



A novel downscaling procedure for compositional data in the Aitchison geometry with application to soil texture data

Federico Gatti¹ · Alessandra Menafoglio¹ · Niccolò Togni¹ · Luca Bonaventura¹ · Davide Brambilla² · Monica Papini² · Laura Longoni²

Accepted: 7 October 2020 / Published online: 29 October 2020
© The Author(s) 2020

Abstract

In this work, we present a novel downscaling procedure for compositional quantities based on the Aitchison geometry. The method is able to naturally consider compositional constraints, i.e. unit-sum and positivity, accounting for the scale invariance and relative scale of these data. We show that the method can be used in a block sequential Gaussian simulation framework in order to assess the variability of downscaled quantities. Finally, to validate the method, we test it first in an idealized scenario and then apply it for the downscaling of digital soil maps on a more realistic case study. The digital soil maps for the realistic case study are obtained from SoilGrids, a system for automated soil mapping based on state-of-the-art spatial predictions methods.

Keywords Geostatistics · Block sequential Gaussian simulation · Area-to-point kriging · Isometric log-ratios

1 Introduction

Uncertainty Quantification (UQ) is a crucial aspect for numerical tools intended to simulate physical processes, since it is important to provide an extensive analysis of the uncertainty of the outputs related to the variability of the inputs. Classical methods to perform this task are based on Monte Carlo (MC) simulations (Kalos and Whitlock 2009). Here, an ensemble of realizations of the input parameters is used to feed a mathematical/numerical model, aiming to assess the distribution of the response in the face of uncertain inputs. In this broad framework, whenever parameters are characterized by a spatial distribution, geostatistical stochastic simulation can be employed to generate input scenarios for the model (Brown et al. 2002). The geostatistical approach allows one to account for the spatial dependence characterizing the input parameters and to model the spatial structure expected for the realizations

(range of variability, degree of smoothness) through a spatial covariance function. Nonetheless, sound geostatistical simulation needs to take into account the possible constraints of the data, particularly when these represent compositional information. For instance, soil moisture retention plays an important role in models that simulate hydrogeological processes and depends on a number of terrain properties, such as the *soil texture*. The latter in turn is determined by particle-size fractions (psfs), i.e. the relative percentages, in terms of soil composition, of clay, silt and sand, the three categories in which grains of fine earth are divided depending on their size, see e.g. Martín et al. (2017). When some sparse samples are available, geostatistical techniques such as Kriging and conditional Gaussian (co)simulation can be used to interpolate the available observations and assess the associated uncertainty. However, neglecting the inherent characteristics of these data may result in inappropriate results, such as prediction of negative components or modeling spurious correlations (Kim 1999). These serious limitations hinder the use of classical geostatistical methods based on the Euclidean geometry in the presence of compositional data (see, e.g. Aitchison 1982; Buccianti and Grunsky 2014).

In the last decades, an increasing attention has been devoted to developing analytics tools able to account for

✉ Alessandra Menafoglio
alessandra.menafoglio@polimi.it

¹ MOX-Department of Mathematics, Politecnico di Milano, Milan, Italy

² Department of Civil and Environmental Engineering, Politecnico di Milano, Milan, Italy

the features of compositional data, starting from the work of Aitchison (1982). Nowadays, Compositional Data Analysis (CoDa, Egozcue et al. 2003; Pawlowsky-Glahn et al. 2015b) is a well-established area of statistics, which studies models and methods for compositional data, grounded on the Aitchison geometry for the simplex. The Aitchison geometry is based on the foundational idea that, in compositional vectors, only the log-ratios among components represent a meaningful information to be accounted for in the statistical analysis. In a geostatistical setting, this foundational idea led to the development of new kriging methods for compositional data, which were successfully applied in several applied studies (e.g., in the mining context, Pawlowsky-Glahn and Olea 2004; Tolosana-Delgado et al. 2019).

In this work, we focus on the problem of geostatistical downscaling of compositional quantities. This is relevant in applications where no (or limited) direct observation is available within the study area—because of cost or environmental constraints—but low-resolution information is available across the region. This is the case of our motivating study, which focuses on the stochastic characterization of soil texture within a mountain river catchment, aiming to model the hydrogeological instability—and consequent natural hazard—of the region. In this case, no direct observation of particle-size fractions is available, but low-resolution data are reported in public databases, such as SoilGrids (Hengl et al. 2014, 2017). In this case, characterizing the spatial distribution of the soil texture requires to operate a change of support of the available (compositional) information and to assess the corresponding uncertainty.

To the authors' knowledge, none of the available methods for (geostatistical) downscaling allows to account for compositional constraints. For instance, methods of area-to-point kriging and stochastic simulation available in the literature (Kyriakidis 2004) inevitably are subject to the limitations of the Euclidean methods. We here propose an extension of Area-To-Point Regression CoKriging (ATPRCoK)—and associated stochastic simulation—to compositional vectors that, based on the Aitchison geometry, allows to overcome such issues and provide stochastic scenarios for the target compositional parameters.

The remaining part of this work is organized as follows. In Sect. 2 we recall the area-to-point regression (co)kriging method; in Sect. 3, we present the downscaling prediction framework for compositional data; in Sect. 4 we recall the definition of psfs, which are used in Sect. 5 to exemplify and test the features of the method in a first to synthetic case and then in a real scenario. Finally, we apply the method to a case study within the Caldane catchment in the Northern Italy city of Lecco, where we show how the

method is able to provide psfs data at a length-scale most of the time very difficult or impossible to be determined.

2 Area-to-point regression kriging

In this section, we recall the main features of Area-To-Point Regression Kriging (ATPRK) and Area-To-Point Regression Cokriging (ATPRCoK); for further details see e.g. Wang et al. (2015), Xiao et al. (2018). Let us consider a scalar continuous random field $\{Z(\mathbf{x}), \mathbf{x} \in D\}$ defined over a geographical region $D \subset \mathbb{R}^d$. Let us discretize $Z(\mathbf{x})$ as

$$Z_j = \frac{1}{|v_j|} \int_{v_j} Z(\mathbf{x}) \eta(d\mathbf{x}),$$

where Z_j denotes the discretized element at pixel j , v_j defines the geographical support of the j -th pixel having center $\mathbf{x}_j \in D$, $|v_j|$ denotes the measure of the support v_j , and η is a positive measure on D (e.g., the Lebesgue measure). We assume the measure of the pixel support to be equal for all the pixels covering the region D and consider two levels of spatial resolution, one coming from a coarse discretization, denoted by the index $K = 1, \dots, M$, and another coming from a fine discretization, denoted by the index $k = 1, \dots, N$. The measure of the coarse support is a multiple of that of the fine support, s.t. we can define an integer number $P = \frac{|v_K|}{|v_k|}$. Moreover, when using a Euclidean geometry for the data—which is the standard setting for which ATPRK is developed—the low-resolution random field is assumed to be obtained as an arithmetic mean of the high-resolution one, i.e., for $K = 1, \dots, M$,

$$Z_K = \frac{1}{P} \sum_{k: \mathbf{x}_k \in v_K} Z_k. \quad (1)$$

Starting from one complete realization of the low-resolution field Z_K , we want to estimate the high-resolution field Z_k , i.e. perform downscaling. ATPRK allows one to compute an estimate of the field Z_k as a (linear) combination of two parts: regression and Area-To-Point Kriging (ATPK, see e.g. Atkinson 2013; Goovaerts 2008; Kyriakidis 2004; Kyriakidis and Yoo 2005). It uses a linear regression model on a set of covariates for the mean term of Z_k , and kriging to interpolate the residuals from the regression model. The ATPRK predictor \hat{Z}_k of the field Z_k at a given fine scale pixel v_k is defined as

$$\hat{Z}_k = \sum_l \beta^l u_k^l + \sum_{\bar{K}} \lambda^{\bar{K}} e_{\bar{K}}, \quad (2)$$

with $\beta^l, l = 1, \dots, L$ and $\lambda^{\bar{K}}, \bar{K} = 1, \dots, M$ unknown real quantities. The first sum in (2) is a classical linear regression term, describing the mean of Z_k

$$\mathbb{E}[Z_k] = \sum_l \beta^l u_k^l, \tag{3}$$

where $u_k^l, l = 1, \dots, L$ and $k = 1, \dots, N$ are a set of known fine-resolution regressors. Given a realization of the field Z_K , one may linearly upscale Eq. (3) to obtain a linear regression model for Z_K , i.e.,

$$\mathbb{E}[Z_K] = \sum_l \beta^l u_K^l, \tag{4}$$

where u_K^l are the upscaled regressors, i.e.

$$u_K^l = \frac{1}{P} \sum_{k: \mathbf{x}_k \in \mathcal{V}_K} u_k^l. \tag{5}$$

By combining Eqs. (1)–(4)–(5), the regression coefficients β^l can be thus estimated by using a standard fitting procedure (e.g. Ordinary Least Squares (OLS) method) on the low-resolution field Z_K , see e.g. Hengl et al. (2007), Minasny and Mcbratney (2007).

The second term in Eq. (2) is the Area-To-Point-Kriging (ATPK) term. It is the best linear unbiased predictor from the coarse residuals e_K , defined as $e_K = Z_K - \mathbb{E}[Z_K]$. In ATPK, the residual at a given fine pixel k is predicted as the best linear combination of the coarse residuals, subject to unbiasedness, i.e., $\hat{e}_k = \sum_{\bar{K}} \lambda^{\bar{K}} e_{\bar{K}}$, where \hat{e}_k is the fine resolution predicted residual and $\lambda^{\bar{K}}$ solve

$$\min_{\lambda^{\bar{K}} \in \mathbb{R}} \mathbb{E}[(\hat{e}_k - e_k)^2] \quad \text{s.t.} \quad \mathbb{E}[\hat{e}_k] = \mathbb{E}[e_k]. \tag{6}$$

In practice, the ATPK predictor is often computed from a subset of $\bar{M} < M$ of residuals (typically selected in a neighborhood of the target pixel), to reduce the computational burden of the procedure. The optimal weights $\lambda = (\lambda_1, \dots, \lambda_{\bar{M}})'$ are computed by minimizing the prediction error variance, which yields the following kriging linear system

$$\begin{bmatrix} \Sigma & \mathbf{1} \\ \mathbf{1}^T & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} \sigma \\ 1 \end{bmatrix}. \tag{7}$$

Here, the element in position (\bar{K}_1, \bar{K}_2) of matrix Σ is the block-block covariance between the residuals at the coarse pixels centered at $\mathbf{x}_{\bar{K}_1}$ and $\mathbf{x}_{\bar{K}_2}$; the \bar{K} -th element of σ is the point-block covariance of the residuals between fine and coarse pixels respectively centered at \mathbf{x}_k and $\mathbf{x}_{\bar{K}}$, and μ is a Lagrange multiplier.

Note that, in practice, neither the residuals nor their covariance are observed, but need to be estimated from the data. Residuals are typically estimated by difference from

the estimated regression term. Estimating the covariance structure is more critical. Under the assumption that the residual field e_k is stationary and isotropic and denoting with C_{k_1, k_2} the covariance between pixels k_1 and k_2 of the residual at fine scale, we can compute the block-block covariance as

$$\Sigma_{\bar{K}_1, \bar{K}_2} = \frac{1}{P^2} \sum_{i=1}^P \sum_{j=1}^P C_{i,j}, \quad \mathbf{x}_i \in \mathcal{V}_{\bar{K}_1}, \mathbf{x}_j \in \mathcal{V}_{\bar{K}_2}. \tag{8}$$

The point-block covariance is then given by

$$\sigma^{\bar{K}} = \frac{1}{P} \sum_{i=1}^P C_{i,k}, \quad \mathbf{x}_i \in \mathcal{V}_{\bar{K}}. \tag{9}$$

The critical point of the ATPK method is thus the determination of the covariance structure at the fine scale, which cannot be directly estimated, as the data are given at the coarse scale only. The problem of estimating the fine-scale semi-variogram γ_{k_1, k_2}

$$\gamma_{k_1, k_2} = C_{k_1, k_1} - C_{k_1, k_2} = \frac{1}{2} \mathbb{E}[(e_{k_2} - e_{k_1})^2],$$

from coarse-scale data is known as a *deconvolution problem*. We shall not focus on this problem in this work, as it is completely analogous to that in the Euclidean setting—for which a number of methods are available. We refer to Goovaerts (2008) for more details on the deconvolution method used in this work.

In case of multi-dimensional random fields, the ATPRK framework changes slightly in order to take into account possible cross-correlations among field components. Generalization of ATPRK to the multivariate case is analogous to cokriging, and yields Area-To-Point-Regression-CoKriging (ATPRCoK), see e.g. Xiao et al. (2018). In ATPRCoK the coarse residuals appearing in (2) are replaced by the residuals of all the components of the multi-dimensional random field, in order to consider possible cross-correlations among their components. If we consider a p -dimensional random field $\{\mathbf{Z}(\mathbf{x}), \mathbf{x} \in D\}$, its ATPRCoK discrete prediction is,

$$\hat{\mathbf{Z}}_k = \sum_l u_k^l \beta^l + \sum_{\bar{K}} A_{\bar{K}} \mathbf{e}_{\bar{K}},$$

where $\beta^l \in \mathbb{R}^p$ are the vectors of the unknown regression coefficients. The matrix $A_{\bar{K}} \in \mathbb{R}^{p \times p}$ contains the unknown cokriging coefficients and $\mathbf{e}_{\bar{K}} \in \mathbb{R}^p$ is the vector of the coarse scale residuals. The optimal weights $A_{\bar{K}}$ are found by solving a system analogous to (7), but considering covariances and cross-covariances within/among fields components, as in a standard cokriging setting.

