

Math Geosci (2009) 41: 869–886
DOI 10.1007/s11004-009-9216-6

Grain-Size Control on Petrographic Composition of Sediments: Compositional Regression and Rounded Zeros

Raimon Tolosana-Delgado · Hilmar von Eynatten

Received: 16 August 2007 / Accepted: 31 October 2008 / Published online: 26 March 2009
© The Author(s) 2009. This article is published with open access at Springerlink.com

Abstract It is well-known that sediment composition strongly depends on grain size. A number of studies have tried to quantify this relationship focusing on the sand fraction, but only very limited data exists covering wider grain size ranges. Geologists have a clear conceptual model of the relation between grain size and sediment petrographic composition, typically displayed in evolution diagrams. We chose a classical model covering grain sizes from fine gravel to clay, and distinguishing five types of grains (rock fragments, poly- and mono crystalline quartz, feldspar and mica/clay). A compositional linear process is fitted here to a digitized version of this model, by (i) applying classical regression to the set of all pairwise log-ratios of the 5-part composition against grain size, and (ii) looking for the compositions that best approximate the set of estimated parameters, one acting as slope and one as intercept. The method is useful even in the presence of several missing values. The linear fit suggests that the relative influence of the processes controlling the relationship between grain size and sediment composition is constant along most of the grain size spectrum.

Keywords Censored data · Compositional Data Analysis · Moore–Penrose generalized inverse · Sedimentary petrography

1 Introduction

It is common knowledge that the apparent composition (both geochemical and petrographic) of a sediment strongly depends on its grain size (e.g., Pettijohn 1957).

R. Tolosana-Delgado (✉) · H. von Eynatten
Department of Sedimentology and Environmental Geology, University of Göttingen,
37077 Göttingen, Germany
e-mail: raimon.tolosana@geo.uni-goettingen.de

H. von Eynatten
e-mail: hilmar.von.eynatten@geo.uni-goettingen.de

Several studies have tried to characterize or quantify that influence, and to distinguish it from other factors that control sediment composition (e.g., Johnsson 1993), such as bedrock geology in the source area (Nesbitt and Young 1984; Grantham and Velbel 1988), mixing of material from different sources (Arribas and Tortosa 2003; Palomares and Arribas 1993) and alteration processes during weathering and transport (Nesbitt and Markovics 1997; Zhang et al. 2002). So far, most studies that aimed to shed light on the quantitative relations between grain size and sediment composition have a low grain size resolution (e.g., in sand-silt-clay, Lim et al. 2006), placed their attention on the sand domain only (e.g., Grantham and Velbel 1988; Nesbitt and Young 1996; Arribas and Tortosa 2003; Solano-Acosta and Dutta 2005; Kiminami and Fujii 2007), and/or use debatable statistical techniques that ignore the compositional nature of the data (e.g., Zhang et al. 2002; Chandrajith et al. 2001). Whitmore et al. (2004) presented both chemical and petrographic data with a higher grain-size resolution (-1 to 4 phi-grade, in 1 phi step) from modern sediments from rivers of Papua New Guinea; however, the fine-grained fractions (silt and clay) are not differentiated. Their main conclusions are that (i) the variation with grain size is much greater than variation with time or between different localities, and (ii) the most important variation with decreasing grain size is decreasing abundance of rock fragments (a conclusion reached by almost all the authors cited previously). This brief literature survey is by no means exhaustive, but nevertheless it is representative. The results of all these papers cannot be used for a global model of sediment evolution with grain size because: (i) only the sand fraction is analyzed in detail, and (ii) almost all study a selected area with specific geologic, climatic and geomorphologic characteristics that cannot account for global-scale Earth surface variations.

With the aim of providing such a general model, Blatt et al. (1972) published a plot on “the probable relationship between grain size and detrital fragment composition, based on the limited data currently available” (Fig. 1). The plot covers a grain size spectrum ranging from fine gravel to clay, and gives the composition on five petrographic grain types: rock fragments, poly-crystalline quartz, mono-crystalline quartz, feldspar, and mica. The latter includes micas s.s. and clay minerals. Grain size is given in ϕ scale, corresponding to the negative logarithm to the basis 2 of the grain diameter in mm. With this figure, Blatt et al. (1972) tentatively summarized a series of processes and controlling factors:

- the typical crystal size of each mineral or mineral aggregate (e.g., monocrystalline quartz and feldspar crystals are smaller than rock fragments, and clay minerals are smaller than quartz and feldspar)
- the global crustal abundance of specific minerals and mineral associations (e.g., pure quartz rocks are less common than polymineralic rocks; thus, at the coarser sizes, unspecified rock fragments are more abundant than polycrystalline quartz)
- the effect of comminution, which may produce from a rock fragment grain smaller grains of any type; from a grain of a given “mineral class” (Q_m , F and M) smaller grains of the same class; and from a grain of polycrystalline quartz either poly- or mono-crystalline quartz
- the effect of chemical alteration, fundamentally generating clay minerals (M class) from feldspar (or from rock fragments)

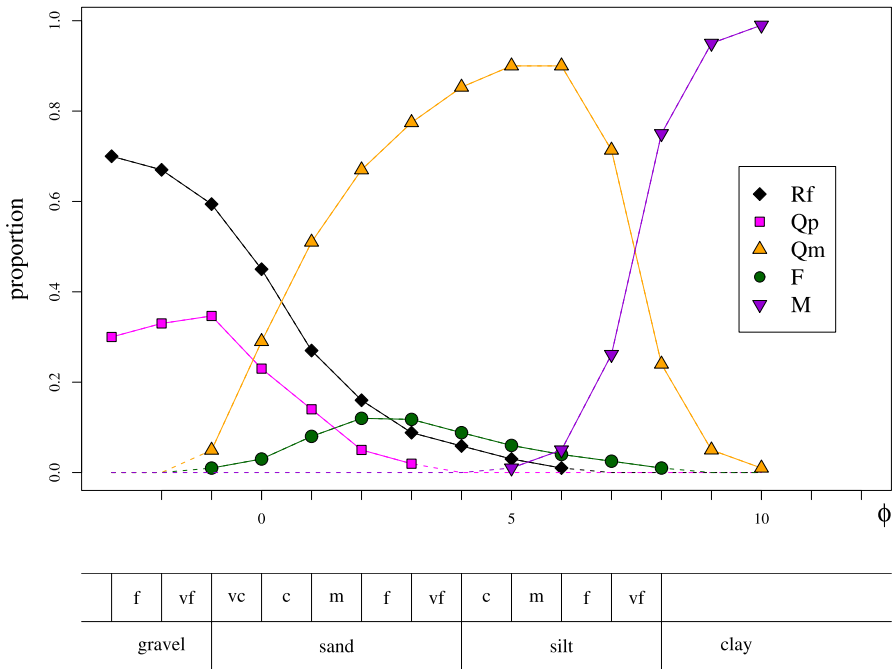


Fig. 1 Data set digitized from the original figure published by Blatt et al. (1972, p. 301). Legend: (R_f)—rock fragments, (Q_p)—poly-crystalline quartz, (Q_m)—mono-crystalline quartz, (F)—feldspar, and (M)—mica including clay minerals. The grain-size scale included below is modified from Wentworth (1922, cited by Pettijohn 1957); abbreviations mean: (vc)—very coarse, (c)—coarse, (m)—medium, (f)—fine and (vf)—very fine

- the balanced effect of mineral dissolution and precipitation, eliminating or adding minerals to/from the sediment body as it is transported from source to sink, and during accumulation in the sedimentary basin.

