# The Whole Versus the Parts: The Challenge of Compositional Data Analysis (CoDA) Methods for Geochemistry

Antonella Buccianti and Caterina Gozzi

**Abstract** In complex geochemical aqueous systems, chemical species are conceptually distinct but empirically related thanks to a large number of interactions taking place at different spatial and/or temporal scales. In this condition, common elements are shared, multiplicative interactions arise, and feedback mechanisms may be able to maintain the system far from the thermodynamical equilibrium, bearing wide fluctuations. Chemical species can have alternative stable states and transitions among them that could produce important consequences for the stability and the resilience of the solutions, also forced by climate changes and with impact on human health. Under the Compositional Data Analysis (CoDA) methodology, it is possible to appreciate the power of some tools able to take a look at the whole instead of the constituting parts, enhancing the understanding of the nature of mutual interactions. In this research work, the role of the perturbation operator governing addition/subtraction in the simplex geometry is explored as a way to trace compositional changes and investigate the system dynamics. The results of our approach on the chemistry of the Arno River waters (Central Italy) highlight the possibility to discover the resilience of chemical species under the pressure of the environmental drivers affecting the catchment. Geochemical mobility (e.g., ionic potential, ionic strength) can be associated with new tools that provide information on either the resistance to change, predisposing the system to critical shifts, or its adaptive capacity, which instead favors gradual changes. This information appears to be fundamental since river water chemistry enables to decipher processes at the boundaries among lithosphere, biosphere, hydrosphere, and atmosphere, all key reservoirs involved in the dynamics of the Earth. This knowledge will be particularly relevant if the pressure of the climatic changes on our planet will continue to increase.

A. Buccianti (✉) · C. Gozzi
Department of Earth Sciences, University of Florence, Florence, Italy
e-mail: antonella.buccianti@unifi.it

C. Gozzi
e-mail: caterina.gozzi@unifi.it

# 1  Introduction

Compositional Data Analysis (CoDA) has been considered, at first, only as a way to open constrained data with the aim of working in the correct sample space. At the beginning, this idea has generated some misunderstanding since data that do not close to a given constant were assumed not to be compositional (Aitchison 1982). Subsequently, the idea of proportionality among the terms of a composition has gained ground and it was clear that the total to which a composition is closed is irrelevant (Aitchison 1986). What really counts in CoDA are the reciprocal relationships among the parts of the composition that are intimately linked to each other (Egozcue et al. 2011). This framework appears to well describe simultaneous chemical reactions that occur in natural processes. A geochemical system comprises, at least, an aqueous solution in which the species of many elements are dissolved, include one or more minerals, a buffering condition with a gas reservoir, and the atmosphere in the surficial environment (Bethke 2008). Water geochemists work daily with equations that describe the equilibrium of several simultaneous chemical reactions among dissolved species, minerals, and gases. Each mass action equation relates the activities of the species to the reaction equilibrium constants. How can we describe the state of such a system? A direct approach would be to write the single reactions among the species of the system, minerals, and gases. To solve equilibrium problems a set of concentrations that simultaneously satisfy the mass action equation should be written for each possible reaction. It is clear that such a system is multicomponent and due to the relationships among the parts and the presence of feedback mechanisms, its dynamics become complex. Reactions and concentrations cannot be decoupled and the reaction rates, relating sink and source terms, can induce nonlinearity (Sanchez-Vila et al. 2007). The evolution of an aqueous system is characterized by the migration of the chemical species between phases in innumerable cycles that form all together a larger, more complex interconnected whole (Kleidon 2010). The interconnections generate properties of the whole that cannot be wholly understood by examining the parts of the system in isolation. Thus, complex systems have properties that depend on the integrity of the whole. In this context, most of the geological systems are complex systems since they are dissipative and self-organizing systems, working in open and dynamic conditions and consisting of a great number of components (e.g., rocks, minerals, and elements) which interact in a nonlinear region far from the equilibrium (Shvartsev 2009). A way to take into account the simultaneous relationships among all the components of a system is to apply methods of multivariate analysis (Krzanowski 2000). If this path is followed, it is necessary to consider the appropriateness of the sample space in which data move, since it determines the variance-covariance structure. Detecting in a correct manner how variability moves is fundamental to intercept the dynamics of a complex system and to understand the response to environmental perturbations. In these terms, the CoDA approach appears to be the theoretical framework to abandon single variables and to focus our attention on the properties of the whole. The treatment of compositional data requires two possible strategies (Aitchison 1986). The first one,