3 Compositional ATPRCoK

In this section, we consider the problem of downscaling compositional data and we propose a method which extends ATPRCoK to the Aitchison geometry, and naturally takes into account the compositional nature of the data.

3.1 Compositional data in the Aitchison simplex

A compositional data point $\mathbf{Z} = (Z_1, \dots, Z_p)$ is typically represented as a vector whose elements are proportions (or percentages) of a whole, named *total*. In this case, compositional vectors are characterized by the unit-sum constraint $\sum_i Z_i = 1$, where we denote with $Z_i \geq 0$ the i -th component of compositional data point. More generally, compositional vectors are data which convey relative information, being subject to a constant-sum constraint. Here the total is typically of no interest for the analysis, in the sense that expressing the data w.r.t. a different total (i.e., in proportion, percentages or ppm) should not change the results of the analysis (i.e., *scale invariance*). In fact, analyses of compositional data should also account for other features of these data, such as their *relative scale*. For a broader discussion, we refer the reader to Pawlowsky-Glahn et al. (2015a). Because of the range limitation and the possible spurious correlation of compositional vectors (Kim 1999; Pawlowsky 1984), the Euclidean-based statistical framework was proved to be ineffective for the spatial prediction of this type of data, although a number of authors have ignored this aspect, see e.g. Delbari et al. (2011). Other works (e.g. Walvoort and De Gruijter 2001) tried to account for the particular nature of regionalised variables expressing relative fractions by proposing an extension of kriging called Compositional Kriging (CK). CK predictions respect the constraints of positivity and constant sum value. However, the CK algorithm is based on empirical considerations rather than a coherent probabilistic model, and is therefore not suited for stochastic simulation. Our developments follow the direction of research on compositional kriging explored by Pawlowsky (1989), Tolosana-Delgado et al. (2011), who formulated geostatistical models and methods based on the Aitchison geometry for the simplex (see Pawlowsky-Glahn and Egozcue 2016; Tolosana-Delgado et al. 2019 for recent reviews).

Presently, the standard approach to the statistical analysis of compositional data is the one pioneered by Aitchison (1982), which is based on the particular geometry of the simplex (Aitchison 1986; Pawlowsky-Glahn et al. 2015b; Pawlowsky-Glahn and Egozcue 2001; Billheimer et al. 2001). A p -dimensional compositional vector $\mathbf{Z} =$

(Z_1, \dots, Z_p) is an element of the p -dimensional standard simplex, \mathbb{S}^p , which is defined as

$$\mathbb{S}^p = \left\{ (Z_1, \dots, Z_p : Z_i \geq 0, \sum_{i=1}^p Z_i = 1) \right\}. \tag{10}$$

In Aitchison (1986), Pawlowsky-Glahn and Egozcue (2001) group operations are defined to give the simplex a structure of a real vector space. These are the perturbation \oplus (sum) and powering \odot (product by a constant) operations, defined, for $\mathbf{x}, \mathbf{y} \in \mathbb{S}^p$, and $\alpha \in \mathbb{R}$, respectively as

$$\mathbf{x} \oplus \mathbf{y} = \mathcal{C}(x_1 y_1, \dots, x_p y_p),$$

$$\alpha \odot \mathbf{x} = \mathcal{C}(x_1^\alpha, \dots, x_p^\alpha).$$

Here, $\mathcal{C}(\cdot)$ denotes the closure operation

$$\mathcal{C}(\mathbf{x}) = \left(\frac{x_1}{\sum_{i=1}^p x_i}, \dots, \frac{x_p}{\sum_{i=1}^p x_i} \right) \quad \mathbf{x} \in \mathbb{R}_+^p.$$

The space \mathbb{S}^p can be equipped with a (finite-dimensional) Hilbert space structure when considering the Aitchison inner product, defined, for $\mathbf{x}, \mathbf{y} \in \mathbb{S}^p$ as

$$\langle \mathbf{x}, \mathbf{y} \rangle_a = \frac{1}{2p} \sum_{i=1}^p \sum_{j=1}^p \ln \frac{x_i}{x_j} \cdot \ln \frac{y_i}{y_j} \quad \mathbf{x}, \mathbf{y} \in \mathbb{S}^p.$$

The inner product induces a norm $\| \cdot \|_a := \sqrt{\langle \cdot, \cdot \rangle_a}$, which in turn induces a distance $d_a(\mathbf{x}, \mathbf{y}) = \| \mathbf{x} \ominus \mathbf{y} \|_a$, $\mathbf{x}, \mathbf{y} \in \mathbb{S}^p$, where $\mathbf{x} \ominus \mathbf{y}$ denotes the perturbation of \mathbf{x} with the reciprocal of \mathbf{y} , i.e., $\mathbf{x} \ominus \mathbf{y} = \mathbf{x} \oplus ((-1) \odot \mathbf{y})$. The Hilbert space structure identified by these operations is called Aitchison geometry, or Aitchison simplex (Pawlowsky-Glahn and Egozcue 2001).

3.2 ATPRCoK in the Aitchison geometry

The statistical approach proposed by Aitchison (1982) and following authors (Pawlowsky-Glahn et al. 2015b; Pawlowsky-Glahn and Egozcue 2002) consists of analyzing compositional data in the context of the Aitchison geometry. Here, a large number of multivariate statistical methods (e.g., principal component analysis, regression) can be properly reformulated to account for the inherent properties of compositional data (e.g., scale invariance, relative scale), research in this field is still ongoing, see e.g. Rodríguez-Díaz et al. (2020). From an operational viewpoint, the standard procedure of analysis consists of transforming the original data by applying an isomorphism from the p -dimensional Aitchison simplex to the classical Euclidean space \mathbb{R}^{p-1} (or in some cases \mathbb{R}^p), perform the statistical analysis on the transformed data and finally back-transform the results in the original space. This strategy was proved to be fully equivalent to working

directly in the Aitchison simplex for a number of statistical methods (see, e.g. Filzmoser et al. 2018). In this section, we shall formulate ATPRCoK in the Aitchison simplex, and then, using the *principle of working on coordinates* (Mateu-Figueras et al. 2011), prove that one may equivalently perform the analysis by relying on the so-called Isometric Log-Ratio (ILR) transformation, which is an isometry that maps the simplex to \mathbb{R}^{p-1} . The latter associates to a compositional vector $\mathbf{z} \in \mathbb{S}^p$ the coordinates of this vector with respect to an orthonormal basis of the simplex $(\psi_1, \dots, \psi_{p-1})$, i.e.,

$$\text{ILR}(\mathbf{z}) = (\langle \mathbf{z}, \psi_1 \rangle_a, \dots, \langle \mathbf{z}, \psi_{p-1} \rangle_a)' \tag{11}$$

Note that the ILR is a linear transformation, see e.g. Egozcue et al. (2003). For several compositional methods (e.g., principal component analysis, regression, see, e.g., Pawlowsky-Glahn et al. 2015b), it was shown that the choice of the basis does not influence the results of the analysis. However, specific choices for the basis can lead to practical advantages. For instance, the basis could be chosen in such a way as to grant uncorrelation of the resulting transformed data, or to ease the interpretation of the results (see, e.g. Fišerová and Hron 2011).

In Dobarco et al. (2016) the authors used ATPCoK to downscale psfs data transformed with Additive-Log Ratio (ALR), using the silt fraction as a reference for the ratios. ALR was there used as a practical solution to account for the compositional nature of the data, but the modelling assumptions were not explicitly stated, and the results were interpreted only in terms of prediction accuracy with respect to a given test set. Recent works highlighted some limitation of the ALR transformation, as this is not isometric, thus does not preserve the Aitchison geometry (see, e.g., Pawlowsky-Glahn et al. 2015b, p. 60). In this section, we propose a general method for the statistical downscaling and simulation of compositional data which extends the ATPRCoK to the context of the Aitchison simplex. We call the method ILR-ATPRCoK to recall the computational strategy we propose to perform ATPRCoK in the Aitchison simplex, which is based on the ILR transformation. Here, we shall also prove that the strategy based on ILR is fully equivalent to working directly in the simplex itself. We remark that this would not be the case for the ALR transformation, as our developments strongly relies on the isometric property of ILR, as further discussed in the following sections.

In the following, we denote by $\{\mathbf{Z}(\mathbf{x}), \mathbf{x} \in D\}$ a random field valued in \mathbb{S}^p , defined over a Euclidean region $D \subset \mathbb{R}^d$. To indicate the Aitchison center of the field (i.e. the mean value in Aitchison geometry), the Aitchison covariance and the integral operator of $\mathbf{Z}(\mathbf{x})$ over a region $v \subset \mathbb{R}^d$, we use respectively the same notation used in

Menafoglio et al. (2014), Pawlowsky-Glahn and Buccianti (2011), that is

- Aitchison center,

$$\boldsymbol{\mu}(\mathbf{x}) = \text{Cen}(\mathbf{Z}(\mathbf{x})) = \operatorname{argmin}_{\mathbf{z} \in \mathbb{S}^p} \mathbb{E}[d_a^2(\mathbf{Z}(\mathbf{x}), \mathbf{z})];$$
- Aitchison covariance operator, acting on a (non-random) element $\mathbf{z} \in \mathbb{S}^p$, for $\mathbf{x}_1, \mathbf{x}_2 \in D$, as

$$C_a(\mathbf{x}_1, \mathbf{x}_2)\mathbf{z} = \text{Cen}[(\mathbf{Z}(\mathbf{x}_1) \ominus \boldsymbol{\mu}(\mathbf{x}_1), \mathbf{z})_a \odot (\mathbf{Z}(\mathbf{x}_2) \ominus \boldsymbol{\mu}(\mathbf{x}_2))];$$
- Aitchison integral over a spatial region v ,

$$\int_v^\oplus \mathbf{Z}(\mathbf{x})\eta(d\mathbf{x}) = \mathcal{C}\left(e^{\int_v \ln(Z_1(\mathbf{x}))\eta(d\mathbf{x})}, \dots, e^{\int_v \ln(Z_p(\mathbf{x}))\eta(d\mathbf{x})} \right),$$

denoting by η a positive measure over D .

We refer to Aitchison (1986), Egozcue et al. (2003) for an insight of the geometry of the random compositions in the Aitchison simplex, and to Bosq (2000) for a recall on covariance operators in Hilbert spaces.

For the element $\mathbf{Z}(\mathbf{x})$ of the compositional field at $\mathbf{x} \in D$, we assume the following model

$$\mathbf{Z}(\mathbf{x}) = \boldsymbol{\mu}(\mathbf{x}) \oplus \mathbf{e}(\mathbf{x}), \tag{12}$$

with $\mathbf{e}(\mathbf{x})$ the residual term. We model the center of the field through a linear model in \mathbb{S}^p

$$\boldsymbol{\mu}(\mathbf{x}) = \bigoplus_l u^l(\mathbf{x}) \odot \boldsymbol{\beta}^l, \tag{13}$$

where $\boldsymbol{\beta}^l \in \mathbb{S}^p, l = 1, \dots, L$ are the vectors of unknown regression coefficients and $u^l(\mathbf{x}) \in \mathbb{R}, \mathbf{x} \in D$, are the covariates. Furthermore, we assume that the residual is second-order stationary, see Tolosana-Delgado et al. (2019), i.e., that the covariance structure (in Aitchison geometry) for a random composition $\mathbf{Z}(\mathbf{x}), \mathbf{x} \in D$, only depends on the increment among locations

$$C_a(\mathbf{x}_1, \mathbf{x}_2) = \tilde{C}_a(\mathbf{x}_1 - \mathbf{x}_2), \quad \mathbf{x}_1, \mathbf{x}_2 \in D.$$

To simplify the notation, we shall indicate the stationary covariance function \tilde{C}_a simply by C_a .

With a notation analogue to that used in Sect. 2, we consider the discretized versions of the field $\{\mathbf{Z}(\mathbf{x}), \mathbf{x} \in D\}$, denoted by \mathbf{Z}_k (resp. \mathbf{Z}_K) and obtained at a fine (resp. coarse) discretization scale, namely

$$\mathbf{Z}_k = \frac{1}{|v_k|} \odot \int_{v_k}^\oplus \mathbf{Z}(\mathbf{x})\eta(d\mathbf{x}), \quad \mathbf{Z}_K = \frac{1}{|v_K|} \odot \int_{v_K}^\oplus \mathbf{Z}(\mathbf{x})\eta(d\mathbf{x}).$$

Here, the powering by $\frac{1}{|v_k|}$ and $\frac{1}{|v_K|}$ is intended as acting element-wise.