Note that no information on the absolute mass of each grain size and type is available. This implies that these processes and relationships cannot be written in terms of mass balances, and one cannot be sure whether mass is added or subtracted from the sediment when “evolving” from coarser to finer grain fractions. For example, the rapid increase of M content on the fine-grained end of the plot could be explained as the precipitation of authigenic clay minerals (adding M); the accumulation of residues of dissolution of all other minerals (eliminating mostly F , but even Q at the finer end of the plot); the disaggregation and comminution of mica crystals from rock fragments, and alteration of feldspar (transfer $R_f, F \rightarrow M$); or any combination of these processes.

Table 1 shows a digitalization at regular intervals of the original plot, that will be used here as data set. Our goal is to fit a model to this data, by means of regression, and interpret the results according to the considerations given above. It is evident that a raw linear fit will utterly fail, as the curves present s-shapes, and some have a maximum. Thus, one would need at least two lines to (poorly) represent them.

Table 1 Data set based on a digitalization of the Blatt et al. (1972) plot. Note that no effort has been made to keep the total sum to 100%. This will be done in the subsequent analysis. Empty cells correspond to very small values, that should be considered zero from a practical point of view

ϕ	R_f	Q_p	Q_m	F	M	Sum
-3	70	30				100
-2	67	33				100
-1	60	35	5	1		101
0	45	23	29	3		100
1	27	14	51	8		100
2	16	5	67	12		100
3	9	2	79	12		102
4	6		87	9		102
5	3		90	6	1	100
6	1		90	4	5	100
7			71	2.5	26	99.5
8			24	1	75	100
9			5		95	100
10			1		100	101

Moreover, the data indeed form a composition (with positive parts and constant sum for each possible ϕ value), and its relative nature has already been pointed out. Thus, it seems more reasonable to follow the log-ratio transformation approach of Aitchison (1986), which has already been successfully applied to sediment evolution modeling (von Eynatten 2004; Noda 2005; Buccianti et al. 2006) and can accommodate all relative information as a linear process.

An early application of regression to compositional data was presented by Daunisi-Estadella et al. (2002). However, the data set analyzed here has lots of small values, most of them zero from a practical point of view. The presence of zeroes has been regarded as an almost intractable problem of compositional analysis based on log-ratios. The typical strategy is to impute the zero value by some small quantity, usually a proportion of the detection limit (Aitchison 1986; Martín-Fernández et al. 2000; Martín-Fernández 2001). Recently, van den Boogaart et al. (2006) proposed an alternative approach based on working with the observed subcompositions. Following this idea, a secondary goal of this contribution is to provide a sensible analysis within the log-ratio framework for the kind of regression problems where the response is a compositional vector with rounded zeroes.

2 Notation and Revision of Basic Concepts

2.1 The Algebraic-Geometric Structure of the Simplex

Denote a D -part composition as a row vector $\mathbf{x} = [x_1, \dots, x_D]$, and its sample space as \mathcal{S}^D , the D -part simplex. According to Billheimer et al. (2001) and Pawłowsky-Glahn and Egozcue (2001), this set can be given an Euclidean space structure with the following operations. For $\mathbf{x}, \mathbf{y} \in \mathcal{S}^D$ compositions, and $\alpha \in \mathbb{R}$ a real value,

- the Abelian group operation, perturbation, is defined as $\mathbf{x} \oplus \mathbf{y} = \mathcal{C}[x_1 y_1, \dots, x_D y_D]$

- the scalar multiplication, powering, is defined as $\alpha \odot \mathbf{x} = \mathcal{C}[x_1^\alpha, \dots, x_D^\alpha]$
- and the scalar or inner product, is defined as

$$\langle \mathbf{x}, \mathbf{y} \rangle_A = \sum_{i=1}^D \text{clr}_i(\mathbf{x}) \cdot \text{clr}_i(\mathbf{y}), \tag{1}$$

where

$$\text{clr}_i(\mathbf{x}) = \ln \frac{x_i}{\sqrt[D]{x_1 \cdots x_D}} \tag{2}$$

denotes the i th component of the centered log-ratio transformation and the closure operation $\mathcal{C}[\cdot]$ closes its argument to total unit sum (Aitchison 1986, 2002). The scalar product (1) has as associated distance

$$d_A^2(\mathbf{x}, \mathbf{y}) = \frac{1}{D} \sum_{i < j} \left(\ln \frac{x_i}{x_j} - \ln \frac{y_i}{y_j} \right)^2, \tag{3}$$

which has the linearity properties of an Euclidean distance, i.e., it is translation invariant, $d_A(\mathbf{x} \oplus \mathbf{z}, \mathbf{y} \oplus \mathbf{z}) = d_A(\mathbf{x}, \mathbf{y})$, and scalable, $d_A^2(\lambda \odot \mathbf{x}, \lambda \odot \mathbf{y}) = \lambda^2 d_A^2(\mathbf{x}, \mathbf{y})$ (Aitchison 1997). Moreover, it is subcompositionally dominant, i.e., the distance between two vectors measured in a subcomposition is always equal or smaller than the distance measured on the full composition.

2.2 The Compositional Linear Trend

To express a compositional response as a (linear) function of a real variable in the scope of this geometric structure one can write

$$\mathbf{z} = \mathbf{b}_0 \oplus x \odot \mathbf{b}_1. \tag{4}$$

Here, $\mathbf{b}_0 = [b_{01}, \dots, b_{0D}]$ plays the role of an intercept, and $\mathbf{b}_1 = [b_{11}, \dots, b_{1D}]$ is a slope. Both vectors are compositions in the same simplex as \mathbf{Z} . Rewriting the compositional vectors in terms of their components, such a trend becomes

$$[z_1, \dots, z_D] = \mathcal{C}[b_{01} \cdot e^{\beta_1 x}, \dots, b_{0D} \cdot e^{\beta_D x}],$$

with $\beta_i = \log(b_{1i})$. This expression can be interpreted as D non-interacting exponential decay/growth processes of rate β_i where no mass transfer occurs between or within the parts (Egozcue and Pawlowsky-Glahn 2006). An important property of this kind of linear processes is that the decay/growth rate of one part against the decay/growth rate of another part is constant along the explanatory variable x (von Eynatten et al. 2003).

2.3 Missing Values and Zeros in Compositional Data

Missing values in multivariate analysis can be assigned to several classes, according to the mechanism that generated them. The essential point is whether missing a value

depends on the observed variables and/or on the missing variable itself. In the present case, we distinguish the following cases (Rubin 1976; Murray 1979), missing-at-random values, where missing a value depends on the observed variables, but does not depend on the missed value itself, and truncated values, from which we know that they are below or above a threshold, but do not know the actual value (also known as censored values). They are a kind of not-missed at random, as the fact of missing a value depends on the value itself. Replacing by a central value would in this case ignore the extra information on the threshold, thus introducing some bias.