called "stay in the simplex" approach, adopts the Aitchison geometry to work in a constrained sample space. The second one is based on the transformation of the data to move cases out from the simplex in the real Euclidean space and it is known as "working in coordinates" (Egozcue et al. 2003; Pawlowsky-Glahn et al. 2015). In this work, our interest focused on the first approach and, in particular, on the use of the perturbation operator. It represents one of the basic tools required to give to the simplex a vector space structure and it could have a strategic role in monitoring changes, giving useful information about the dynamics of a system. An application example to the waters of the Arno (Central Italy) riverine geochemical system will be used to illustrate the procedure. The research question is whether the perturbation operator is able to provide insight about similarity among chemical elements and the coherence in their behavior within the dynamics of the system, thus offering a new geochemical tool of investigation.

## 2 Material and Methods

### 2.1 Tracing Change Through Perturbations

The vector $\mathbf{x} = (x_1, x_2, \ldots, x_D)$ indicates a $D$-part composition in the simplex:

$$S^D = \left[ \mathbf{x} = (x_1, x_2, \ldots, x_D) : x_i > 0 (i = 1, 2, \ldots, D), \sum_{i=1}^{D} x_i = \kappa \right] \quad (1)$$

where $\kappa$ is a given positive constant whose value depends on the measure unit. Two basic operations on the simplex, called perturbation and powering that are able to induce a real vector space structure on the simplex, have been introduced (Aitchison 1986). Furthermore, the introduction of an inner product, with its associated norm and distance, has been used to obtain a $D - 1$-dimensional Hilbert space structure (Billheimer et al. 2001; Pawlowsky-Glahn and Egozcue 2001). In this framework, a perturbation $\mathbf{p} = (p_1, p_2, \ldots, p_D)$ is a differential scaling operator that, when applied to the composition $\mathbf{x} = (x_1, x_2, \ldots, x_D)$, yields the composition:

$$\mathbf{y} = \mathbf{p} \oplus \mathbf{x} = C(p_1 x_1, \ldots, p_D x_D) \quad (2)$$

where $C$ is the closure operator that scales elements to remain in the simplex sample space:

$$C(\mathbf{x}) = \left[ \frac{\kappa \cdot x_1}{\sum_{i=1}^{D} x_i}, \frac{\kappa \cdot x_2}{\sum_{i=1}^{D} x_i}, \ldots, \frac{\kappa \cdot x_D}{\sum_{i=1}^{D} x_i} \right]. \quad (3)$$
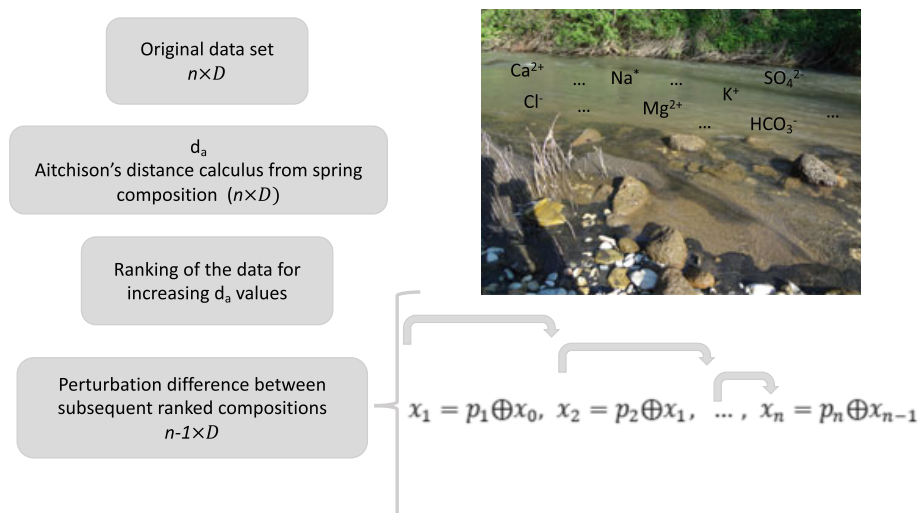
**Fig. 1** Steps of the procedure to rank the cases by increasing Aichison's distance from the spring composition and to calculate the perturbation difference for subsequent compositions, $x_{i+1} - x_i$