Given the realization of the coarse-scale field \mathbf{Z}_K , and by analogy with (2), we define the ILR-ATPRCoK predictor of the fine-scale field $\mathbf{Z}_k \in \mathbb{S}^p$ as

$$\widehat{\mathbf{Z}}_k = \bigoplus_l u_k^l \odot \boldsymbol{\beta}^l \oplus \bigoplus_K A_K \square \mathbf{e}_K, \tag{14}$$

where $A_K \in \mathbb{R}^{p \times p}$ is a matrix of ATPCoK unknown weights to be optimized, \square is the matrix-by-composition multiplication—consistent with perturbation and powering, as defined in Pawlowsky-Glahn et al. (2015b) (p. 55), i.e., denoting by $e_{K,i}, i = 1, \dots, p$ the elements of $\mathbf{e}_K \in \mathbb{S}^p$ and by $\lambda_{K,i,j}, i, j = 1, \dots, p$ the elements of A_K

$$A_K \square \mathbf{e}_K = \mathcal{C} \left[\prod_{j=1}^p e_{K,1,j}^{\lambda_{K,1,j}}, \dots, \prod_{j=1}^p e_{K,p,j}^{\lambda_{K,p,j}} \right].$$

The residuals $\mathbf{e}_K \in \mathbb{S}^p$ represent the upscaled residuals of (12), defined as

$$\mathbf{e}_K = \frac{1}{|v_K|} \odot \int_{v_K}^{\oplus} (\mathbf{Z}(\mathbf{x}) \ominus \boldsymbol{\mu}(\mathbf{x})) \eta(d\mathbf{x}). \tag{15}$$

Under the assumption that the regression coefficients $\boldsymbol{\beta}^l$ and the covariance function C_a are known, the optimal weights A_K in (14) are found as to guarantee that the ILR-ATPRCoK is the Best Linear Unbiased Predictor (BLUP) in \mathbb{S}^p , i.e., by solving the following constrained minimization problem:

Area-to-Point Regression Cokriging in the Aitchison simplex

$$\arg \min_{A_K \in \mathbb{R}^{p \times p}} \mathbb{E} [d_a^2(\bigoplus_K A_K \square \mathbf{e}_K, \mathbf{e}_k)] \quad \text{s.t.} \tag{16}$$

$$\text{Cen}(\bigoplus_K A_K \square \mathbf{e}_K) = \bar{\boldsymbol{\mu}},$$

where $\bar{\boldsymbol{\mu}}$ is the spatially constant residual mean.

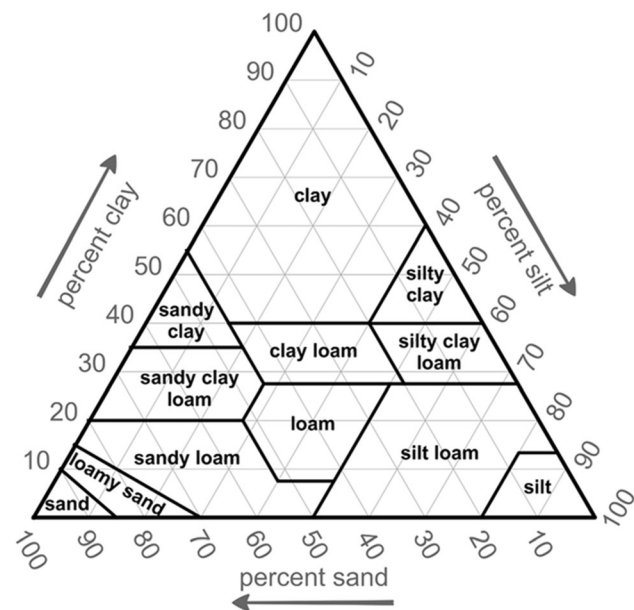


Fig. 1 Soil texture triangle: soil texture classification according to the USDA classification system, based on relative fractions of clay, silt and sand. Figure taken from Groenendyk et al. (2015)

Note that the operator defined by A_K in Eq. (16) is a linear endomorphism, hence enjoys of the properties described in Pawlowsky-Glahn et al. (2015a) (p. 56). In particular, the rows and columns of A_K add to zero, for further details see “Appendix A”. The following result states that, by applying the ILR transformation (i.e. mapping from the Aitchison simplex to an Euclidean space, via an isometric isomorphism), one obtains an equivalent formulation of problem (16) that can be solved using the standard ATPRCoK presented in Sect. 2. The proof of Proposition 1 is given in “Appendix A”.

Proposition 1 Given a compositional random field $\mathbf{Z}(\mathbf{x})$ valued in \mathbb{S}^p and a random field $\mathbf{Y}(\mathbf{x})$ valued in \mathbb{R}^{p-1} defined as $\mathbf{Y}(\mathbf{x}) = \text{ILR}(\mathbf{Z}(\mathbf{x}))$ for $\mathbf{x} \in D$, the BLUP in \mathbb{S}^p for \mathbf{Z}_k —found by solving (16)—coincides with the ILR-back-transformed ATPRCoK predictor for \mathbf{Y}_k defined in (2), i.e.,

$$\widehat{\mathbf{Z}}_k = \text{ILR}^{-1}(\widehat{\mathbf{Y}}_k). \tag{17}$$

Even though $\boldsymbol{\beta}^l$ is rarely *a priori* known, an estimate of $\boldsymbol{\beta}^l$ can be obtained by back-transforming the corresponding estimate of the coefficient vectors $\boldsymbol{\beta}_Y^l$ referred to the ILR-transformed field $\mathbf{Y}(\mathbf{x})$ (see, e.g., Pawlowsky-Glahn et al. 2015b). Similarly, an estimate of the covariance operator C_a can be obtained from the estimated (Euclidean) covariance operator C_Y of the vector field $\mathbf{Y}(\mathbf{x})$. In this work, for $\boldsymbol{\beta}_Y^l$ we shall consider OLS estimates, whereas for C_Y the estimates obtained by Goovaerts’ deconvolution (Goovaerts 2008) of classical cross-variograms.

Note that the equivalence between the ATPRCoK in the Aitchison simplex and the Euclidean ATPRCoK on ILR-transformed data (as stated in Proposition 1), implies the possibility to analogously perform Block Sequential Gaussian Simulation (BSGS, Benndorf 2004, 2003), as BSGS grounds on the same hypothesis as ATPRCoK, and it is indeed based on the latter method. In the context of Uncertainty Quantification (UQ), BSGS is key to propagate the uncertainty in numerical models that take as input downscaled compositional data, as we discuss in Sect. 6.

Finally, one should note that, since the assumptions are made with respect to the Aitchison geometry, the *mass-preserving* property as stated in the Euclidean framework, i.e. (see Kyriakidis 2004),

$$\mathbf{Z}_K = \frac{1}{P} \sum_{k:\mathbf{x}_k \in v_K} \widehat{\mathbf{Z}}_k,$$

does not hold. In a discrete prediction setting, as in Sect. 2, the Aitchison geometry predictions respect the following *centre-preserving* property

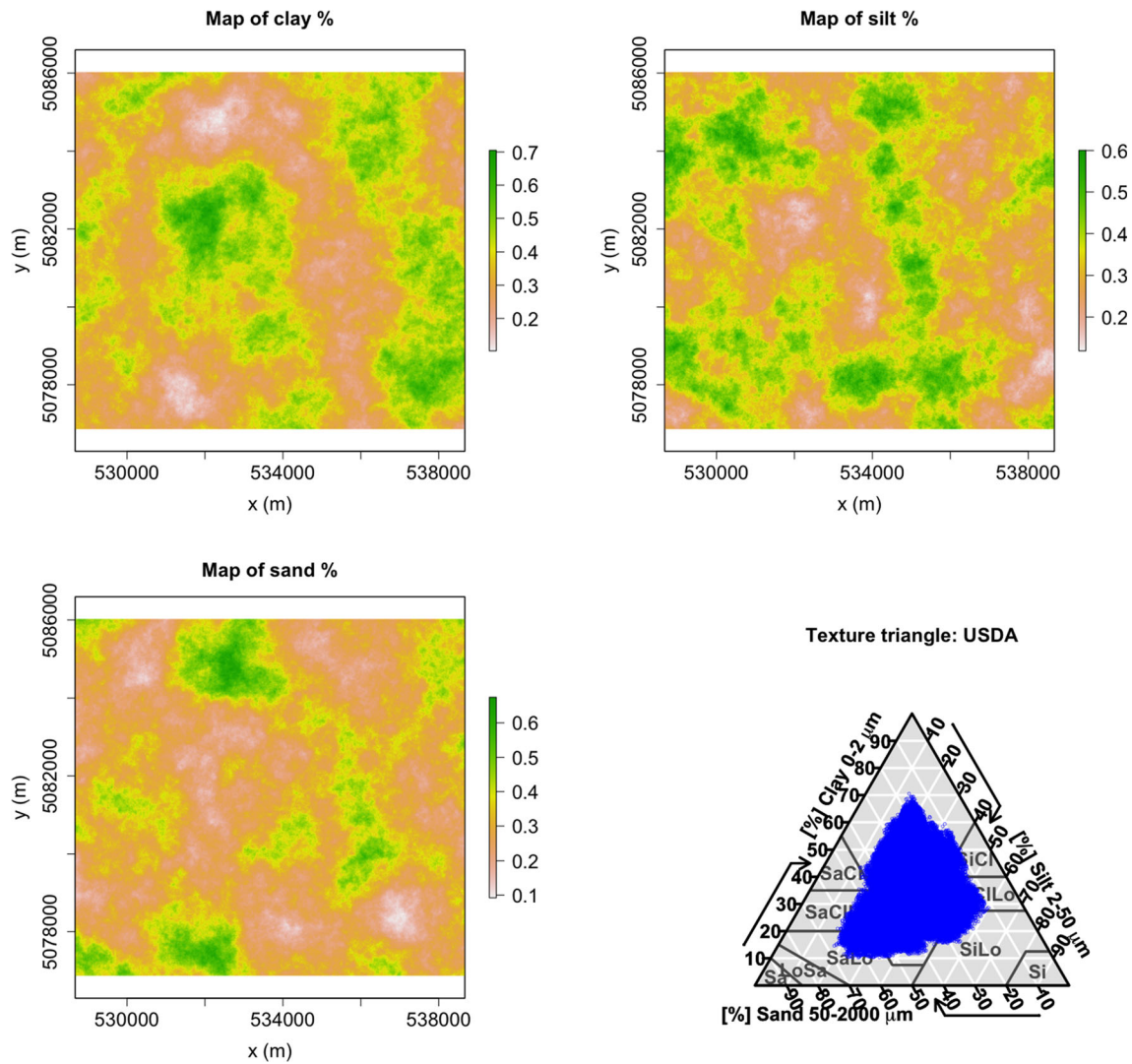


Fig. 2 One realization of the initial psfs field $\mathbf{Z}(\mathbf{x})$ with $Cen(\mathbf{Z}(\mathbf{x})) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and $\sigma^2 = 0.1$

$$\mathbf{Z}_K = \frac{1}{P} \odot \bigoplus_{k: \mathbf{x}_k \in v_K} \hat{\mathbf{Z}}_k. \tag{18}$$

Indeed, since $P = \frac{|v_K|}{|v_k|} \in \mathbb{Z}^+$, as defined in Sect. 2, one has

$$\mathbf{Z}_K = \mathcal{C} \left(\left(\prod_{k=1}^P \hat{\mathbf{Z}}_{1,k} \right)^{\frac{1}{P}}, \dots, \left(\prod_{k=1}^P \hat{\mathbf{Z}}_{n,k} \right)^{\frac{1}{P}} \right). \tag{19}$$

This means that, in the Aitchison simplex, coarse areal data coincide with the geometric mean of the predicted fine areal values (normalized to having unit-sum).

In the following sections we exemplify the proposed methodology through its application to particle-size fractions, whose definition is recalled in the next Sect. 4.

4 Particle size fractions

Soil texture is a classification instrument used to determine soil classes. More specifically, soil texture is quantitatively determined on the basis of the relative fractions of the fine particles of different sizes that compose the terrain. Soil particles under 2 mm are divided in three groups

- clay: particles with a diameter less than 2 μm ;
- silt: particles with a diameter between 2 and 50 μm ;
- sand: particles with a diameter between 50 μm and 2 mm.

Fractions of clay, silt and sand are usually indicated as particle-size fractions (psfs). Soil texture classes are determined by the relative percentages of psfs, according to a standard that may vary depending on the country.