As the compositional operations given before are based on logarithms (equations (1)–(3)), zero values add further complexity to the analysis of a compositional data set. Two kinds of zeroes might appear. Rounded zeros, or values below the instrumental detection limit, corresponding to the truncated values of classical analysis. The classical strategy in this case is to replace the zero value by a suitable small quantity, e.g., $2/3$ of the detection limit of the zero variable (Martín-Fernández et al. 2000; Martín-Fernández 2001). A variation on that idea uses an EM algorithm to replace each rounded zero by a randomly-generated value, conditional on the non-zero parts and the detection limit (Palarea-Albaladejo et al. 2007). The second kind of zeros are structural zeros, values which are actually zero due to a physical process or limitation (e.g., chemical incompatibility of quartz and feldspatoids: if one is present, the other must have a concentration of zero). This is not a missing value, and in this case, the recommended treatment is to split the sample in two populations, according to that variable with zeros (Aitchison 1986).

Following a different line, van den Boogaart et al. (2006) put forward an approach for the estimation of the mean and variance of a compositional data set with missing values, suitable to incorporate a limited amount of zeros of any kind. The idea is that each datum contributes only to the determination of these statistics in the observed subcomposition. For instance (see Table 1), the datum at $\phi = 4$ would only contribute to the subcomposition $[R_f, Q_m, F]$. This paper will follow a similar approach, each datum will contribute to those pairwise log-ratios where both parts are observed. Then, one can distinguish three situations:

- If both parts are observed, this gives rise to an “observed” log-ratio, a datum, e.g., for $\phi = 4$, all $R_f/F = 6/9$, $R_f/Q_m = 6/87$, $Q_m/F = 87/9$, and their inverse ratios are observed data.
- If only one of both parts is observed and the other is below the detection limit, one can still say something about the log-ratio: it is a censored value, e.g., with $Q_m = 87\%$, if the detection limit for Q_p (not observed) is taken as 1%, we can say that $\ln(Q_p/Q_m)$ must be below $\ln(0.01/0.87)$ in $\phi = 4$.
- If both parts are not observed (either below the detection limit or missing-at-random), one could actually have any possible value for the log-ratio. The ratio is thus behaving quite like a missing-at-random value of conventional multivariate analysis, e.g., knowing the detection limits of Q_p and M , say $n\%$ and $d\%$, we can say that the true non-observed values will respectively be $\alpha \cdot n\%$ and $\beta \cdot d\%$ with arbitrary values α, β in $(0, 1)$, but the log-ratio $\ln(\alpha \cdot n/\beta \cdot d) = \ln(\alpha/\beta) + \ln(n/d)$ can take any value, as $\ln(\alpha/\beta)$ freely varies between $-\infty$ and $+\infty$ and $\ln(n/d)$ has a fixed value (in our case, it is 0, as $n = d = 0.01$). Thus, in the absence of information about the probability law of the composition, the log-ratio of two rounded zeroes is a missing-at-random value.

3 Regression Analysis for Compositions

3.1 The Fully-Observed Case

In compositional regression, the goal is to predict a compositional random vector \mathbf{Z} as a (linear) function of a real variable. Taking (4) and adding the residuals $\epsilon_n = [\epsilon_{n1}, \dots, \epsilon_{nD}]$, the linear regression model becomes

$$\mathbf{z}_n = \mathbf{b}_0 \oplus x_n \odot \mathbf{b}_1 \oplus \epsilon_n.$$

Following Daunis-i-Estadella et al. (2002), least-squares estimates of the compositional parameters are obtained minimizing the expected squared Aitchison distance (3) between the observations \mathbf{z}_n and the predictions $\hat{\mathbf{z}}_n = \mathbf{b}_0 \oplus x_n \odot \mathbf{b}_1$,

$$\begin{aligned} \hat{L} &= \frac{1}{D} \sum_{i=1}^{D-1} \sum_{j>i}^D \left[\frac{1}{N} \sum_{n=1}^N \left(\ln \frac{z_{ni}}{z_{nj}} - \ln \frac{b_{0i}}{b_{0j}} - x_n \cdot \ln \frac{b_{1i}}{b_{1j}} \right)^2 \right] \\ &= \frac{1}{D} \sum_{i=1}^{D-1} \sum_{j>i}^D \left[\frac{1}{N} \sum_{n=1}^N (y_{n(ij)} - \beta_{0(ij)} - x_n \cdot \beta_{1(ij)})^2 \right], \end{aligned} \tag{5}$$

where $y_{n(ij)} = \ln(z_{ni}/z_{nj})$, $\beta_{0(ij)} = \ln(b_{0i}/b_{0j})$ and $\beta_{1(ij)} = \ln(b_{1i}/b_{1j})$. From here on, the pair $(ij) \equiv k$ behaves like a single index from 1 to $K = D(D - 1)/2$, running along all the pairs (i, j) such that $1 \leq i < j \leq D$. Being equivalent, we use either the notation k or (ij) .

Given that (5) is a sum of positive terms of the form

$$\hat{L}_k = \sum_n (y_{nk} - \beta_{0k} - x_n \cdot \beta_{1k})^2,$$

Daunis-i-Estadella et al. (2002) suggest to minimize the empirical \hat{L} by minimizing each \hat{L}_k , i.e., to obtain estimates $\hat{\beta}_{0k}$ and $\hat{\beta}_{1k}$ by using standard regression for each \hat{L}_k , and retrieve from these estimates compatible values of $\hat{\mathbf{b}}_0$ and $\hat{\mathbf{b}}_1$. The authors state that the same solution is obtained by working just with the set of $D - 1$ log-ratios taking $j = D$, with less computational effort and no need of special software. But for the case presented in this contribution, it is instrumental to show when these compatible values exist, and how to obtain them.

Denote with $\mathbf{\Delta}$ a matrix of $D \times K$ elements, where the columns give all possible pairwise differences. For instance, in the case study this (5×10) -element matrix is

$$\mathbf{\Delta} = \begin{pmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & -1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & -1 & 0 & 0 & -1 & 0 & -1 & 0 & 1 \\ 0 & 0 & 0 & -1 & 0 & 0 & -1 & 0 & -1 & -1 \end{pmatrix}. \tag{6}$$

The Moore–Penrose generalized inverse of $\mathbf{\Delta}$, denoted as $\mathbf{\Delta}^-$, is easily shown to be

$$\mathbf{\Delta}^- = \frac{1}{D} \mathbf{\Delta}^t. \tag{7}$$

The appendix contains a proof of this equality, and a brief summary of those properties of the Moore–Penrose inverse most useful for this contribution. These two matrices allow us to relate the compositions and their set of K log-ratios, through

$$\begin{aligned} \mathbf{y}_n &= \ln \mathbf{z}_n \cdot \mathbf{\Delta} = \text{clr}(\mathbf{z}_n) \cdot \mathbf{\Delta}, \\ \mathbf{z}_n &= \mathcal{C}(\exp(\mathbf{y}_n \cdot \mathbf{\Delta}^-)) = \text{clr}^{-1}(\mathbf{y}_n \cdot \mathbf{\Delta}^-). \end{aligned}$$

This inversion holds exactly, because $\mathbf{\Delta}$ and $\mathbf{\Delta}^-$ are linear transformations with rank $D - 1$, the dimension of \mathcal{S}^D .

