired predictions. In summary,
259 spatial compositional data analysis consists in the following steps (Tolosana-Delgado and Boogaart,
260 2013):

- 261 1. transform the D -part compositional vectors into $(D - 1)$ -dimensional real
262 vectors by means of a convenient isometric log-ratio (ilr) transformation;
- 263 2. apply any standard geostatistical technique to the vectors obtained;
- 264 3. back-transform interpolated and/or simulated scores back using the ilr
265 inverse.

266 To model necessary variograms and cross-variograms

- 267 • compute the matrix-valued variation matrix and adjust a Linear Model of
268 Coregionalisation;
- 269 • apply the corresponding matrix transformation to obtain the desired matrix-
270 valued variogram (containing variograms in the diagonal and cross-variograms
271 off-diagonal) of the ilr transformation used before.

272 Details of the procedure can be found in Tolosana-Delgado and Boogaart (2013).

273 Note that, as stated in Tolosana-Delgado et al. (2008a), the proposed pro-
274 cedure leads to BLU estimators when performing cokriging.

275 *4.4. Simplicial Indicator Kriging*

276 The recognition of the Euclidean vector space structure of the sample space
277 of compositional data and the understanding that probabilities can be con-
278 sidered to be a composition allowed to solve the problems intrinsic to Indicator
279 Kriging (Pawłowsky-Glahn et al., 2006; Tolosana-Delgado, 2006; Tolosana-Delgado et al.,
280 2008c,b). By construction, Simplicial Indicator Kriging avoids all the known
281 problems associated with usual Indicator Kriging (Journel, 1983), namely neg-
282 ative predictions, order relation violations, or predictions larger than one.

283 *4.5. Further developments — 2015: Cokriging r -compositions with a total*

284 As mentioned before, compositional data are multivariate positive real data
285 that carry only relative information, and can be represented simply taking clo-
286 sure, i.e. taking proportions or concentrations. In this case, the information

287 about their total sum is lost. In some cases, additionally to the composition,
288 the sum of some of the positive variables, called total, can be informative or
289 of interest. Consequently, the need of a joint analysis of composition and to-
290 tal arises. Some possibilities were studied in Pawlowsky-Glahn et al. (2015a)
291 which concluded that the chosen total can be included as an additional coor-
292 dinate to those coming from the composition. This applies to r-compositions
293 where some regionalized total is of interest. The geostatistical analysis can be
294 conducted by cokriging of compositional ilr-coordinates, jointly with the coor-
295 dinate of the total.

296 A first application of this procedure was performed in Pawlowsky-Glahn et al.
297 (2015b) although the main goal was dimension reduction of a geochemical data
298 set. The problem appears when applying compositional techniques of dimension
299 reduction since, after orthogonal projections, the original units of the compo-
300 sition are lost. In order to recover original units, cokriging of ilr-coordinates
301 of the composition is carried out jointly with the sum of initial concentrations.
302 This joint cokriging of ilr-coordinates with supplementary real variables appears
303 to be a promising technique in compositional geostatistics.

304 **5. Other approaches**

305 Not many attempts have been made to find spatial interpolation meth-
306 ods for regionalized compositional data. Methods that comply with nonneg-
307 ativity and the representation as data constraint to a constant sum include
308 nearest neighbor interpolation, triangulation, local sample (arithmetic) mean,
309 and inverse distance interpolation, which are described in Isaaks and Srivastava
310 (1989). Another approach, called *compositional kriging*, was introduced by
311 Walwoort and de Gruijter (2001). All of them are implicitly based on the as-
312 sumption that the sample space of compositional data is the simplex as a con-

313 straint subset of real space, and that they obey the induced geometry, i.e. the
314 standard Euclidean geometry. This fact implies the assumption that composi-
315 tional data carry absolute and not relative information, a decision that lies with
316 the researcher analysing the data. Furthermore, as stated by Walwoort and de Gruijter
317 (2001), the former methods do not take the spatial structure into account, but
318 neither does *compositional kriging* completely, as it does not take into account
319 cross-correlations, and thus cross-variograms, to avoid problems with spurious
320 correlation. As shown by Pawlowsky-Glahn and Egozcue (2002), even using
321 the alr representation of compositional data leads to BLU estimators within the
322 Aitchison geometry of the simplex (Pawlowsky-Glahn and Egozcue, 2001), and
323 numerical comparisons of results based on different assumptions for the struc-
324 ture of the sample space make no sense. Whichever is the assumption made
325 by the scientist, spatial interpolation using cokriging will be optimal within the
326 assumed geometry.

327 6. Conclusions and comments

328 Reviewing the early developments in the spatial analysis of compositional
329 data, and in the analysis of compositional data in general, one can see the
330 evolution of the way of thinking on that type of data. One typical example is
331 the statement in Pawlowsky (1984) that $\mathbf{Z}(x) - \mathbf{Z}(x+h)$ is an r-composition for
332 any $x \in \mathcal{R}$ and any lag h . This is clearly not true, as it always yields at least
333 some non-positive numbers.