The most common classification is that used by the United States Department of Agriculture (USDA Schaefer

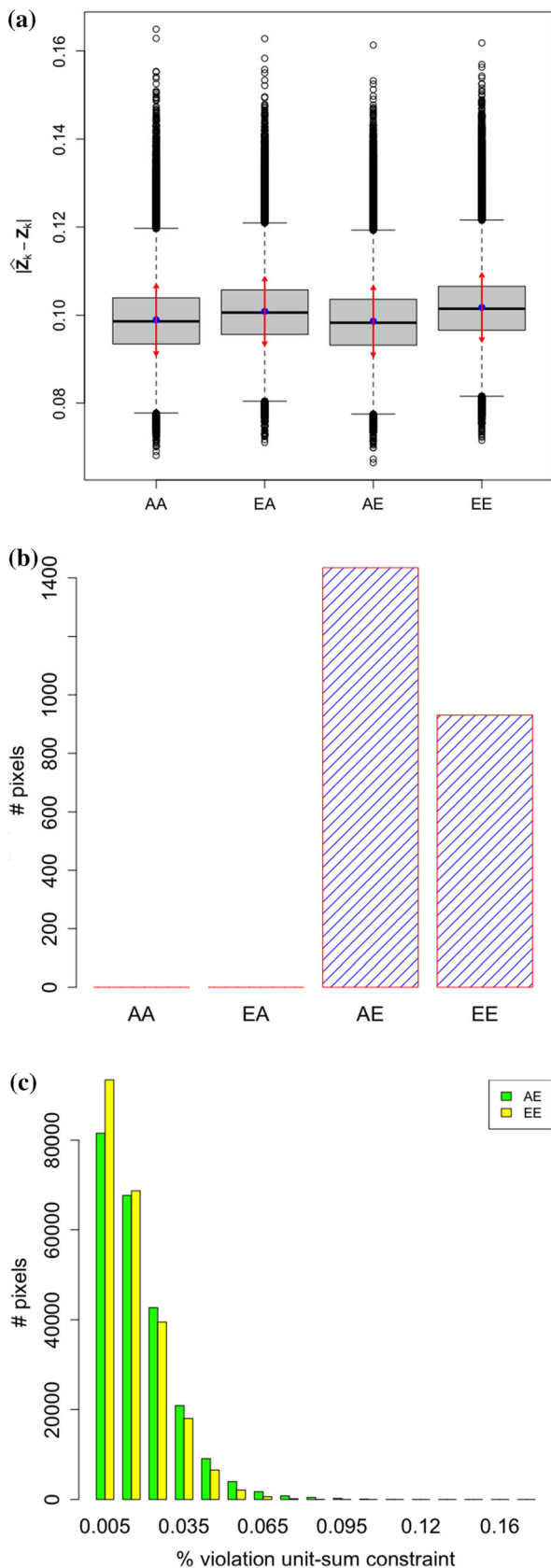


Fig. 3 Synthetic data results: **a** boxplots of the sample mean error. Blue points and red segments are respectively the spatial mean and spatial standard deviation of the sample mean error. **b** the mean (across realizations) number of pixels that violate the positivity and unit-sum constraint for the four methods considered. **c** relative violation of the unit-sum constraint. A value of 1% on the x-axis indicates that the reconstructed psfs sums e.g. to 1.01. The histograms related to methods AA and EA are not reported since the ILR-ATPRCoK method by construction guarantees that resulting psfs sum to 1

et al. 2007), which distinguishes twelve major soil texture classes shown in Fig. 1. The classes are typically named after the primary constituent particle-size or a combination of the most abundant particles sizes, e.g. sandy clay or silty clay. A fourth term, “loam”, is used to describe equal proportions of sand, silt, and clay in a soil sample, and leads to the naming of even more classes, e.g. clay loam or silt loam.

5 Validation

The geostatistical method outlined in the previous sections has been implemented in R-3.6 (R Development Core Team 2008) using the libraries *gstat* (Gräler et al. 2016; Pebesma 2004) for ATPK and geostatistical simulation and *compositions* (Boogaart and Tolosana-Delgado 2008) for the analysis of compositional data. In particular, for the variogram deconvolution we use the Goovaerts’ procedure (Goovaerts 2008, 2010). We define a continuous random field $\mathbf{Z}(\mathbf{x}) = (Z_1(\mathbf{x}), Z_2(\mathbf{x}), Z_3(\mathbf{x}))$ where

$Z_1(\mathbf{x}) =$ % of part 1 at location \mathbf{x} of the domain D ;

$Z_2(\mathbf{x}) =$ % of part 2 at location \mathbf{x} of the domain D ;

$Z_3(\mathbf{x}) =$ % of part 3 at location \mathbf{x} of the domain D .

In our analysis, compositional data are transformed using the function ILR of the package *compositions*, see e.g. Boogaart and Tolosana-Delgado (2008). The basis used for the transformation is the one introduced in Egozcue et al. (2003), based on the partition of the vector of compositional variables in two sub-compositions, i.e.,

$$\{\psi_1, \psi_2\} = \left\{ \mathcal{C} \left[\exp \left(\sqrt{\frac{1}{2}}, -\sqrt{\frac{1}{2}}, 0 \right) \right], \mathcal{C} \left[\exp \left(\sqrt{\frac{1}{2p}}, \sqrt{\frac{1}{2p}}, -\sqrt{\frac{2}{p}} \right) \right] \right\}$$

For the sake of illustration and in view of the motivating study, we shall interpret $\mathbf{Z}(\mathbf{x})$ as the psfs at \mathbf{x} (i.e., $Z_1(\mathbf{x})$, $Z_2(\mathbf{x})$, $Z_3(\mathbf{x})$ represent the composition in clay, silt and sand, respectively). Nonetheless, the validity of the

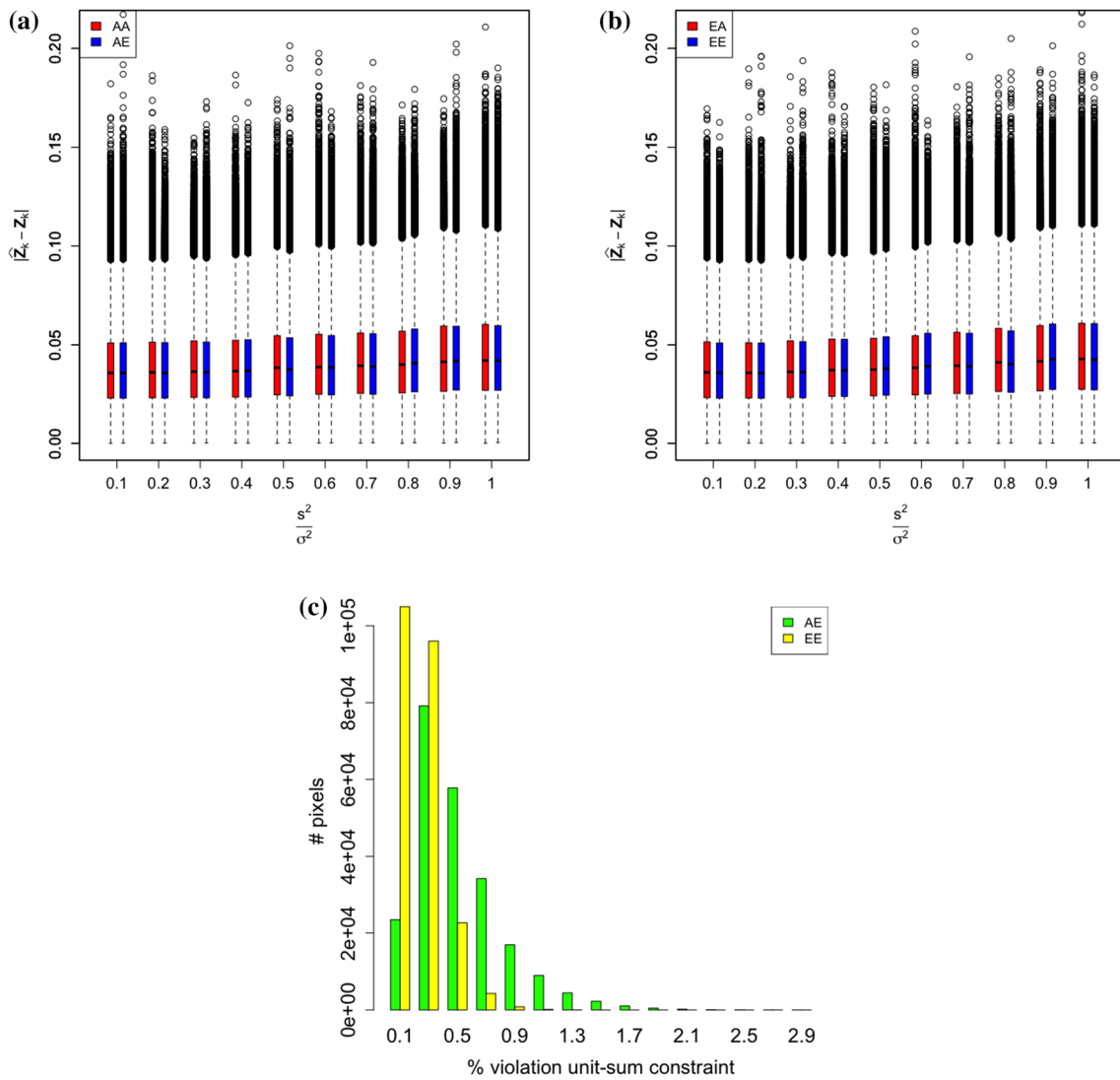


Fig. 4 Validation on synthetic data. **a, b** Boxplots of the errors between initial and predicted psfs data. **c** Histograms of the maximum of the violation of the unit-sum constraint experienced during the

tested variances s^2 . In panel **c** the histograms corresponding to AA and EA are not reported, since the ILR-ATPRCoK method by construction guarantees that the resulting psfs sum to 1

simulation study here presented is clearly not limited to the specific case of psfs.

5.1 Synthetic data

To assess the performance of the proposed method, we consider a simulated dataset $\mathbf{Z}(\mathbf{x}), \mathbf{x} \in D$, with support measure $|v_k| = 20 \times 20 \text{ m}^2$, on a given rectangular domain D with area $|D| = 10000 \times 9160 \text{ m}^2$. The compositional vector $\mathbf{Z}(\mathbf{x})$ is modelled as a process with a given spatially-constant center $Cen(\mathbf{Z}(\mathbf{x})) = \boldsymbol{\mu}$ and stationary-isotropic covariance structure. From the operational viewpoint, the mean $\boldsymbol{\mu} = \mathcal{C}[(\mu_1, \mu_2, \mu_3)']$ is set based on independent uniform distributions $\mu_i \sim U(0, 1), i = 1, 2, 3$. Compositions

were simulated by back-transforming through ILR^{-1} two-dimensional Gaussian random vectors \mathbf{Y} , with constant mean $\boldsymbol{\mu}_Y = ILR(\boldsymbol{\mu})$, and stationary-isotropic marginal variograms from the spherical model without nugget (Chilès and Delfiner 2012; Cressie 1993). For the following simulations, the components of \mathbf{Y} are always assumed to be uncorrelated, and the marginal ranges are both set to 2000 m. In each simulation, the common sill σ^2 is sampled according to a uniform distribution $U[0.025, 2.5]$. In Fig. 2 we show an example of realization of the psfs distribution.

Starting from this set of synthetic psfs, we perform a sequence of upscaling-downscaling procedures, as follows. Downscaling is done using either ATPRCoK or ILR-ATPRCoK and upscaling either in Euclidean or Aitchison

geometry, so that four different possibilities arise. In the following, we call AA the upscaling in the Aitchison simplex and downscaling via ILR-ATPRCoK, EE the upscaling in the Euclidean space and downscaling via ATPRCoK whereas EA, AE are the mixed methods, referring respectively to upscaling in the Euclidean space and to downscaling via ILR-ATPRCoK and upscaling in the Aitchison geometry and downscaling via ATPRCoK. Given that, usually, little information is available on how the coarse-scale map relates with the fine-scale map (i.e., how a block value is obtained from smaller cells), we test the performance of both ATPRK and ILR-ATPRCoK on both Aitchison and Euclidean upscaling.

For each method, we consider a set of 100 realizations of the fine scale compositional field, each yielding a reconstructed field after the upscaling-downscaling process. The upscaling factor P is set each time by randomly

and independently sampling in the discrete range $\{2^2, 3^2, \dots, 30^2\}$. For each method and each realization we compute the sample mean error, i.e. the average of the Euclidean distance between initial and reconstructed psfs fields—the average being taken over the realizations. We use the Euclidean distance as a metric for comparison, since the ATPRCoK produces vectors that do not necessarily belong to the simplex due to the absence of the closure operation for this method. In those cases, the Aitchison distance is not well defined, and thus cannot be used for comparison purposes.