Finally, regarding the regression parameters, if all log-ratios are observed, the set of estimates $\{\hat{\beta}_{0(ij)}, \hat{\beta}_{1(ij)}\}$ are compatible with two compositions $\hat{\mathbf{b}}_0$ and $\hat{\mathbf{b}}_1$, through

$$\begin{aligned} \boldsymbol{\beta}_a &= \ln \mathbf{b}_a \cdot \mathbf{\Delta} = \text{clr}(\mathbf{b}_a) \cdot \mathbf{\Delta}, \\ \mathbf{b}_a &= \mathcal{C}(\exp(\boldsymbol{\beta}_a \cdot \mathbf{\Delta}^-)) = \text{clr}^{-1}(\boldsymbol{\beta}_a \cdot \mathbf{\Delta}^-), \end{aligned} \tag{8}$$

for $a = 0$ or $a = 1$.

Sketch of a proof: If one stacks the two vectors $\hat{\mathbf{b}}_0$ and $\hat{\mathbf{b}}_1$ together in a $2 \times D$ matrix \mathbf{B} , and defines the $N \times 2$ matrix \mathbf{X} containing a column of N ones and a column with all observed values of the predictand variable x , one can compute $\text{clr}(\mathbf{B}) = (\mathbf{X}^t \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^{-1} \cdot [\text{clr}(\mathbf{Z})]$. The same can be done ordering the estimates $\{\hat{\beta}_{a(ij)}\}$ in a $2 \times K$ matrix $\boldsymbol{\beta}$, which can be estimated as $\boldsymbol{\beta} = (\mathbf{X}^t \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^{-1} \cdot [\text{clr}(\mathbf{Z}) \cdot \mathbf{\Delta}]$. The equalities sought after are immediately obtained from these two equations.

3.2 The Undersampled Case

In the case that some parts are lost (either because the variable is a structural zero, it is below the detection limit or is missing at random), some error terms in (5) will not be computable. This may affect some log-ratios preferentially, invalidating the direct estimation with Moore–Penrose pseudo-inverses. This section presents a way to compensate for the different “sample size” of each involved log-ratio. We follow the same steps proposed by Daunis-i-Estadella et al. (2002).

1. Solve the regression for each log-ratio independently. Here, one obtains estimates of the linear regression coefficients $\{\hat{\beta}_{0(ij)}, \hat{\beta}_{1(ij)}\}$, with standard regression,

$$\hat{\beta}_{1(ij)} = \frac{s_{xy(ij)}}{s_x^2(ij)}, \quad \text{and} \quad \hat{\beta}_{0(ij)} = \bar{y}_{(ij)} - \hat{\beta}_{1(ij)} \cdot \bar{x}_{(ij)},$$

where $s_{xy(ij)}$ is the covariance of $x = \phi$ and $y = \log(z_i/z_j)$, $\bar{x}_{(ij)}$ and $s_x^2(ij)$ are the mean and variance of ϕ for the subset of samples where the log-ratio y is observed, and $\bar{y}_{(ij)}$ is the average of this log-ratio. It is important to remind here that, for each (ij) , the data set used in the computation of these statistics is different, even for the explanatory variable x . As explained before, the following situations arise: samples where parts i and j are observed give an observed datum; samples where neither part i nor j are observed behave like a missing-at-random value; it is thus

- expected that such values should not induce bias in the estimation, but just increase its uncertainty (decrease its reliability); samples where either part i or j are not observed are truncated values, not missing at random. Ignoring this information could yield some bias, in particular when the proportion of such samples is high.
- Combine the partial estimates. The goal is now to combine these estimates $\{\hat{\beta}_{a(ij)}\}$, accounting for their different reliability. Taking (8) as approximately fulfilled, one could write

$$\hat{\beta}_{a(ij)} = \log \frac{\hat{b}_{ai}}{\hat{b}_{aj}} + \Delta\beta_{a(ij)} = (\text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta})_{(ij)} + \Delta\beta_{a(ij)},$$

where $\Delta\beta_{a(ij)}$ would be the residual error of the approximation. Summing these residuals errors, squared and as many times as they are “observed”, one obtains a measure of global discrepancy between the $\beta_{a(ij)}$ estimates of the first step and the global $\hat{\mathbf{b}}_a$ sought

$$\begin{aligned} \epsilon &= \sum_{(ij)} N_{(ij)} \Delta\beta_{a(ij)}^2 \\ &= (\hat{\beta}_a - \text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta}) \cdot \text{diag}(\mathbf{N}) \cdot (\hat{\beta}_a - \text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta})^t, \end{aligned} \tag{9}$$

where $\text{diag}(\mathbf{N})$ is a diagonal matrix with elements $\{N_{(ij)}\}$, the number of times the log-ratio (i, j) is observed. Minimizing discrepancy (9), by differentiation with respect to $\hat{\mathbf{b}}_a$ and equating to zero, one gets

$$\mathbf{0} = \frac{\partial \epsilon}{\partial \hat{\mathbf{b}}_a} = (\hat{\beta}_a - \text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta}) \cdot \text{diag}(\mathbf{N}) + \text{diag}(\mathbf{N}) \cdot (\hat{\beta}_a - \text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta})^t,$$

which, given that $\text{diag}(\mathbf{N})$ is diagonal, is simplified to

$$\begin{aligned} \mathbf{0} &= (\hat{\beta}_a - \text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta}) \cdot \text{diag}(\mathbf{N}), \\ \hat{\beta}_a \cdot \text{diag}(\mathbf{N}) &= \text{clr}(\hat{\mathbf{b}}_a) \cdot \mathbf{\Delta} \cdot \text{diag}(\mathbf{N}). \end{aligned}$$

By taking $\mathbf{\Delta}_N = \mathbf{\Delta} \cdot \text{diag}(\mathbf{N})$, this last equation can be rewritten as

$$\hat{\mathbf{b}}_a = \text{clr}^{-1}((\beta_a \cdot \text{diag}(\mathbf{N})) \cdot \mathbf{\Delta}_N^-). \tag{10}$$

Note that $\mathbf{\Delta}_N$ is the matrix of differences (6), where each column is multiplied by the number of times that difference was observed. In this sense, (10) as an estimator of $\hat{\mathbf{b}}_a$ is a kind of weighted average, where the value $\hat{\beta}_{a(ij)}$ observed $N_{(ij)}$ times has a contribution of $+N_{(ij)}$ to the estimation of $\log(b_{ai})$ and of $-N_{(ij)}$ to the estimation of $\log(b_{aj})$.

4 Practical Case

The method proposed before is here applied to the Blatt et al. (1972) data set.