334 Another hurdle were the problems related to the alr transformation. After
335 understanding that the alr represents the data in an oblique basis of the sim-
336 plex, one can recognise two ways of proceeding: (1) to avoid the alr and use
337 only isometric log-ratio transformations, or (2) to take into account the oblique
338 nature and use appropriate matrix transformations to obtain consistent results.

339 The first approach is straightforward and safe, the second requires more care.
340 It is up to the researcher to choose which transformation is better suited for the
341 case he or she is dealing with.

342 One of the characteristics of cokriging ilr-coordinates is that the modelling
343 of cross-variograms can be afforded modelling the variation variograms, thus
344 avoiding the always difficult cross-variogram modelling.

345 The main conclusion is that analysing compositional data, regionalized or
346 not, is nowadays summarised by the *principle of working on coordinates*; it
347 transforms the compositional analysis into a standard geostatistical problem
348 where well known procedures can be applied without additional difficulties.

349 **Acknowledgements**

350 The authors thank two anonymous reviewers for their suggestions and com-
351 ments. This research has been supported by the *Spanish Ministry of Educa-*
352 *tion and Science* under project ‘METRICS’ (Ref. MTM2012-33236); and from
353 the *Agència de Gestió d’Ajuts Universitaris i de Recerca* of the *Generalitat de*
354 *Catalunya* under the project Ref. 2009SGR424.

355 **References**

- 356 Aitchison J. The statistical analysis of compositional data (with discussion).
357 Journal of the Royal Statistical Society, Series B (Statistical Methodology)
358 1982;44(2):139–77.
- 359 Aitchison J. The Statistical Analysis of Compositional Data. Monographs on
360 Statistics and Applied Probability. London (UK): Chapman & Hall Ltd.,
361 London (UK). (Reprinted in 2003 with additional material by The Blackburn
362 Press), 1986. 416 p.

- 363 Aitchison J. The one-hour course in compositional data analysis or composi-
364 tional data analysis is simple. In: Pawlowsky-Glahn V, editor. Proceedings
365 of IAMG'97 - The III Annual Conference of the International Association
366 for Mathematical Geology. Barcelona (E): International Center for Numerical
367 Methods in Engineering (CIMNE), Barcelona (E), 1100 p; volume I, II and
368 addendum; 1997. p. 3–35.
- 369 Aitchison J, Barceló-Vidal C, Egozcue JJ, Pawlowsky-Glahn V. A concise guide
370 for the algebraic-geometric structure of the simplex, the sample space for com-
371 positional data analysis. In: Bayer U, Burger H, Skala W, eds. Proceedings
372 of IAMG'02 – The VIII Annual Conference of the International Association
373 for Mathematical Geology. Selbstverlag der Alfred-Wegener-Stiftung, Berlin,
374 1106 p; volume I and II; 2002. p. 387–92.
- 375 Aitchison J, Shen SM. Logistic-normal distributions. Some properties and uses.
376 *Biometrika* 1980;67(2):261–72.
- 377 Barceló-Vidal C, Martín-Fernández JA, Pawlowsky-Glahn V. Mathematical
378 foundations of compositional data analysis. In: Ross G, editor. Proceedings
379 of IAMG'01 – The VII Annual Conference of the International Association
380 for Mathematical Geology. Cancun (Mex); 2001. p. 20 p.
- 381 Billheimer D, Guttorp P, Fagan W. Statistical interpretation of species compo-
382 sition. *Journal of the American Statistical Association* 2001;96(456):1205–14.
- 383 Chayes F. On correlation between variables of constant sum. *Journal of Geo-*
384 *physical Research* 1960;65(12):4185–93.
- 385 Chilès JP, Delfiner P. *Geostatistics - Modeling Spatial Uncertainty*. 2nd ed.
386 *Probability and Statistics*. United States of America: Wiley, 2012.

387 Egozcue JJ, Barceló-Vidal C, Martín-Fernández JA, Jarauta-Bragulat E, Díaz-
388 Barrero JL, Mateu-Figueras G. Elements of simplicial linear algebra and
389 geometry. In: Pawlowsky-Glahn and Buccianti (2011); 2011. p. 141–57. 378
390 p.

391 Egozcue JJ, Pawlowsky-Glahn V. Groups of parts and their balances in com-
392 positional data analysis. *Mathematical Geology* 2005;37(7):795–828.

393 Egozcue JJ, Pawlowsky-Glahn V, Mateu-Figueras G, Barceló-Vidal C. Isometric
394 logratio transformations for compositional data analysis. volume 35; 2003. p.
395 279–300.

396 Isaaks EH, Srivastava RM. *An Introduction to Applied Geostatistics*. Oxford
397 University Press, New York, NY (USA), 1989. 592 p.

398 Journel AG. Nonparametric estimation of spatial distributions. *Mathematical*
399 *Geology* 1983;15(3):445–68.

400 Mateu-Figueras G, Pawlowsky-Glahn V, Egozcue JJ. The principle of working
401 on coordinates. In: Pawlowsky-Glahn and Buccianti (2011); 2011. p. 31–42.
402 378 p.