Even if the distribution of the sample mean error between the considered methods, reported in Fig. 3a, would suggest a substantial equivalence among the methods, Fig. 3b, c clearly show that, unlike ATPRCoK, ILR-ATPRCoK is able to produce psfs maps that are consistent with the unit-sum and positivity constraints. Indeed, the

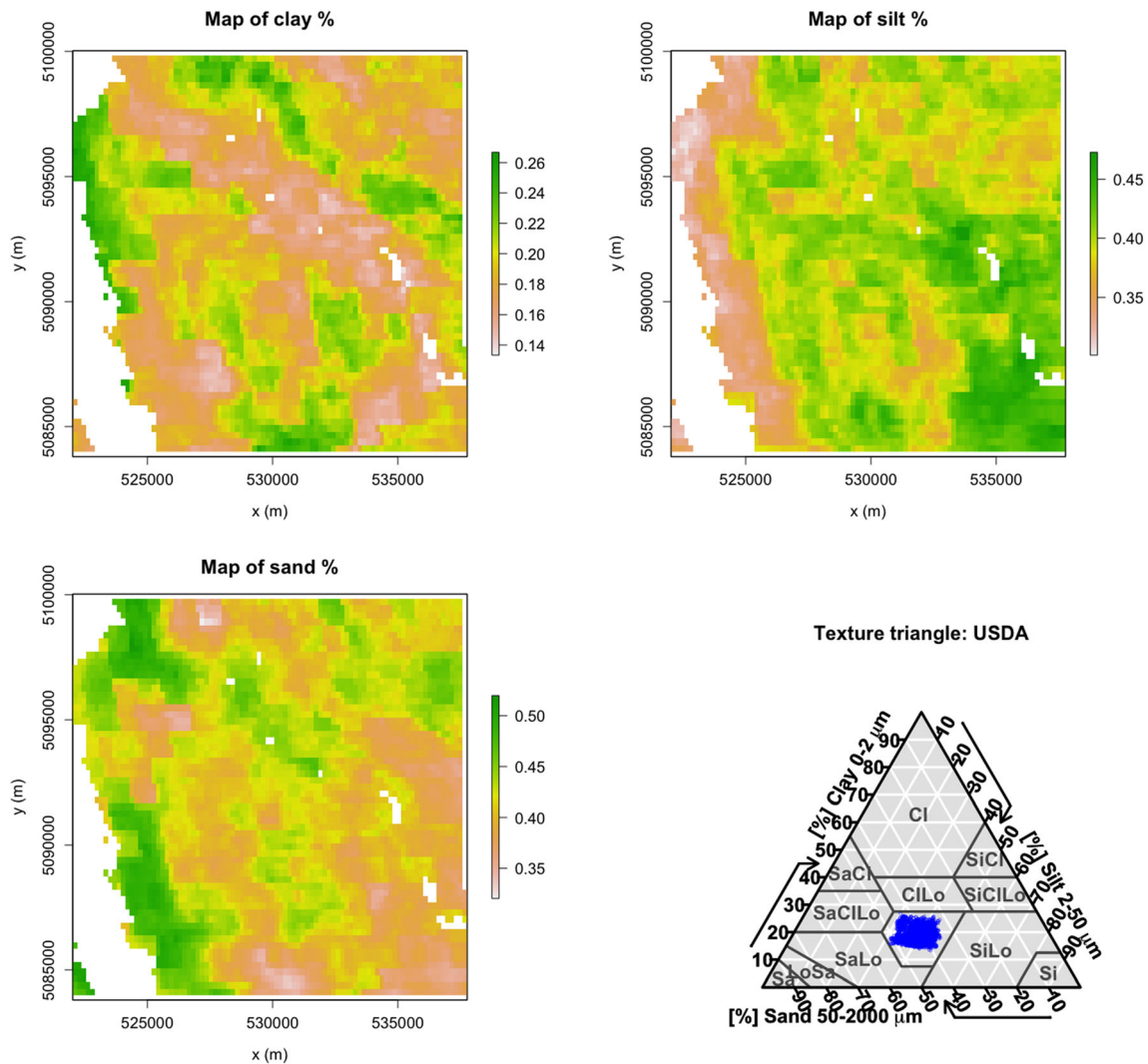


Fig. 5 SoilGrids psfs data within the study region. The considered data refer to the mean value of psfs as reconstructed in the SoilGrids repository

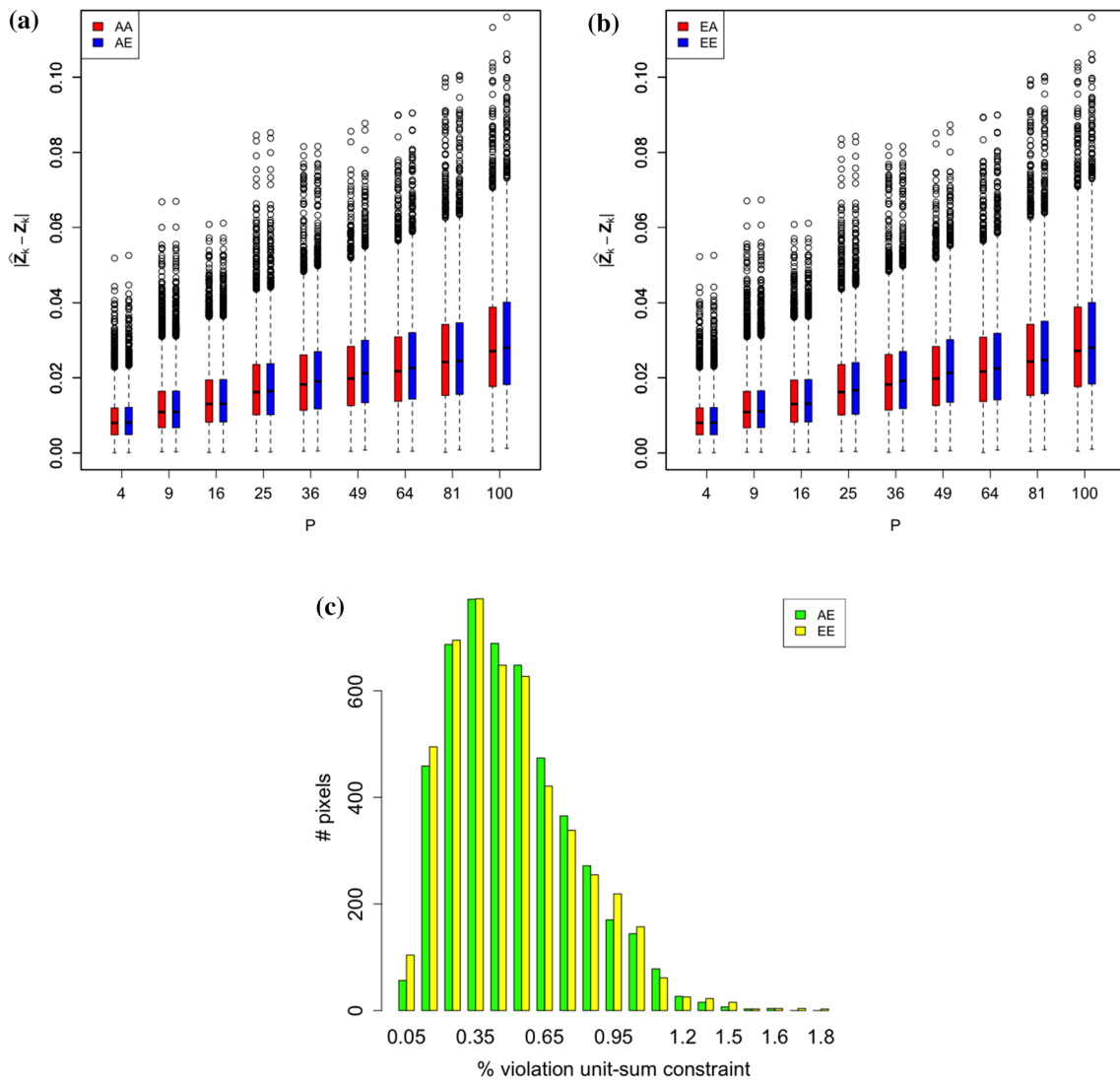


Fig. 6 Downscaling on SoilGrids data. **a, b** Boxplots of the errors between SoilGrids and predicted psfs data. **c** Histograms of the maximum of the violation of the unit-sum constraint experienced during the tested upscaling factors. In panel **c** the histograms

ATPRCoK shows a violation of the positivity constraint in about 10^3 pixels on average, representing roughly 0.5% of the study area.

We then perform a sensitivity analysis with respect to random perturbations of the initial data of the downscaling procedure, for the four methods described above. This case is representative of input data characterized by a given degree of uncertainty. We thus consider a realization of the synthetic psfs field in case of $\mu = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and sill $\sigma^2 = 0.1$ (i.e. as in Fig. 2) and we set the upscaling factor to $P = 225$. Let us indicate with $K = 1, \dots, M$ the elements of the coarse maps. The upscaled maps Z_K are then perturbed with a set of i.i.d. Gaussian random errors ϵ_K . Similarly as before, these perturbations were generated on the ILR

corresponding to AA and EA are not reported, since the ILR-ATPRCoK method by construction guarantees that resulting psfs sum to 1

transforms, by adding a zero-mean independent Gaussian error with variance s^2 , ranging from 10 to 100% of the sill σ^2 .

In Fig. 4a, b we report the boxplots of the errors for each value of s^2 . The errors are computed, for each pixel, as the Euclidean distance between initial and predicted psfs. We note that both ATPRCoK and ILR-ATPRCoK are quite robust even in case of relatively high perturbations of the initial data. In Fig. 4c we show the histograms of the maximum, across simulations, of the violation of the unit-sum constraint. For instance, a vertical bar in correspondence of the range [1,1.2] indicates the count of pixels whose maximum discrepancy (across simulations) from unity of the sum of psfs is between 1% and 1.2% (i.e., the sum is in [1.01,1.012] or [0.988,0.99]). These results

clearly show the ability of ILR-ATPRCoK method to produce results consistent with the unit-sum constraint, unlike the ATPRCoK method which yields maps with a significant violation of the aforementioned constraint. Note that, beside the constraints, ILR-ATPRCoK allows accounting for the relative scale of compositional data, avoiding to model spurious correlations (Pawlowsky-Glahn et al. 2015b, p. 23).

5.2 Downscaling SoilGrids data

In this section, we test the performances of the proposed method in downscaling psfs from soil digital maps publicly available. This case is considered to analyse compositional random fields having a realistic spatial distribution. For this purpose and in view of our case study, we consider SoilGrids, which is a system for automated digital soil mapping based on state-of-the-art spatial predictions methods (Hengl et al. 2014, 2017) released in 2014 by ISRIC (International Soil Reference and Information Centre)—World Soil Information, a non-profit organization funded by the Dutch government. SoilGrids predictions are based on globally fitted models using soil profile and environmental covariate data. When first released, SoilGrids provided a collection of soil properties and class maps of the world at 1 km spatial resolutions, produced using

automated soil mapping based on statistical regression models. In 2017, the resolution has been increased to 250 m and the accuracy of the predictions has been greatly improved by using machine learning algorithms instead of the previously employed linear regression (Hengl et al. 2017). In 2020, SoilGrids released a version where, among other updates, soil map predictions are provided with a mean value together with an uncertainty level map. SoilGrids data are available publicly under the Open DataBase License. Among SoilGrids predicted variables, relevant to this work are clay, silt and sand percentages at different soil depths. In this section, the values considered for geostatistical downscaling are those referred to the topsoil, i.e. depth of 0 cm.

We focus on a geographical domain with area $|D| = 15750 \times 16000 \text{ m}^2$, located in a pre-Alpine area, more specifically the basin of Pioverna river in the Lombardy region in Northern Italy near the city of Lecco. This region was selected as it is similar, from the geomorphological viewpoint, to the area analyzed in the case study presented in Sect. 6. The psfs as available in SoilGrids are reported in Fig. 5. Based on these data, we test the performance of the ILR-ATPRCoK method, at different levels of the upscaling factor P . Following the procedure described in Sect. 5.1, we consider a sequence of upscaling-downscaling operations, both in Aitchison and Euclidean geometry, of the SoilGrids data in Fig. 5. For each upscaling factor P in the range $\{2^2, 3^2, \dots, 10^2\}$ and for each pixel in D , we compute the Euclidean distance of the psfs estimates from the initial SoilGrids data, yielding a set of error (one for each value of P). These are displayed through boxplots in Fig. 6a, b. We note that, mainly at high upscaling factors, the ILR-ATPRCoK method shows slightly better behaviour with respect to the classical ATPRCoK, producing solutions that are closer to the reference ones w.r.t the results produced via ATPRCoK downscaling.

In Fig. 6c we report the histograms of the maximum of the violation of the unit-sum constraint experienced during the set of upscaling factors, for the four cases being considered. Interpretation of these histograms is fully analogous to that in Fig. 4. These results confirm that the ILR-ATPRCoK method is able to produce downscaled maps that are consistent with the unit-sum constraint, as opposed to the ATPRCoK downscaling method. Finally, we do not report any violation of the positivity constraint for ATPRCoK, differently from what shown in the tests reported in Sect. 5.1.