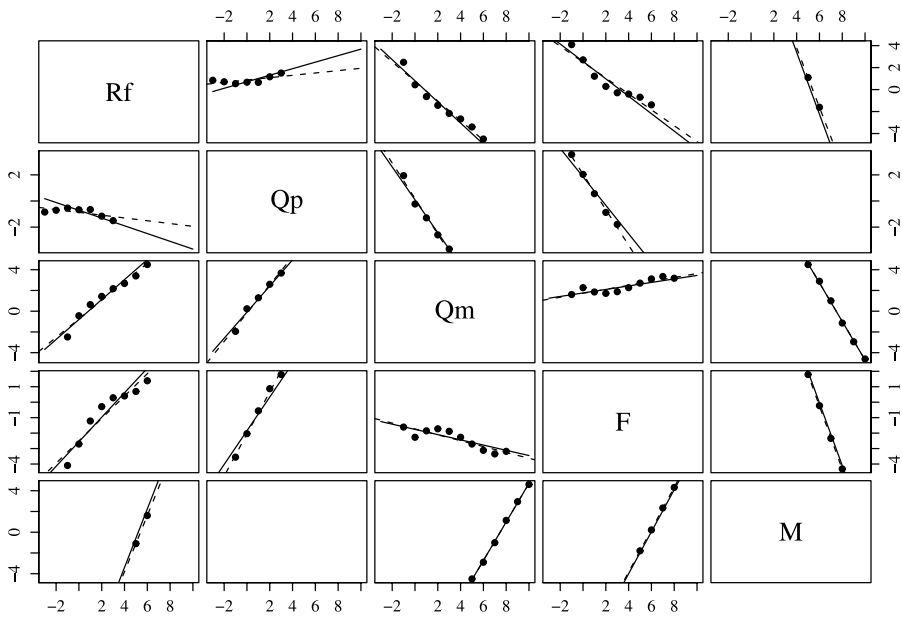


Fig. 2 Scatter plots of all observed log-ratios against grain size, with indication of each partial regression line (*dashed line*) and the final model obtained with the proposed procedure (*solid line*). The (i, j) plot represents $\log(z_i/z_j)$ against ϕ . Note that all plots in a row share the same vertical scale

Table 2 Number of points where each pair of parts was simultaneously observed (left), and number of points where only one part of the log-ratio was available, out of the 14 samples of Table 1: the differences between 14 and the sum of these two tables is the number of missing-at-random values. The upper triangle of the left table will be used afterwards

	Fully observed data					Censored data				
	R_f	Q_p	Q_m	F	M	R_f	Q_p	Q_m	F	M
R_f	10	7	8	8	2	0	3	6	4	12
Q_p	7	7	5	5	0	3	0	9	7	0
Q_m	8	5	12	10	6	6	9	0	2	6
F	8	5	10	10	4	4	9	2	0	8
M	2	0	6	4	6	12	0	6	8	0

1. Solve the regression for each log-ratio independently. Results are shown in Fig. 2, containing $D(D - 1) = 20$ plots: each one is a scatter plot of the (i, j) -log-ratio against ϕ , where i is the row and j is the column. Note that the plot is “antisymmetric”, e.g., the y axes of figures (1, 2) and (2, 1) are equal, but with changed signs. Note also that the ratio Q_p/M is never observed, and the corresponding plot is empty. Recall that each plot is obtained with a different set of data: Table 2 gives the number of available data used in each regression, and reports also the number of censored values.

Fig. 3 Final regression parameter estimates expressed as compositions and plotted as a bar plot

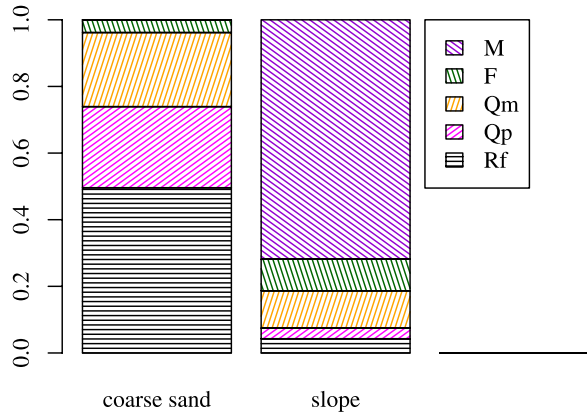


Table 3 Final regression parameter estimates expressed as compositions and as enrichment/depletion factors with respect to quartz. The table also reports the predicted composition of an hypothetical rock block following the linear trend at $\phi \approx -2$

	Composition of		Slope ($\hat{\mathbf{b}}_1$)	Multiplicative factors	
	Rock block	Coarse sand ($\hat{\mathbf{b}}_0$)		Enrichment	Depletion
R_f	4.49×10^{-1}	4.96×10^{-1}	0.04	0.38	2.62
Q_p	5.37×10^{-1}	2.43×10^{-1}	0.03	0.28	3.52
Q_m	1.12×10^{-2}	2.22×10^{-1}	0.11	1	1
F	3.21×10^{-3}	3.83×10^{-2}	0.09	0.84	1.18
M	4.04×10^{-11}	2.10×10^{-7}	0.72	6.41	0.16

2. Combine the partial estimates. One needs the matrix $\text{diag}(\mathbf{N})$, containing the number of samples observed for each log-ratio, which can be extracted from Table 2, its upper triangle, without the diagonal elements. Multiplying then the columns of matrix $\mathbf{\Delta}$ by these numbers, one obtains the weighted difference matrix

$$\mathbf{\Delta}_N = \begin{pmatrix} 7 & 8 & 8 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ -7 & 0 & 0 & 0 & 5 & 5 & 0 & 0 & 0 & 0 \\ 0 & -8 & 0 & 0 & -5 & 0 & 0 & 10 & 6 & 0 \\ 0 & 0 & -8 & 0 & 0 & -5 & 0 & -10 & 0 & 4 \\ 0 & 0 & 0 & -2 & 0 & 0 & -0 & 0 & -6 & -4 \end{pmatrix}.$$

Finally, computing its Moore–Penrose generalized inverse (which has no nice form) and plugging into (10), we get the parameter estimates $\hat{\mathbf{b}}_0$ and $\hat{\mathbf{b}}_1$ with least discrepancy.

Resulting estimates are included in Fig. 3 and Table 3. The compositional linear process is thus described with an “original” coarse sand composition (for $\phi = 0$, the intercept at the origin, which is the threshold between very coarse and coarse sand fractions) of approximately 50% rock fragments, 45% quartz (evenly distributed among mono- and poly-crystalline grains) and 5% feldspar, and a trend towards finer