As any composition can be expressed as a result of a perturbation on any other composition, the operator acquires a fundamental role in tracing compositional changes. This is especially true if a reference composition is available and differences with respect to this one can be calculated or when compositions are linked each other by a subsequent evolutive pattern. Thus, if $\mathbf{y} = \mathbf{p} \oplus \mathbf{x}$ corresponds to addition in the $\mathbb{R}$, the $\mathbf{y} = \mathbf{p} \ominus \mathbf{x}$ corresponds to the perturbation difference when obtained by a component-wise division of the elements of the $\mathbf{x}$ and $\mathbf{y}$ vectors (Aitchison 1986). Our approach is based on the transformation of the original dataset $n \times D$, related to the chemistry of waters collected in a river, in a perturbations matrix $n - 1 \times D$. In this matrix, each row is related to the perturbation difference among two subsequent compositions, corresponding to the difference $x_{i+1} - x_i$, after having ranked cases under some geochemical hypothesis. Thus, in order to calculate the perturbation matrix, the original dataset was first ranked by considering the increasing value of the Aitchison distance from the chemical composition of the spring, to be considered as a pristine water. The procedure is illustrated in Fig. 1. The (squared) Aitchison distance is a simplicial metric given by

$$d_a^2(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{D} \left[ \ln \frac{x_i}{g_m(\mathbf{x})} - \ln \frac{y_i}{g_m(\mathbf{y})} \right]^2 \tag{4}$$

where $g_m(\cdot)$ is the geometric mean of the components, calculated considering the parts of the vectors $\mathbf{x}$ and $\mathbf{y}$ (Pawlowsky-Glahn and Egozcue 2001). The explorative map of the Aitchison distance could show where the compositional difference with respect to the spring is higher at the catchment scale. This permits the identification of homogeneous areas or spatial patterns to be related to other non-compositional environmental factors. On the other hand, the perturbation matrix informs, for each

sample, about the chemical species that mostly take charge of the change offering an interesting tool to trace and interpret the behavior of the composition as a whole.

## 2.2 The Dataset

The catchment area of the Arno River Basin is entirely located in Tuscany, Central Italy and has a surface of 8,228 km$^2$ with an average elevation of 353 m. The river's headwater spring is in the Northern Apennines at 1650 m, and flows for about 242 km toward the Ligurian Sea, 10 km West of Pisa and 110 km of Florence, respectively. The drainage network follows NW-SE trending tectonic structures that form six main sub-basins from East to West: (1) Casentino, (2) Chiana Valley, (3) Sieve, (4) Upper
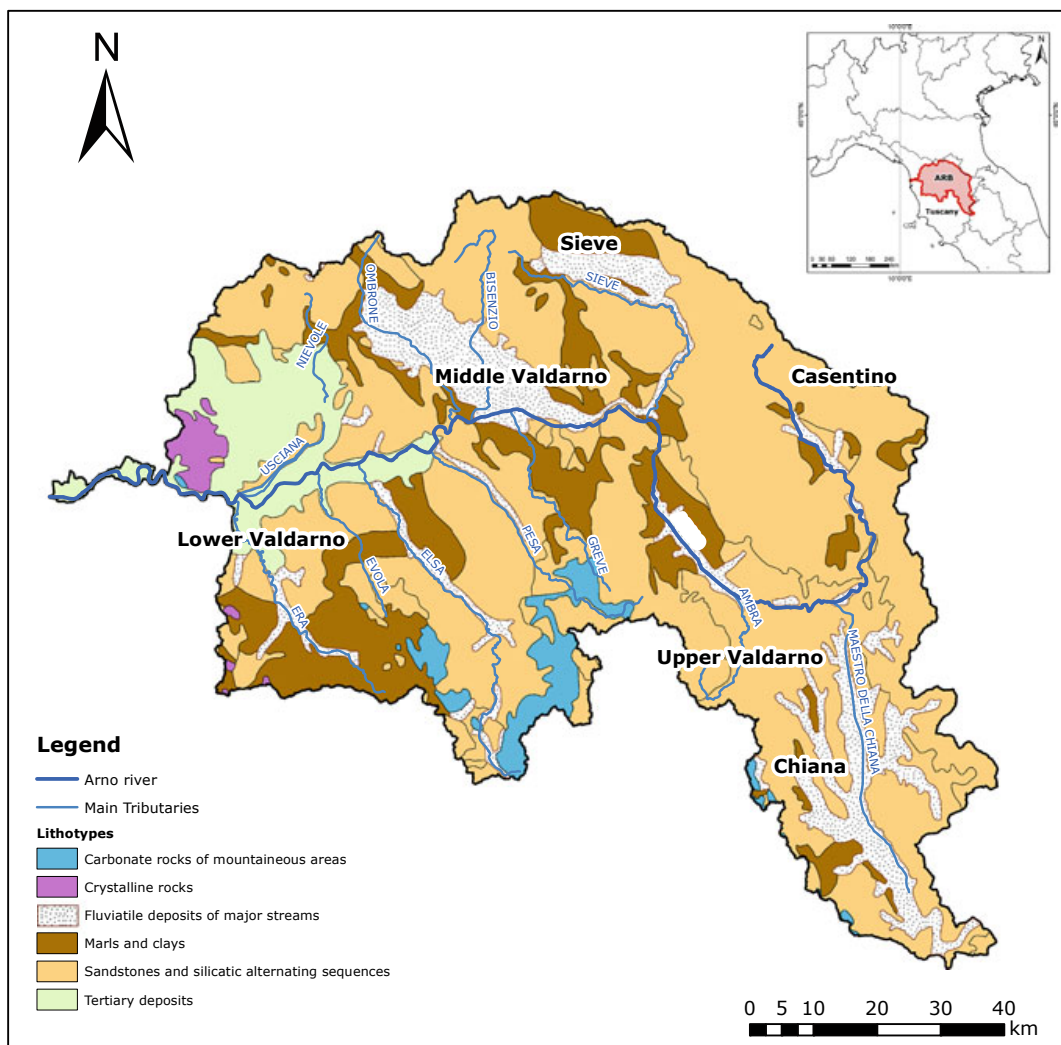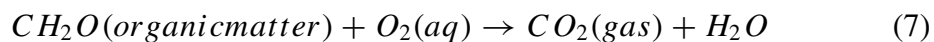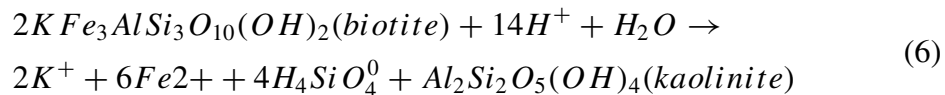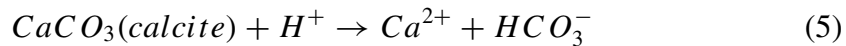