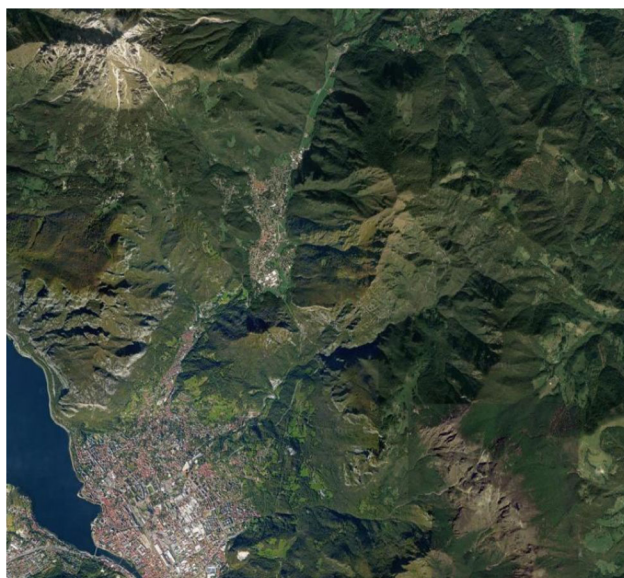


Fig. 7 Aerial view of the case study area

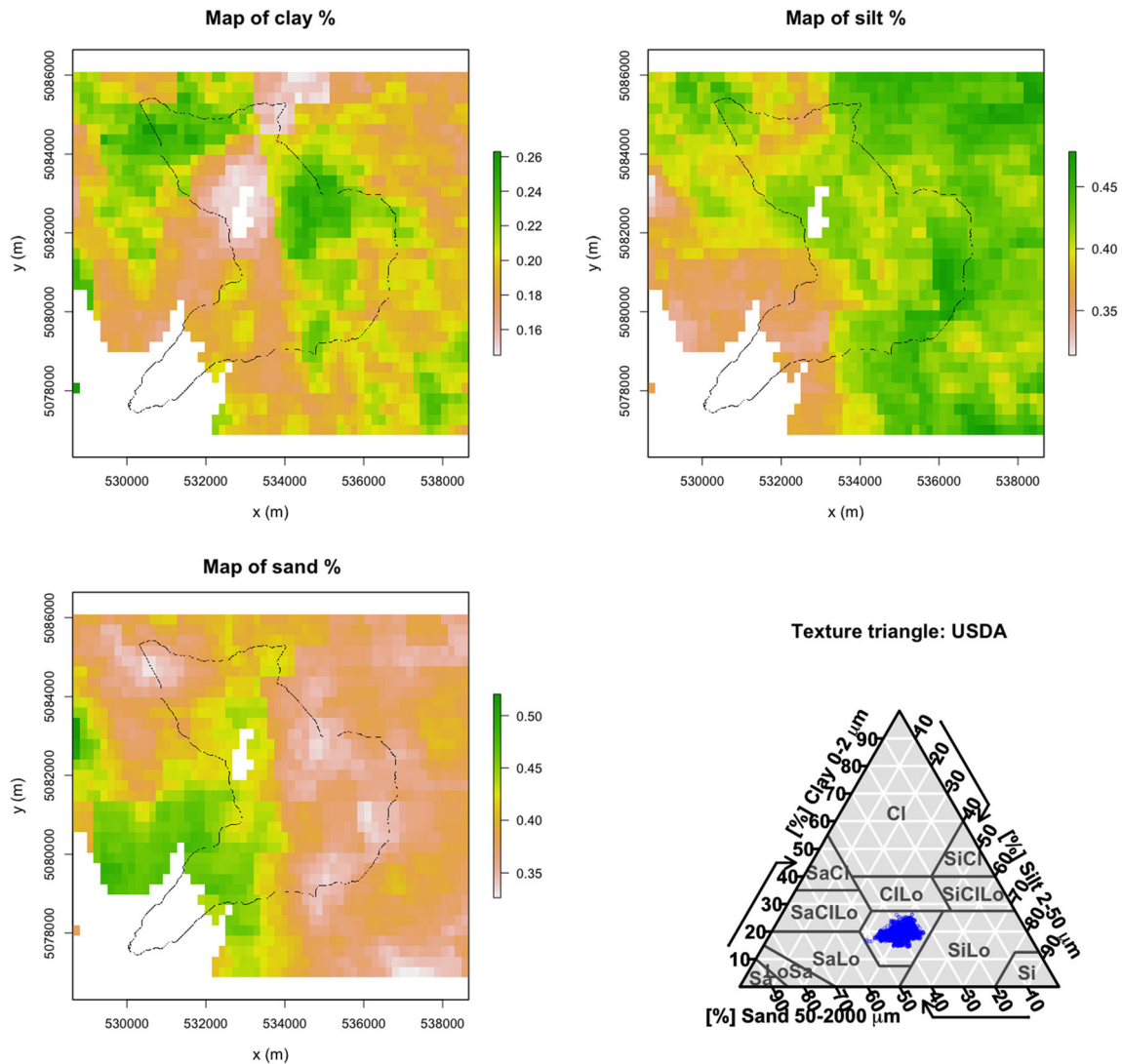


Fig. 8 Clay, silt and sand maps coming from SoilGrids. In black is shown the polygon delimiting the hydrographic Caldone basin

6 A case study

Our case study considers an application to a domain D centered on the city of Lecco, located in the Lombardy region in Northern Italy, which is crossed by three streams (Bione, Caldone, and Gerenzone) that have the typical characteristics of the pre-Alpine area (see Fig. 7). The hydrographic basin of the Caldone water course is 24 km² wide, with an altitude ranging from 197 m a.m.s.l. to 2118 m a.m.s.l. at the top of Grigna Meridionale mountain. Geologically, the basin is characterized by rocky outcrops in the higher part (mainly limestone and clastic rock), while downstream towards the city the river flows through a floodplain. The average precipitation over the city of Lecco is about 1400 mm/year.

The combination between a short hydrologic response time, high slope, intense sediment transport and flow within a densely urban area makes the Caldone river a suitable case study for hydrogeological instability and hazard. This motivates the development of numerical models intended to simulate hydrogeological processes, such as the SMART-SED simulation tool (Sustainable Management of sediment transpORT in responSE to climate change conDitions, Gatti et al. 2020) which is able to simulate sediment transport resulting from slope erosion. These models typically need to be initialized with psfs maps, with a resolution consistent with the Digital Terrain Model (DTM), to be able to model properly the hydrological processes taking place in the study region. However, field measurements of psfs are not available at the

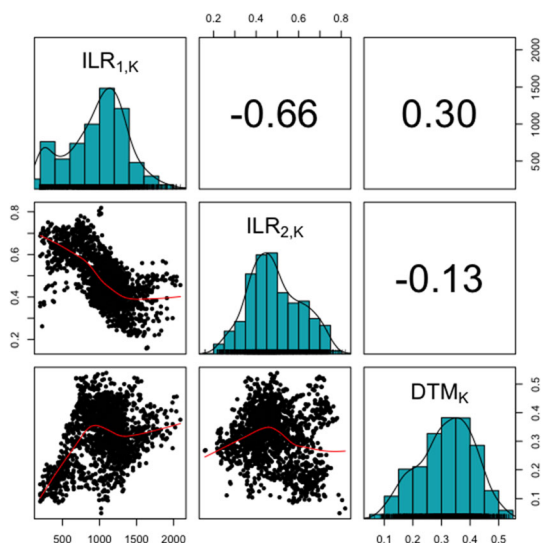


Fig. 9 Scatter plots, histograms and Pearson coefficients of the ILRs and the DTM at coarse resolution

study site, which motivates the use of public repositories to obtain indirect information on these input data (Fig. 8).

SoilGrids psfs at the study region have a pixel support with measure $|v_K| = 250 \times 250 \text{ m}^2$. In terms of the USDA classification, the soil texture of the SoilGrids data for the present case study falls into the loam category. This kind of soil texture, according to Rosso (2004), is classified as fairly permeable soil with moderate infiltration rates and moderate runoff potential. In the following, these coarse-scale data are downscaled to the resolution of the DTM employed for the SMART-SED model, i.e. 5 m, using ILR-ATPRCoK, following the methodology described in Sect. 3. Together with ILR-ATPRCoK results, we here aim to provide random realizations of the psfs fields—obtained via Block Sequential Gaussian Simulation (BSGS)—as demonstration of the ability of the method to produce stochastic compositional maps compatible with coarse scale data.

Based on SoilGrids psfs data, we define the following coarse resolution maps

$$(\text{ILR}_{1,K}, \text{ILR}_{2,K})' = \text{ILR}((Z_{1,K}, Z_{2,K}, Z_{3,K})'),$$

$$K = 1, \dots, M.$$

In the ILR-ATPRCoK model, we consider as covariates u_k^l , $l = 1, \dots, L$, the elevation DTM_k —as given by the DTM at the fine resolution of 5 m—and its square, driven by the parabolic relation displayed in the scatterplot in Fig. 9. For the fine map predictions $\widehat{\text{ILR}}_{1,k}, \widehat{\text{ILR}}_{2,k}$ we thus consider the model

$$\widehat{\text{ILR}}_{1,k} = \beta_0^{(1)} + \beta_1^{(1)} \cdot \text{DTM}_k + \beta_2^{(1)} \cdot \text{DTM}_k^2 + \sum_K \lambda_K e_{1,K},$$

$$\widehat{\text{ILR}}_{2,k} = \beta_0^{(2)} + \beta_1^{(2)} \cdot \text{DTM}_k + \beta_2^{(2)} \cdot \text{DTM}_k^2 + \sum_K \lambda_K e_{2,K}.$$

The fitted values are plotted against the observed values in Fig. 10.

To perform ILR-ATPK, the spatial correlation structure of the fine residuals $e_{1,k}$ and $e_{2,k}$ is estimated by applying the Goovaerts’ deconvolution procedure to the variograms fitted to the coarse residuals $e_{1,K}$ and $e_{2,K}$ (based on a spherical model with nugget), and by assuming $e_{1,k}$ and $e_{2,k}$ to be uncorrelated, see Fig. 11. The latter assumption is supported by the residuals’ analysis, see Fig. 10c, d. Once the fine variograms of the residuals have been estimated, it is possible to solve the ATPK linear system, according to (7). The downscaled ILR are then backtransformed in the Aitchison space in order to get downscaled psfs, see Fig. 12, left column.

Finally, in Fig. 12, right column, we show a sample realization for the downscaled psfs, obtained via BSGS. These stochastic maps shall form the cornerstone to evaluate the propagation of the uncertainty associated with the psfs through the SMART-SED model, and eventually assess the sensitivity of the sediment transport model to this information.

7 Conclusions

We have presented a novel downscaling method for compositional data, based on the ATPRCoK method in the Aitchison geometry, with application to the geostatistical downscaling of psfs data. We have tested the method first in the case of synthetic data and then on a dataset from the SoilGrids online repository. In particular, we have shown the ability of the method to automatically handle the compositional nature of the considered data. Indeed, the proposed method produces maps that respect the unit-sum and positivity constraints and the relative scale of compositions property, as opposed to the classical ATPRCoK method that produces maps which are not consistent with the compositional constraints.

Validation on both synthetic and SoilGrids data show good performances of the method in downscaling, as well as robustness to the uncertainty of the input data. This is critical to the use of data from public repositories in local analyses, when point observations are not available, as they

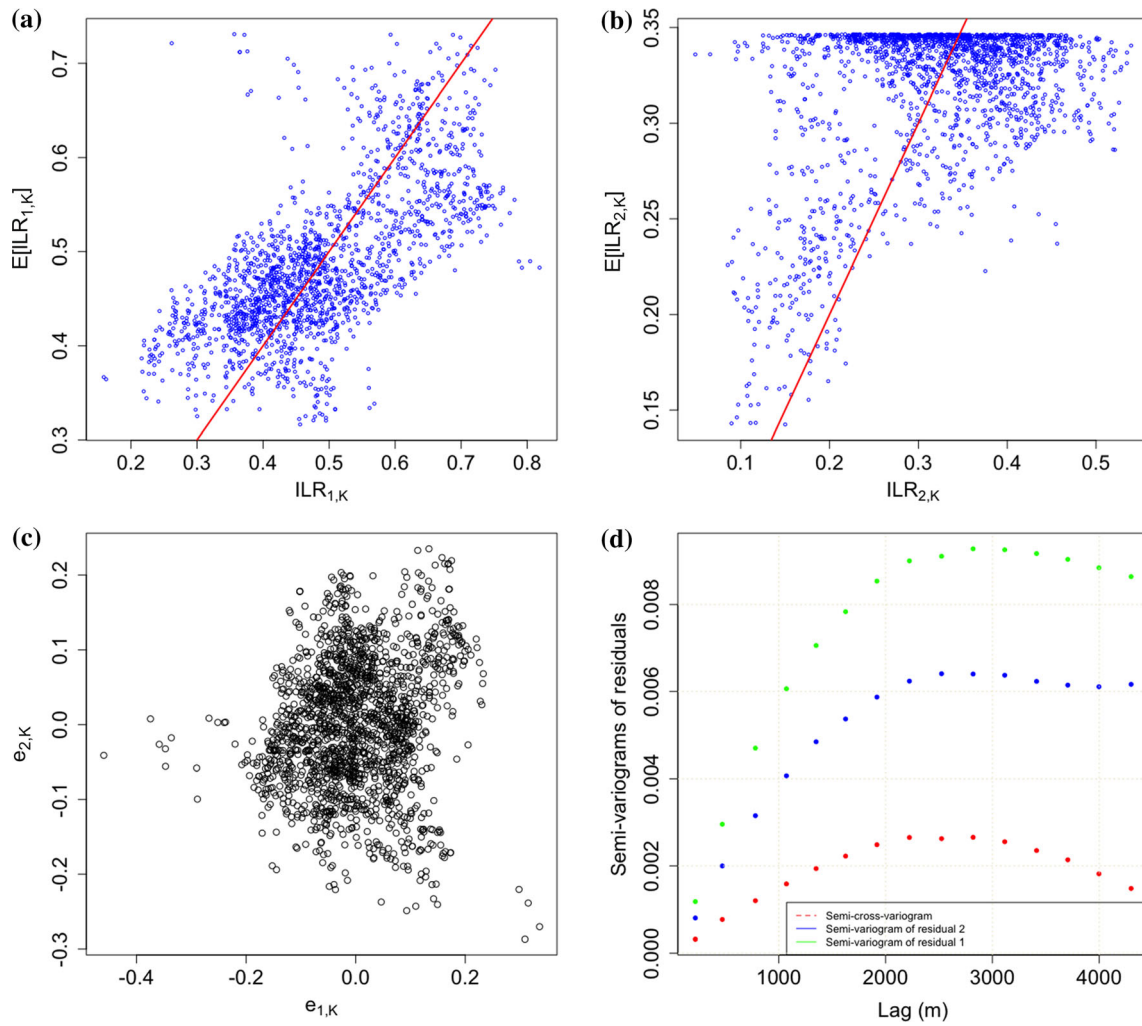


Fig. 10 **a, b** Scatter plots of the observed values and the fitted values of the regression model. In red, the line of equation $\mathbb{E}[\widehat{\text{ILR}}_{i,K}] = \text{ILR}_{i,K}, i = 1, 2$. The Pearson coefficient is 0.67 for

$\mathbb{E}[\widehat{\text{ILR}}_{1,K}], \text{ILR}_{1,K}$ and 0.43 for $\mathbb{E}[\widehat{\text{ILR}}_{2,K}], \text{ILR}_{2,K}$. **c** Scatter plots of the coarse residuals $e_{1,K}$ and $e_{2,K}$. **d** Empirical semi-variograms and cross-variograms of coarse residuals $e_{1,K}$ and $e_{2,K}$

are naturally prone to uncertainty at a fine scale. While a full account of SoilGrids uncertainty will be the scope of future work, a relevant feature of ILR-ATPRCoK method—similarly as ATPRCoK in the Euclidean setting—is the possibility to easily incorporate point observations collected at the site, thus anchoring the downscaled

maps (either kriged or simulated) to such observations (Park 2013). For instance, at the time of writing, a campaign of data acquisition is under way in the Caldane basin, and will support the definition of (possibly improved) random psfs maps, to be used as input to the SMART-SED model discussed in Sect. 6.