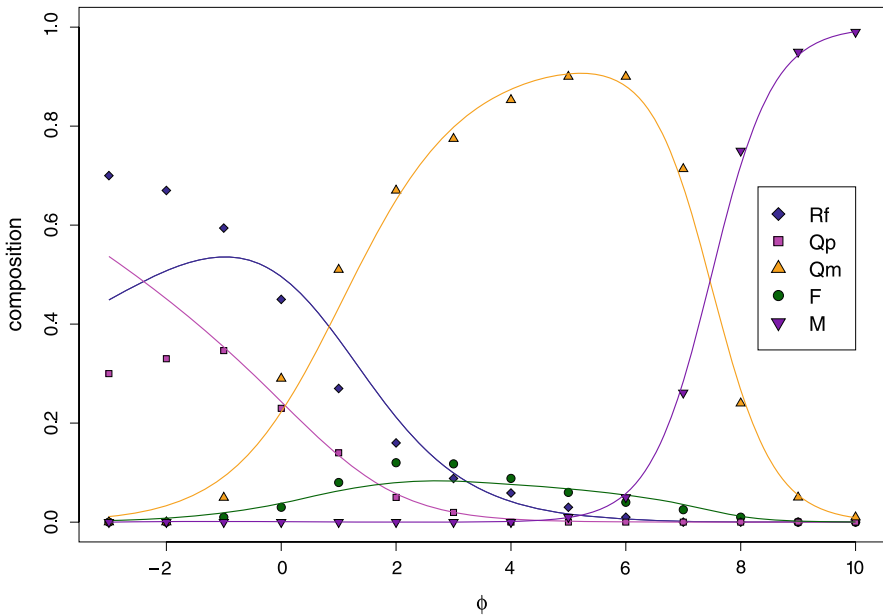


Fig. 4 Predicted and observed compositions, as a diagram of evolution with ϕ

grain sizes of a very steep enrichment on mica (6.5 times that of monocrystalline quartz enrichment), and depletion of all the other grain types (2.5 times stronger for rock fragments, 3.5 times for polycrystalline quartz, and for feldspar a slightly quicker depletion than that of monocrystalline quartz).

The obtained parameters can be transformed back to log-ratios with (8). This yields some $\hat{\beta}_{0(ij)}$ and $\hat{\beta}_{1(ij)}$ different from the original ones, but mutually compatible. The lines associated with these compatible parameters are also included in Fig. 2. The agreement between original and compatible lines is very high, with the exception of the ratio R_f/Q_p . This ratio shows an extreme around $\phi = 0$, whereas the original regression looks quite flat around the average of the log-ratio (and does not capture in any way its non-linear behavior), the compatible one closely follows the decreasing tail of the trend (for $\phi \geq 0$).

Given that regression goodness-of-fit tests are in this case rather meaningless (as they depend on the sample size, i.e., the number of points we decided to use in the digitalization), the least one must do is visually assess that results fit the data. This is shown in Fig. 4, where one clearly sees that most features of the original plot are captured: monocrystalline grains (Q_m , F and M) are perfectly described by the model, whereas polycrystalline grains (Q_p and R_f) are only adequately modeled for $-1 < \phi$. Finally, one can look at several subcompositions, in order to assess the degree to which the relative information of the data set is adequately captured by the model. Figure 2 already showed all possible 2-part compositions. Figure 5 reports the 3-part subcompositions in ternary diagrams where both the data set and the model are plotted. Data and model fit very well, even for those diagrams with zeroes, where samples visually fall on the border.

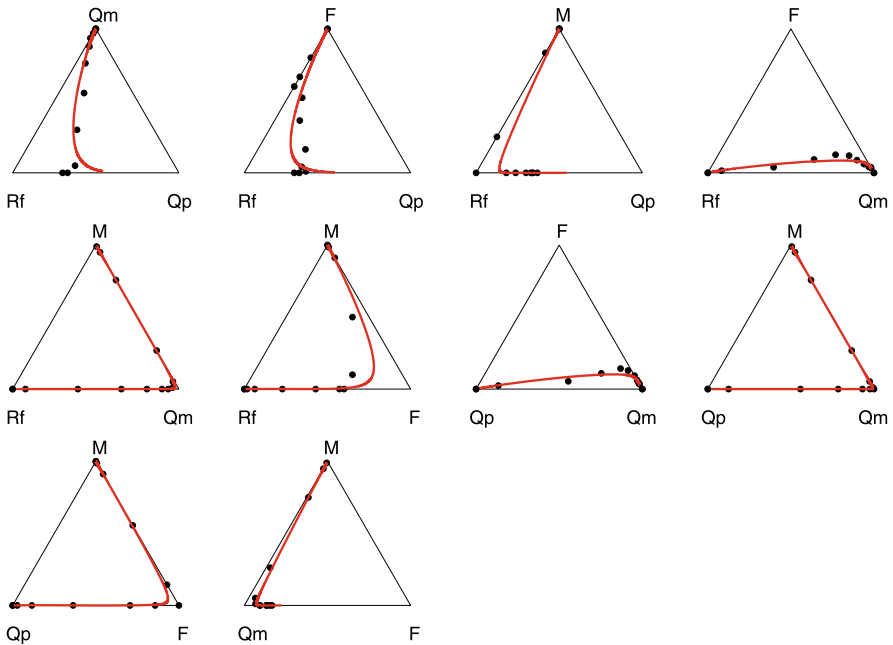


Fig. 5 Predicted and observed compositions, displayed in all possible ternary diagrams

5 Discussion

The rather complex plot in Fig. 1 can be described by a compositional linear process, with an intercept composition (for $\phi = 0$) of approximately 50% rock fragments, 45% quartz (evenly distributed among mono- and polycrystalline grains) and 5% feldspar, and a trend of strong enrichment in mica towards finer grain sizes, and quick depletion of polycrystalline quartz and rock fragments. Feldspar grains are slightly less stable than monocrystalline quartz grains, and polycrystalline grains of quartz decrease slightly faster than rock fragments. Assuming a global average sand size of $\phi = 1.5-2$, the average sand due to our model has a composition of $R_f = 29-20$, $Q_p = 9-6$, $Q_m = 55-65$, and $F = 7-8$ (in percentages). Previous estimates of global average sand composition are in the range of $R_f = 23-24$, $Q_t =$ (total quartz) $63-68$ and $F = 10-12$. These are based on (i) a rough compilation of the frequency of common sandstone types worldwide and their average composition (Pettijohn et al. 1987), and (ii) a more substantiated study of the average sand of South America (Potter 1994). These estimates and our model are in very good agreement, though the feldspar content is slightly lower in our model (and thus, the ratio F/Q). This is an effect of linear regression, which does not capture the small-scale variations on the ratios including F (Fig. 2).

The agreement between model and “data” is very good, specially for sands and finer sediments ($\phi \geq -1$), and far beyond what might be obtained with a raw, classical regression. The quality of this agreement is surprising, given that the plot is in fact a gathering of disperse information, mostly semi-quantitative or qualitative.

Therefore, we conclude that the Blatt et al. (1972) understanding of the processes controlling the relationship between sediment grain type and grain size can be perfectly described by a one-dimensional compositional linear process (in the Aitchison geometry), for sands and finer sediments. Keeping this in mind, one can interpret the intercept composition as a global average “original” coarse to very coarse sand (for $\phi = 0$), and the slope composition as the resultant of the relative influences of all processes controlling sediment composition. One must then conclude that the resultant of these effects acting on the sediment body is constant along the grain size spectrum, as this is the implication of a linear process. This does not necessarily contradict the classical idea that chemical weathering is more intense in finer fractions than in coarser ones, when compared with physical alteration processes (crushing, comminution, abrasion), as these processes evolve with time, whereas our linear model describes the time-independent changes of petrographic composition when changing grain size. Further research is nevertheless necessary to better understand this issue.

For coarse sediments ($\phi < -1$), the model does not adequately capture the observed distribution between rock fragments and polycrystalline quartz, the only two grain types (typically) present in these size fractions. The distribution observed seems mainly inherited from the source rock characteristics, as rock fragments purely made of quartz become less common for larger blocks. Moreover, the very difference between these two grain classes is quite arbitrary for grains of this size. In summary, these arguments suggest that the trend can be safely used as a baseline for the petrographic composition of sands and finer sediments, and should not be extrapolated to coarser grain fractions.