403 Matheron G. *Traité de Géostatistique Appliquée*. volume I of *Mémoires du*
404 *Bureau de Recherches Géologiques et Minières (France)*, 14. Technip, Paris
405 (F), 1962. 333 p.

406 Matheron G. *Les Variables Régionalisées et leur Estimation—une Application*
407 *de la Théorie des Fonctions Aléatoires aux Sciences de la Nature*. Masson et
408 Cie., Paris (F), 1965. 305 p.

409 Pawlowsky V. On spurious spatial covariance between variables of constant
410 sum. *Science de la Terre, Sér Informatique* 1984;21:107–13.

- 411 Pawlowsky V. Räumliche Strukturanalyse und Schätzung ortsabhängiger Kom-
412 positionen mit Anwendungsbeispielen aus der Geologie. Ph.D. thesis; Fach-
413 bereich Geowissenschaften, Freie Universität Berlin, Berlin (D); 1986. 170
414 p.
- 415 Pawlowsky-Glahn V. Statistical modelling on coordinates. In: Thió-Henestrosa
416 S, Martín-Fernández JA, eds. Proceedings of CoDaWork'03, The 1st Composi-
417 tional Data Analysis Workshop. Girona (E): Universitat de Girona, ISBN
418 84-8458-111-X, <http://ima.udg.es/Activitats/CoDaWork2003/>; 2003. .
- 419 Pawlowsky-Glahn V, Buccianti A, eds. Compositional Data Analysis: Theory
420 and Applications. John Wiley & Sons.; 2011. 378 p.
- 421 Pawlowsky-Glahn V, Burger H. Spatial structure analysis of regionalized com-
422 positions. *Mathematical Geology* 1992;24(6):675–91.
- 423 Pawlowsky-Glahn V, Egozcue JJ. Geometric approach to statistical analysis
424 on the simplex. *Stochastic Environmental Research and Risk Assessment*
425 (SERRA) 2001;15(5):384–98.
- 426 Pawlowsky-Glahn V, Egozcue JJ. BLU estimators and compositional data.
427 volume 34; 2002. p. 259–74.
- 428 Pawlowsky-Glahn V, Egozcue JJ, Lovell D. Tools for composi-
429 tional data with a total. *Statistical Modelling* 2015a;15(2):175–90.
430 doi:10.1177/1471082X14535526, Nov. 25, 2014.
- 431 Pawlowsky-Glahn V, Egozcue JJ, Olea RA, Pardo-Igúzquiza E. Cokriging of
432 compositional balances including a dimension reduction and retrieval of origi-
433 nal units. *The Journal of The Southern African Institute of Mining and*
434 *Metalurgy, SAIMM* 2015b;115:59–72.

435 Pawlowsky-Glahn V, Egozcue JJ, Tolosana-Delgado R. Modeling and Analysis
436 of Compositional Data. Statistics in practice. John Wiley & Sons, Chichester
437 UK, 2015c. 272 pp.

438 Pawlowsky-Glahn V, Olea RA. Geostatistical Analysis of Compositional Data.
439 Number 7 in Studies in Mathematical Geology. Oxford University Press, 2004.
440 Editor: DeGraffenreid, Jo Anne.

441 Pawlowsky-Glahn V, Tolosana-Delgado R, Egozcue JJ. Simplicial indicator
442 kriging. 2006. 4 p.

443 Tolosana-Delgado R. Geostatistics for Constrained Variables:
444 Positive Data, Compositions and Probabilities. Application
445 to Environmental Hazard Monitoring. 2006. p. 198. URL:
446 http://www.tesisenxarxa.net/TDX-0123106-122444/index_an.html;
447 198 p.

448 Tolosana-Delgado R, Boogaart KGvd. Joint consistent mapping of high-
449 dimensional geochemical surveys. Mathematical Geosciences 2013;45:983-
450 1004.

451 Tolosana-Delgado R, Boogaart KGvd, Pawlowsky-Glahn V. Geostatistics for
452 compositions. In: Pawlowsky-Glahn and Buccianti (2011); p. 73-86. 378 p.

453 Tolosana-Delgado R, Egozcue JJ, Pawlowsky-Glahn V. Cokriging of composi-
454 tions: log-ratios and unbiasedness. In: Ortiz JM, Emery X, eds. Geostatistics
455 Chile 2008. Gecamin Ltd., Santiago, Chile, 2 vols, 1188 p.; 2008a. p. 299-308.

456 Tolosana-Delgado R, Pawlowsky-Glahn V, Egozcue JJ. Indicator kriging with-
457 out order relation violations. Mathematical Geosciences 2008b;40:327-47.

458 Tolosana-Delgado R, Pawlowsky-Glahn V, Egozcue JJ. Simplicial indicator
459 kriging. Journal of China University of Geosciences 2008c;19:65-71.

- 460 Walwoort DJ, de Gruijter JJ. Compositional kriging: a spatial interpolation
461 method for compositional data. *Mathematical Geology* 2001;33(8):951–66.
- 462 Wang H, Liu Q, Mok HMK, Fu L, Tse WM. A hyperspherical transformation
463 forecasting model for compositional data. *European Journal of Operational*
464 *Research* 2007;179:459–68.