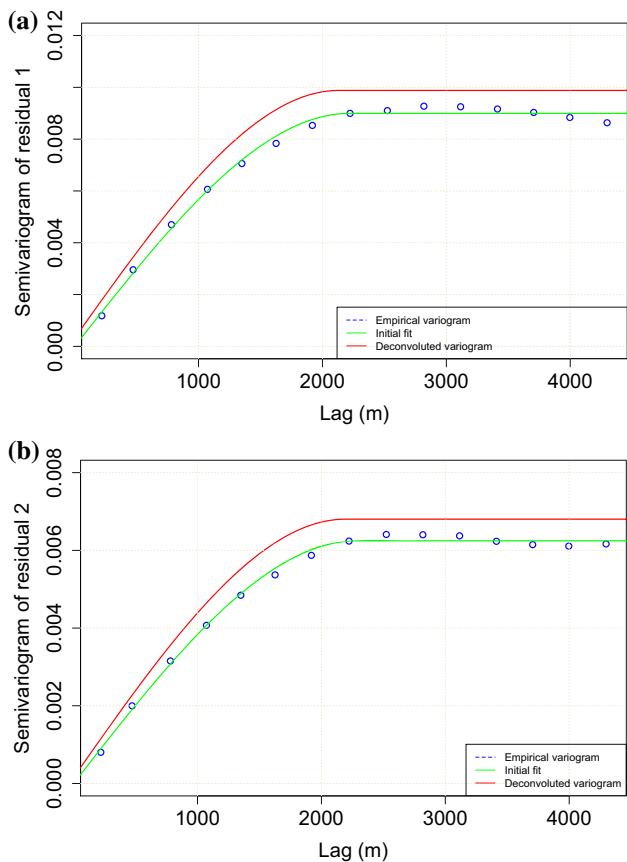


Fig. 11 Results of the Goovaerts’ deconvolution (red line) procedure applied to the empirical variograms (blue dots) starting from an initial fit (green line). The empirical variograms are fitted to a spherical variogram model. Fitted models: **a** Sill: 0.00956, Range: 2130 m, Nugget: 0.00032; **b** Sill: 0.00665, Range: 2190 m, Nugget: 0.00016

Acknowledgements The authors gratefully acknowledge the financial support of Fondazione Cariplo in the framework of the SMART-SED project, Grant Number 2017-0722.

Funding Open access funding provided by Politecnico di Milano within the CRUI-CARE Agreement.

Compliance with ethical standards

Conflict of interest The authors declare that they have no conflict of interest.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Proof of ILR-ATPRCoK proposition

In the following we propose a proof of Proposition 1, i.e. the equivalence of the ILR-ATPRCoK predictor (in the simplex \mathbb{S}^p) to the classical ATPRCoK predictor (in the space \mathbb{R}^{p-1}) by applying an isometric isomorphism. The equivalence must be intended in the predictor, unbiasedness and optimality conditions. In the following we make extensive use of the ILR properties defined in Pawlowsky-Glahn et al. (2015b), pp. 37–43, and of the perturbation-linear combination of compositions (a matrix product), defined as follows (Pawlowsky-Glahn et al. 2015b, pp. 54–55). If we consider a column vector $\mathbf{y} \in \mathbb{R}^{p-1}$ and a matrix Ψ such that each row belongs to \mathbb{S}^p , then

$$\mathbf{y}' \odot \Psi = \mathcal{C} \left[\prod_{i=1}^{p-1} \psi_{i,1}^{y_i}, \dots, \prod_{i=1}^{p-1} \psi_{i,p}^{y_i} \right].$$

We shall also use the fact that $ILR : \mathbb{S}^p \rightarrow \mathbb{R}^{p-1}$ extract the Fourier coordinates of a basis projection for the vector $\mathbf{z} \in \mathbb{S}^p$, i.e.,

$$\mathbf{y} = ILR(\mathbf{z})$$

$$\mathbf{z} = \bigoplus_{i=1}^{p-1} \langle \mathbf{z}, \psi_i \rangle_a \odot \psi_i = (\mathbf{y}' \odot \Psi)' = ILR^{-1}(\mathbf{y})$$

where the rows of $\Psi = [\psi_{ij}]_{i=1, \dots, p-1}^{j=1, \dots, p}$ are (compositional) vectors identifying an orthonormal basis of the simplex $\{\psi_1, \dots, \psi_{p-1}\}$ and $\mathbf{y} = [y_i]_{i=1, \dots, p-1} \in \mathbb{R}^{p-1}$ is the vector of coordinates (i.e. of the Fourier coefficients) of the identified basis of the simplex.

Let us start with the predictor, applying the ILR to the ATPRCoK predictor defined in the Aitchison space \mathbb{S}^p (14), we get

$$\hat{\mathbf{Y}}_k = \sum_l u_l^k \beta_l^{\mathbf{Y}} + \sum_K ILR(\Lambda_K \square((\mathbf{e}_K^{\mathbf{Y}})' \odot \Psi)').$$

where $\beta^l = ILR^{-1}(\beta_l^{\mathbf{Y}})$ and $\mathbf{e}_K = ((\mathbf{e}_K^{\mathbf{Y}})' \odot \Psi)'$. Being $e_{K,i}^{\mathbf{Y}}$, the i -th element of the vector $\mathbf{e}_K^{\mathbf{Y}}$, $i = 1, \dots, p - 1$, we have

$$\begin{aligned} & ILR(\Lambda_K \square((\mathbf{e}_K^{\mathbf{Y}})' \odot \Psi)') \\ &= ILR(((\mathbf{e}_K^{\mathbf{Y}})' \odot \Psi \square \Lambda_K')') \\ &= ILR(((\mathbf{e}_K^{\mathbf{Y}})' \odot [\psi_i' \square \Lambda_K]_{i=1, \dots, p-1})') \\ &= ILR\left(\left(\bigoplus_{i=1}^{p-1} e_{K,i}^{\mathbf{Y}} \odot \psi_i' \square \Lambda_K'\right)'\right) \\ &= \left(\sum_{i=1}^{p-1} e_{K,i}^{\mathbf{Y}} ILR(\psi_i' \square \Lambda_K')\right)' \\ &= ((\mathbf{e}_K^{\mathbf{Y}})' [ILR(\psi_i' \square \Lambda_K')]_{i=1, \dots, p-1})' \\ &= [\langle \Lambda_K \square \psi_i, \psi_j \rangle_{a, i,j=1, \dots, p-1}] \mathbf{e}_K^{\mathbf{Y}} = \Lambda_K^{\mathbf{Y}} \mathbf{e}_K^{\mathbf{Y}}. \end{aligned}$$

Note that the matrix Λ_K represents an endomorphism in

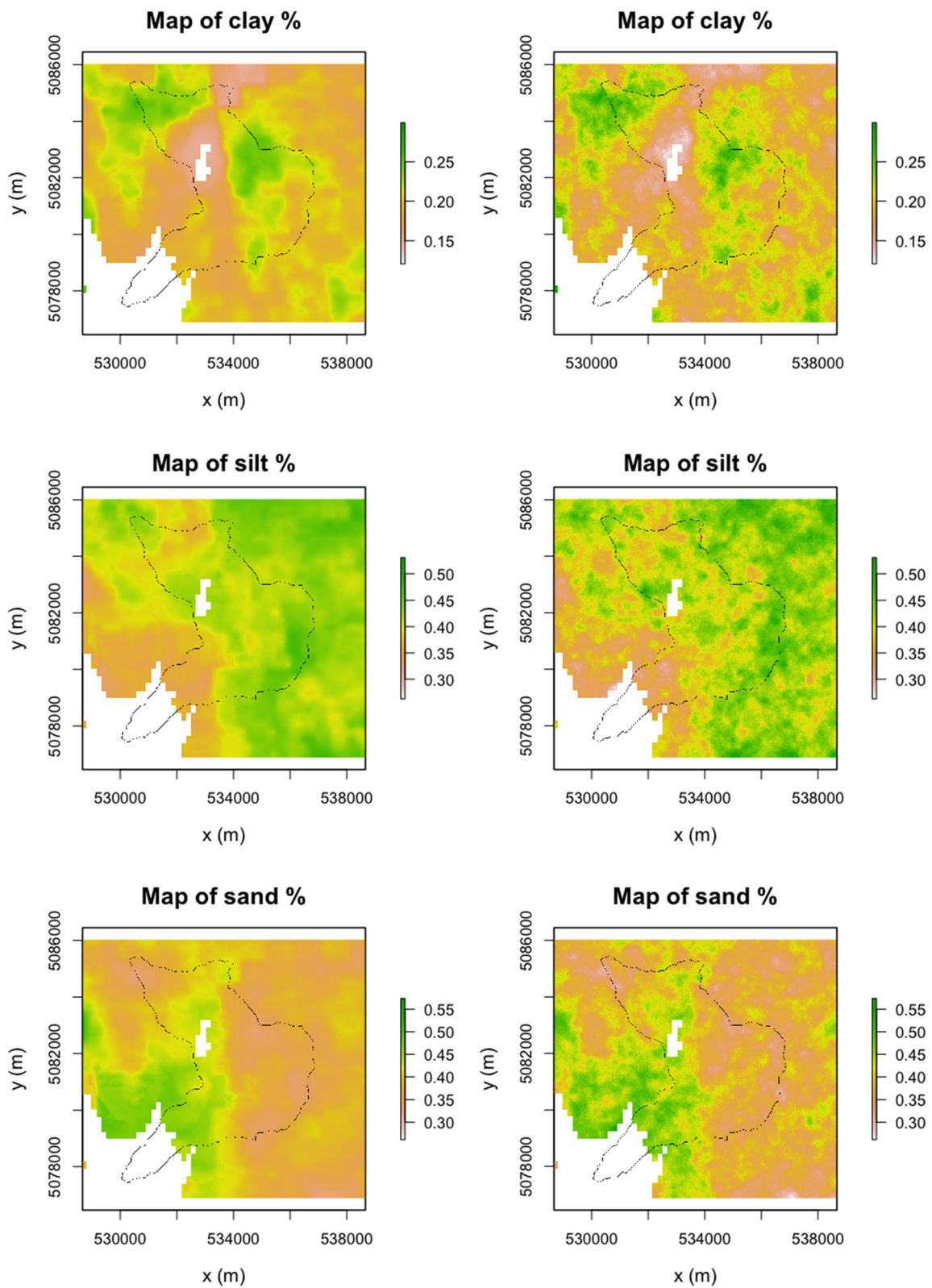


Fig. 12 Clay, silt and sand maps results of the ILR-ATPRCoK on the left and on the right results of BSGS. The black polygon delimits the hydrographic Caldane basin

the simplex, as it can be linked to A_K^Y , through a linear transformation. Indeed, consider the CLR transformation (Pawlowsky-Glahn et al. 2015b, p. 35), defined as

$$CLR(\mathbf{x}) = \left[\frac{\ln(x_i)}{\left(\prod_{j=1}^p x_j\right)^{1/p}} \right]_{i=1,\dots,p},$$

and define the contrast matrix Ψ_c , whose columns are the $CLR(\psi_i), i = 1, \dots, p - 1$ (see also Pawlowsky-Glahn et al. 2015b, p. 36). In this context, A_K^Y is recast as

$$\begin{aligned} A_K^Y &= [\langle A_K \square \psi_i, \psi_j \rangle_a]_{i,j=1,\dots,p-1} \\ &= [\langle A_K CLR(\psi_i), CLR(\psi_j) \rangle]_{i,j=1,\dots,p-1} \\ &= [CLR(\psi_j)' A_K CLR(\psi_i)]_{i,j=1,\dots,p-1} = \Psi_c' A_K \Psi_c \end{aligned}$$

The properties of an endomorphism matrix are listed in Pawlowsky-Glahn et al. (2015b), p. 56.