6 Conclusions

Regression with a compositional response can be equivalently obtained from two approaches: either one minimizes the average Aitchison distance between the predicted and the observed compositions (as a function of two parameter compositions playing the role of slope and intercept); or one applies regression to all possible log-ratios and combines the estimated slopes or intercepts, to obtain the parameter slope or intercept composition. If no missing value is present, this contribution shows that both approaches give the same result. In the presence of zeroes, however, only the second approach provides sensible estimates. In this case, only the ratios where both numerator and denominator were observed do actually contribute to the global regression.

For the case study, we conclude that the Blatt et al. (1972) understanding of the processes controlling the relationship between sediment grain type and grain size can be perfectly described by a one-dimensional compositional linear process (in the Aitchison geometry), for sands and finer-grained sediments. This linear trend globally balances all processes acting on the sediment body. The fact that the fitted model is linear implies that the relative influence of all controlling processes is constant along the grain sizes studied here. As this trend is not based on specific data but is aimed to describe a global model, it may be taken as a baseline for testing local case studies against it, and in this case one can split the study in the explanation of the trend itself, and the explanation of the discrepancies found.

Acknowledgements This research was funded by the Department of Universities, Research and Information Society (DURSI, grant 2005 BP-A 10116) of the *Generalitat de Catalunya*, and the German Research Foundation (DFG, grant EY23/11-1). The authors acknowledge fruitful discussion with G. van den Boogaart and the Girona Group on Compositional Data Analysis, in particular with J. Daunis-i-Estadella, J.J. Egozcue, and the reviews of V. Pawlowsky-Glahn and A. Buccianti.

Open Access This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

Appendix A: The Moore–Penrose Generalized Inverse

Definition 1 For any matrix \mathbf{A} of N columns and M rows, the Moore–Penrose generalized inverse is the unique matrix \mathbf{A}^- satisfying the following set of pseudo-inversion properties:

1. $\mathbf{A} \cdot \mathbf{A}^- \cdot \mathbf{A} = \mathbf{A}$.
2. $\mathbf{A}^- \cdot \mathbf{A} \cdot \mathbf{A}^- = \mathbf{A}^-$.
3. $(\mathbf{A} \cdot \mathbf{A}^-)^t = \mathbf{A} \cdot \mathbf{A}^-$.
4. $(\mathbf{A}^- \cdot \mathbf{A})^t = \mathbf{A}^- \cdot \mathbf{A}$.

where the superscript t is transposition (or conjugate transposition if the matrix is complex).

Recall that the Moore–Penrose generalized inverse of a matrix may be computed with its singular value decomposition: if $\mathbf{A} = \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{V}^t$, then $\mathbf{A}^- = \mathbf{V} \cdot \mathbf{D}^- \cdot \mathbf{U}^t$, where \mathbf{U} and \mathbf{V} columns form two sets of orthonormal vectors (respectively called left and right singular vectors), \mathbf{D} is a diagonal matrix of non-negative singular values, and \mathbf{D}^- is the generalized inverse of \mathbf{D} , containing in the diagonal the inverse of each non-zero singular value (and a zero for each zero singular values).

Property 1 For any matrix \mathbf{A} and its generalized inverse \mathbf{A}^- , the following properties hold:

1. $(\mathbf{A}^-)^- = \mathbf{A}$ (pseudo-inversion is reversible).
2. $(\alpha \cdot \mathbf{A})^- = \alpha^{-1} \cdot \mathbf{A}^-$.
3. $\mathbf{A} = \mathbf{A}^t \cdot (\mathbf{A} \cdot \mathbf{A}^t)^-$.

This is adapted from Ben-Israel and Greville (2003).

Appendix B: A Proof

Proof of equation (7): The first step of the proof is to show that

$$\mathbf{\Delta} \cdot \mathbf{\Delta}^t = \begin{pmatrix} D-1 & -1 & \dots & -1 \\ -1 & D-1 & \dots & -1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \dots & D-1 \end{pmatrix} = D \cdot \mathbf{I} - \mathbf{1} =: \mathbf{B}, \tag{11}$$

because the pseudo-inverse of \mathbf{B} is easy to compute. Here, \mathbf{I} is the identity, and $\mathbf{1}$ is a $D \times D$ matrix full of ones. We can see that (11) holds by looking at the structure of the rows of $\mathbf{\Delta}$: the element (i, j) of matrix \mathbf{B} is, in fact, the scalar product of two of these rows, $b_{ij} = \langle \mathbf{\Delta}_i, \mathbf{\Delta}_j \rangle$. If $i = j$, then b_{ii} is the sum of as many $(\pm 1)^2 = 1$ as different log-ratios can be built involving the i th part; thus, $b_{ii} = D - 1$ (the log-ratio of a part with itself is not useful; and if we counted a given log-ratio, the log of its inverse ratio is not counted). If $i \neq j$, then $b_{ij} = -1$, given that two parts only coincide in one single log-ratio, and they have always different signs, being one in the numerator and the other in the denominator.

Now define the matrix $\mathbf{C} = (1/D) \cdot \mathbf{B}$. The second step is to show that it is a pseudo-identity

$$\begin{aligned} \mathbf{C} \cdot \mathbf{C} &= \left(\mathbf{I} - \frac{1}{D} \mathbf{1} \right) \cdot \left(\mathbf{I} - \frac{1}{D} \mathbf{1} \right) = \mathbf{I} \cdot \mathbf{I} - \frac{1}{D} \mathbf{I} \cdot \mathbf{1} - \frac{1}{D} \mathbf{1} \cdot \mathbf{I} + \frac{1}{D^2} \mathbf{1} \cdot \mathbf{1} \\ &= \mathbf{I} - \frac{2}{D} \mathbf{1} + \frac{1}{D^2} D \mathbf{1} = \mathbf{I} - \frac{1}{D} \mathbf{1} = \mathbf{C}. \end{aligned}$$

Now, by simple recurrence of these calculations and attending to the definition of generalized inverse, one easily deduces that $\mathbf{C}^- = \mathbf{C}$.

The third step is to show that $\mathbf{B}^- = (1/D) \cdot \mathbf{C}$, since then one will be able to apply Property 1.3. But, given that \mathbf{C} is a pseudo-identity, this is immediate thanks to Property 1.2

$$\mathbf{B}^- = (D \cdot \mathbf{C})^- = \frac{1}{D} \cdot \mathbf{C}^- = \frac{1}{D} \cdot \mathbf{C}.$$

The last step is the application of Property 1.3, giving

$$\mathbf{\Delta}^- = \mathbf{\Delta}^t \cdot \mathbf{B}^- = \mathbf{\Delta}^t \cdot \frac{1}{D} \cdot \mathbf{C} = \frac{1}{D} \cdot \mathbf{\Delta}^t \left(\mathbf{I} - \frac{1}{D} \mathbf{1} \right) = \frac{1}{D} \cdot \mathbf{\Delta}^t - \frac{1}{D^2} \cdot \mathbf{\Delta}^t \cdot \mathbf{1}.$$

This yields (7), by realizing that the term $\mathbf{\Delta}^t \cdot \mathbf{1}$ is a matrix full of the sums of the elements in given columns of $\mathbf{\Delta}$, which are identically zero.