**Fig. 2** Lithological map of the Arno River catchment (Tuscany, central Italy). Geological layer modified from SINAnet (ISPRA Ambiente 2019)

Valdarno, (5) Middle Valdarno, and (6) Lower Valdarno. The outcropping rocks are predominantly sedimentary folded and faulted Mesozoic and Tertiary units derived from the weathering of the Apennine Mountains. A schematic lithological map is shown in Fig. 2. The database is related to a 2-year sampling campaign (2002–2003) performed by taking into account the seasonal variability. However, no significant influence on compositional changes compared to other environmental factors (Nisi et al. 2008) appears to be attributable to seasonality. The main determined chemical species were $Na^+$, $K^+$, $Ca^{2+}$, $Mg^{2+}$, $Cl^-$, $HCO_3^-$, $SO_4^{2-}$, and $SiO_{2(aq)}$. The Arno River Basin, similarly to other European watersheds (Berner and Berner 1996), has suffered from the past century an increased industrial and agricultural development, so that the contribution of chemical weathering is mixed with anthropic inputs (Arrighi et al. 2018). Further details about sampling and analytical methodologies can be found in Nisi et al. (2008).

## 3   Results and Discussion

The analyzed dataset is given by 474 cases and 8 variables ($Na^+$, $K^+$, $Ca^{2+}$, $Mg^{2+}$, $Cl^-$, $HCO_3^-$, $SO_4^{2-}$, $SiO_{2(aq)}$) measured in mg/L. The chemical composition of the spring is characterized by a low conductivity of about 0.12 mS/cm and a TDS (total dissolved solids) equal to 315 mg/L. Values higher than 6.0 mS/cm are related to polluted areas of the catchment, while the maximum value of about 27.66 mS/cm is measured near the river mouth where seawater interacts with the fresh one (TDS equal to 8700 mg/L). Aitchison's distance values of each sample from the composition of the spring, considered representative of a pristine water, have been calculated. The hypothesis is that, starting from this point, weathering processes begin to modify the chemistry of the solution. The spring water with few dissolved constituents interacts with atmospheric air, minerals, and solid organic matter generating a disequilibrium condition. In this situation, dissolution reactions occur and new components are added to the water modifying its composition step by step. Typical reactions that can develop are given by