In this way, we obtain the ATPRCoK predictor in the Euclidean space \mathbb{R}^{p-1} , i.e.,

$$\hat{\mathbf{Y}}_k = \sum_l u_k^l \beta_Y^l + \sum_K A_K^Y \mathbf{e}_K^Y.$$

Regarding the unbiasedness and optimality conditions, calling $\bar{\boldsymbol{\mu}}^Y = ILR(\bar{\boldsymbol{\mu}})$, $d(\cdot, \cdot)$ the Euclidean distance and considering ILR properties, one easily obtains that the same conditions hold in \mathbb{R}^{p-1} ,

$$\mathbb{E} [d_a^2(\bigoplus_K A_K \square \mathbf{e}_K, \mathbf{e}_k)] = \mathbb{E} \left[d^2 \left(\sum_K A_K^Y \mathbf{e}_K^Y, \mathbf{e}_k^Y \right) \right], \quad (20)$$

$$\begin{aligned} ILR(Cen(\bigoplus_K A_K \square \mathbf{e}_K)) \\ = \mathbb{E} [ILR(\bigoplus_K A_K \square \mathbf{e}_K)] = \mathbb{E} \left[\sum_K A_K^Y \mathbf{e}_K^Y \right] = \bar{\boldsymbol{\mu}}^Y. \end{aligned} \quad (21)$$

The first equality (20) (i.e. the optimality condition) derives from the property of the ILR to preserve distances (as opposed e.g. to ALR); in this way the optimality condition holds in \mathbb{R}^{p-1} .

Finally from an “energy” point of view the two formulation are equivalent. Indeed if we consider the quadratic form associated with the covariance structure $C_a(\mathbf{x}_1, \mathbf{x}_2)$, $\mathbf{x}_1, \mathbf{x}_2 \in D$, denote by \mathbf{z} a non-random element of \mathbb{S}^p , and use the ILR properties, we obtain

$$\begin{aligned} \xi &= \langle \mathbf{Z}(\mathbf{x}_1) \ominus \boldsymbol{\mu}(\mathbf{x}_1), \mathbf{z} \rangle_a = \langle ILR(\mathbf{Z}(\mathbf{x}_1) \ominus \boldsymbol{\mu}(\mathbf{x}_1)), \mathbf{y} \rangle = \\ &= \langle \mathbf{Y}(\mathbf{x}_1) - \boldsymbol{\mu}_Y(\mathbf{x}_1), \mathbf{y} \rangle. \end{aligned}$$

Using the latter expression, one has

$$\begin{aligned} \langle C_a(\mathbf{x}_1, \mathbf{x}_2) \mathbf{z}, \mathbf{z} \rangle_a &= \langle ILR(C_a(\mathbf{x}_1, \mathbf{x}_2) \mathbf{z}), \mathbf{y} \rangle \\ &= \langle \mathbb{E}[\xi ILR(\mathbf{Z}(\mathbf{x}_2) \ominus \boldsymbol{\mu}(\mathbf{x}_2)), \mathbf{y}] \rangle \\ &= \langle \mathbb{E}[\xi (\mathbf{Y}(\mathbf{x}_2) - \boldsymbol{\mu}_Y(\mathbf{x}_2))], \mathbf{y} \rangle \\ &= \langle C_Y(\mathbf{x}_1, \mathbf{x}_2) \mathbf{y}, \mathbf{y} \rangle. \end{aligned}$$

Hence, the covariance structure in the Euclidean space \mathbb{R}^{p-1} , reads, $C_Y(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}[\langle \mathbf{Y}(\mathbf{x}_1) - \boldsymbol{\mu}_Y(\mathbf{x}_1), \mathbf{y} \rangle \langle \mathbf{Y}(\mathbf{x}_2) - \boldsymbol{\mu}_Y(\mathbf{x}_2) \rangle]$. This means that the knowledge of the covariance structure in the Aitchison simplex C_a implies the knowledge of the covariance structure in the Euclidean space C_Y and viceversa. This result, together with the relation stated above among the regression coefficients ($\boldsymbol{\beta}' = ((\boldsymbol{\beta}_Y^l)' \odot \Psi')'$), and the equivalence of predictor, optimality and unbiasedness conditions, are sufficient to prove Proposition 1.

References

Aitchison J (1982) The statistical analysis of compositional data. *J R Stat Soc Ser B* 44(2):139–177

Aitchison J (1986) The statistical analysis of compositional data. Chapman and Hall Ltd, London

Atkinson P (2013) Downscaling in remote sensing. *Int J Appl Earth Obs Geoinf* 22:106–114

Benndorf J (2003) Conditional joint simulation of random fields on block support. MPhil, University of Queensland, Brisbane, p 160

Benndorf J (2004) Large scale stochastic simulations for long term scheduling formulations. MPhil, University of Queensland, Brisbane, p 224

Billheimer D, Guttorp P, Fagan WF (2001) Statistical interpretation of species composition. *J Am Stat Assoc* 96(456):1205–1214

Boogaart KGVD, Tolosana-Delgado R (2008) “compositions”: a unified R package to analyze compositional data. *Comput Geosci* 34(4):320–338

Bosq D (2000) Linear processes in function spaces: theory. Lecture Notes in Statistics and applications. Springer, New York

Brown DG, Goovaerts P, Burnicki A, Li MY (2002) Stochastic simulation of land-cover change using geostatistics and generalized additive models. *Photogramm Eng Remote Sens* 68(10):1051–1062

Buccianti A, Grunsky E (2014) Compositional data analysis in geochemistry: Are we sure to see what really occurs during natural processes? *J Geochem Explor* 141:1–5

Chilès JP, Delfiner P (2012) Geostatistics: modeling spatial uncertainty. Wiley Series in Probability and Statistics. Wiley, New York

Cressie N (1993) Statistics for spatial data. Wiley Series in Probability and Statistics. Wiley, New York

Delbari M, Afrasiab P, Loiskandl W (2011) Geostatistical analysis of soil texture fractions on the field scale. *Soil Water Resour* 6:172–189

Dobarco RM, Orton GT, Arrouays D, Lemercier B, Paroissien JB, Walter C, Saby N (2016) Prediction of soil texture using descriptive statistics and area-to-point kriging in Region Centre (France). *Geoderma Reg* 7(3):279–292

Egozcue JJ, Pawlowsky-Glahn V, Mateu-Figueras G, Barceló-Vidal C (2003) Isometric logratio transformations for compositional data analysis. *Math Geol* 35(3):279–300

Filzmoser P, Hron K, Templ M (2018) Analyzing compositional data using R. In: Applied compositional data analysis. Springer, pp 17–34

Fišerová E, Hron K (2011) On the interpretation of orthonormal coordinates for compositional data. *Math Geosci* 43(4):455–468

- Gatti F, Bonaventura L, Menafoglio A, Papini M, Longoni L (2020) Preliminary results from the SMART-SED basin scale sediment yield model, volume 4: understanding and reducing landslide disaster risk (World landslide Forum 5). Springer Nature Switzerland AG
- Goovaerts P (2008) Kriging and semivariogram deconvolution in presence of irregular geographical units. *Math Geol* 40:101–128
- Goovaerts P (2010) Combining areal and point data in geostatistical interpolation: applications to soil science and medical geography. *Math Geosci* 42:535–554
- Gräler B, Pebesma E, Heuvelink G (2016) Spatio-temporal interpolation using gstat. *RFID J* 8:204–218
- Groenendyk D, Ferré T, Thorp K, Rice A (2015) Hydrologic-process-based soil texture classifications for improved visualization of landscape function. *PLoS One* 10:e0131299
- Hengl T, Heuvelink GBM, Rossiter DG (2007) About regression-kriging: from equations to case studies. *Comput Geosci* 33(10):1301–1315
- Hengl T, De Jesus JM, MacMillan RA, Batjes NH, Heuvelink GBM, Ribeiro E, Samuel-Rosa A, Kempen B, Leenaars JGB, Walsh MG, Gonzalez MR (2014) Soilgrids1km: global soil information based on automated mapping. *PLoS One* 9(8):1–17
- Hengl T, De Jesus JM, Heuvelink GBM, Ruiperez-Gonzalez M, Kilibarda M (2017) Soilgrids250m: global gridded soil information based on machine learning. *PLoS One* 12(2):1–40
- Kalos MH, Whitlock PA (2009) Monte Carlo methods. Wiley, New York
- Kim JH (1999) Spurious correlation between ratios with a common divisor. *Stat Probab Lett* 44(4):383–386
- Kyriakidis P (2004) A geostatistical framework for area-to-point spatial interpolation. *Geograph Anal* 36(3):259–289
- Kyriakidis PC, Yoo EH (2005) Geostatistical prediction and simulation of point values from areal data. *Geograph Anal* 37(2):124–151
- Martín MA, Pachepsky Y, Baez C, Reyes M (2017) On soil textural classifications and soil texture-based estimations. *Solid Earth Discuss* 9:54
- Mateu-Figueras G, Pawlowsky-Glahn V, Egozcue JJ (2011) The principle of working on coordinates, chap 3. Wiley, Chichester, pp 29–42
- Menafoglio A, Guadagnini L, Secchi P (2014) A kriging approach based on aitchison geometry for the characterization of particle-size curves in heterogeneous aquifers. *Stoch Environ Res Risk Assess* 28:1835–1851
- Minasny B, Mcbratney A (2007) Spatial prediction of soil properties using EBLUP with the Matérn covariance function. *Geoderma* 140:324–336
- Park NW (2013) Spatial downscaling of TRMM precipitation using geostatistics and fine scale environmental variables. *Adv Meteorol*. <https://doi.org/10.1155/2013/237126>
- Pawlowsky V (1984) On spurious spatial covariance between variables of constant sum. *Sciences de la terre Informatique géologique* 21:107–113
- Pawlowsky V (1989) Cokriging of regionalized compositions. *Math Geol* 21(5):513–521
- Pawlowsky-Glahn V, Buccianti A (2011) *Compositional data analysis*. Wiley, Chichester
- Pawlowsky-Glahn V, Egozcue JJ (2001) Geometric approach to statistical analysis on the simplex. *Stoch Environ Res Risk Assess* 15(5):384–398
- Pawlowsky-Glahn V, Egozcue JJ (2002) Blu estimators and compositional data. *Math Geol* 34:259–274
- Pawlowsky-Glahn V, Egozcue JJ (2016) Spatial analysis of compositional data: a historical review. *J Geochem Explor* 164:28–32
- Pawlowsky-Glahn V, Olea RA (2004) *Geostatistical analysis of compositional data*, vol 7. Oxford University Press, Oxford
- Pawlowsky-Glahn V, Egozcue JJ, Olea RA, Pardo-Iguzquiza E (2015a) Cokriging of compositional balances including a dimension reduction and retrieval of original units. *J South Afr Inst Min Metall* 115:59–72
- Pawlowsky-Glahn V, Egozcue JJ, Tolosana-Delgado R (2015b) *Modeling and analysis of compositional data*. Wiley, Chichester
- Pebesma EJ (2004) Multivariable geostatistics in S: the gstat package. *Comput Geosci* 30:683–691
- R Development Core Team (2008) R: a language and environment for statistical computing. <http://www.R-project.org>
- Rodríguez-Díaz J, Rivas-Lopez M, Santos-Martin M, Marinas-Collado I (2020) Optimal designs for a linear-model compositional response. *Stoch Environ Res Risk Assess* 34(1):139–148
- Rosso R (2004) Mappatura dell' indice di assorbimento e del massimo volume specifico di ritenzione potenziale del terreno. *Relazione Finale Progetto SHAKEUP-2, ARPA Lombardia*
- Schaefer GL, Cosh MH, Jackson TJ (2007) The USDA natural resources conservation service soil climate analysis network (SCAN). *J Atmos Ocean Technol* 24(12):2073–2077
- Tolosana-Delgado R, van den Boogaart KG, Pawlowsky-Glahn V (2011) Geostatistics for compositions. In: Pawlowsky-Glahn V, Buccianti A (eds) *Compositional data analysis: theory and applications*. Wiley, Chichester, pp 73–86
- Tolosana-Delgado R, Mueller U, Van Den Boogaart KG (2019) Geostatistics for compositional data: an overview. *Math Geosci* 51(4):485–526
- Walvoort DJJ, De Gruijter JJ (2001) Compositional kriging: a spatial interpolation method for compositional data. *Math Geol* 33(8):951–966
- Wang Q, Shi W, Atkinson PM, Zhao Y (2015) Downscaling MODIS images with area-to-point regression kriging. *Remote Sens Environ* 166:191–204
- Xiao M, Zhang G, Breitkopf P, Villon P, Zhang W (2018) Extended co-kriging interpolation method based on multi-fidelity data. *Appl Math Comput* 323:120–131