References

Aitchison J (1986) The statistical analysis of compositional data. Monographs on statistics and applied probability. Chapman & Hall, London. (Reprinted in 2003 with additional material by The Blackburn Press)

Aitchison J (1997) The one-hour course in compositional data analysis or compositional data analysis is simple. In: Pawlowsky-Glahn V (ed) Proceedings of IAMG'97—The third annual conference of the International Association for Mathematical Geology, vol 1. International Center for Numerical Methods in Engineering (CIMNE), Barcelona, pp 3–35

Aitchison J (2002) Simplicial inference. In: Viana MAG, Richards DSP (eds) Algebraic methods in statistics and probability. Contemporary mathematics series, vol 287. American Mathematical Society, Providence, pp 1–22

Arribas J, Tortosa A (2003) Detrital modes in sedimenticlastic sands from low-order streams in the Iberian Range, Spain: the potential for sand generation by different sedimentary rocks. *Sedimentary Geol* 159:275–303

- Ben-Israel A, Greville TNE (2003) *Generalized inverses: theory and applications*, 2nd edn. Springer, New York
- Billheimer D, Guttorp P, Fagan W (2001) Statistical interpretation of species composition. *J Am Stat Assoc* 96(456):1205–1214
- Blatt H, Middleton GV, Murray RC (1972) *Origin of sedimentary rocks*. Prentice-Hall, Englewood Cliffs
- Buccianti A, Mateu-Figueras G, Pawlowsky-Glahn V (eds) (2006) *Compositional data analysis: from theory to practice*. Special publication, vol 264. The Geological Society, London
- Chandrajith R, Dissanayake CB, Tobschall HJ (2001) Application of multi-element relationships in stream sediments to mineral exploration: a case study of Walawe Ganga Basin, Sri Lanka. *Appl Geochem* 16:339–350
- Dauinis-i-Estadella J, Egozcue JJ, Pawlowsky-Glahn VV (2002) Least squares regression in the simplex. In: Bayer U, Burger H, Skala W (eds) *Proceedings of IAMG'02—The eighth annual conference of the International Association for Mathematical Geology*. Selbstverlag der Alfred-Wegener-Stiftung, Berlin, pp 411–416
- Egozcue JJ, Pawlowsky-Glahn V (2006) Simplicial geometry for compositional data. See Buccianti, Mateu-Figueras, and Glahn (2006), pp 145–160
- Grantham JH, Velbel MA (1988) The influence of climate and topography on rock-fragment abundance in modern fluvial sands of the southern Blue Ridge Mountains, North Carolina. *J Sedimentary Res* 58:219–227
- Johnsson M (1993) The system controlling the composition of clastic sediments. *Geol Soc Am Spec Paper* 284:1–19
- Kiminami K, Fujii K (2007) The relationship between major element concentration and grain size within sandstones from four turbidite sequences in Japan. *Sedimentary Geol* 195:203–215
- Lim DI, Jung HS, Choi JY, Yang S, Ahn KS (2006) Geochemical compositions of river and shelf sediments in the Yellow Sea: grain-size normalization and sediment provenance. *Cont Shelf Res* 26:15–24
- Martín-Fernández JA (2001) *Medidas de diferencia y clasificación no paramétrica de datos composicionales (Measures of difference and non-parametric classification of compositional data)*. PhD thesis, Universitat Politècnica de Catalunya, Barcelona
- Martín-Fernández JA, Barceló-Vidal C, Pawlowsky-Glahn V (2000) Zero replacement in compositional data sets. In: Kiers H, Rasson J, Groenen P, Shader M (eds) *Studies in classification, data analysis, and knowledge organization (Proceedings of the 7th conference of the International Federation of Classification Societies (IFCS'2000), University of Namur, Namur, 11–14 July, Springer, Berlin*, pp 155–160
- Murray GD (1979) The estimation of multivariate normal density functions using incomplete data. *Biometrika* 66:375–380
- Nesbitt H, Markovics G (1997) Weathering of granodioritic crust, long-term storage of elements in weathering profiles and petrogenesis of siliciclastic sediments. *Geochim Cosmochim Acta* 61:1653–1670
- Nesbitt H, Young G (1984) Prediction of some weathering trends of plutonic and volcanic rocks based on thermodynamic and kinetic considerations. *Geochim Cosmochim Acta* 41:1523–1534
- Nesbitt H, Young G (1996) Petrogenesis of sediments in the absence of chemical weathering: Effects of abrasion and sorting on bulk composition and mineralogy. *Sedimentology* 43:341–358
- Noda A (2005) Texture and petrology of modern river, beach and shelf sands in a volcanic back-arc setting, northeastern Japan. *Isl Arc* 14:687–707
- Palarea-Albaladejo J, Martín-Fernández JA, Gómez-García J (2007) A parametric approach for dealing with compositional rounded zeros. *Math Geol* 359(7):625–645
- Palomares M, Arribas J (1993) Modern stream sands from compound crystalline sources: Composition and sand generation index. *Geol Soc Am Spec Paper* 284:313–322
- Pawlowsky-Glahn V, Egozcue JJ (2001) Geometric approach to statistical analysis on the simplex. *Stoch Environ Res Risk Assess* 15(5):384–398
- Pettijohn F (1957) *Sedimentary rocks*. Harper, New York
- Pettijohn F, Potter P, Siever R (1987) *Sand and sandstone*, 2nd edn. Springer, New York
- Potter P (1994) Modern sands of South America: composition, provenance and global significance. *Int J Earth Sci* 83(1):212–232
- Rubin DB (1976) Inference and missing data. *Biometrika* 63:592–581
- Solano-Acosta W, Dutta PK (2005) Unexpected trend in the compositional maturity of second-cycle sand. *Sedimentary Geol* 178:275–283
- van den Boogaart KG, Tolosana-Delgado R, Bren M (2006) Concepts for handling zeroes and missing values in compositional data. In: Pirard E, Dassargues A, Havenith HB (eds) *Proceedings of IAMG'06—The XI annual conference of the International Association for Mathematical Geology*. University of Liège, Belgium, CD-ROM

- von Eynatten H (2004) Statistical modelling of compositional trends in sediments. *Sedimentary Geol* 171:79–89
- von Eynatten H, Barceló-Vidal C, Pawlowsky-Glahn V (2003) Modelling compositional change: the example of chemical weathering of granitoid rocks. *Math Geol* 35(3):231–251
- Wentworth C (1922) A scale of grade and class terms for clastic sediments. *J Geol* 30:377–392
- Whitmore GP, Crook KAW, Johnson DP (2004) Grain size control of mineralogy and geochemistry in modern river sediment, New Guinea Collision, Papua New Guinea. *Sedimentary Geol* 171:129–157
- Zhang C, Wang L, Li G, Dong S, Yang J, Wang X (2002) Grain size effect on multi-element concentrations in sediments from the intertidal flats of Bohai Bay, China. *Appl Geochem* 17:59–68