$$CaCO_3(calcite) + H^+ \rightarrow Ca^{2+} + HCO_3^- \tag{5}$$

$$2KFe_3AlSi_3O_{10}(OH)_2(biotite) + 14H^+ + H_2O \rightarrow$$
$$2K^+ + 6Fe2+ + 4H_4SiO_4^0 + Al_2Si_2O_5(OH)_4(kaolinite) \tag{6}$$

$$CH_2O(organic matter) + O_2(aq) \rightarrow CO_2(gas) + H_2O \tag{7}$$

representing the weathering of carbonate, silicate minerals, and the oxidation of the organic matter, respectively. The effectiveness of Aitchison's distance in capturing what is occurring during weathering processes is revealed by its clear positive rela-
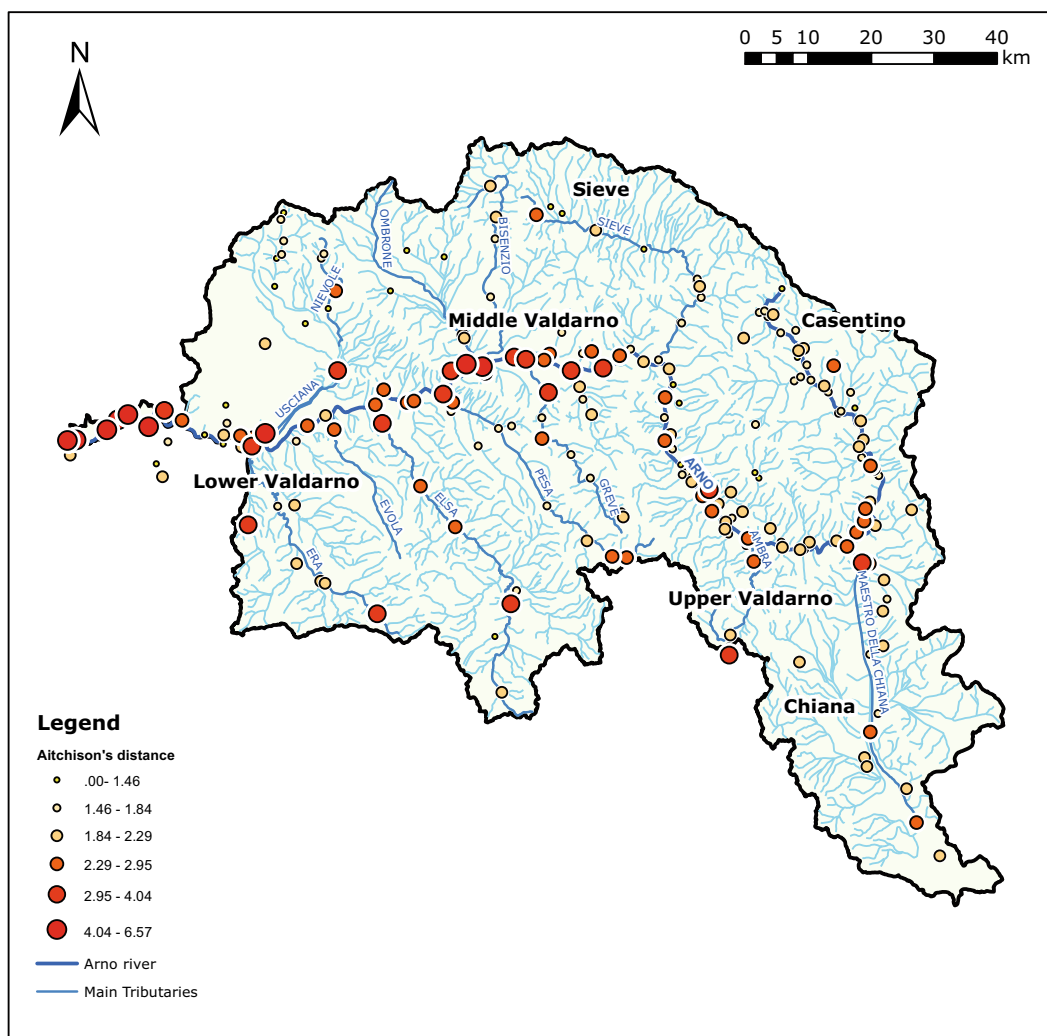
**Fig. 3** Map of Aitchison's distance values of each sample from the spring water considered representative of a pristine condition

tionship with the conductivity, thus representing a possible marker of geochemical changes. The spatial distribution of the distance values is reported in Fig. 3. As we can see from the map, its values increase along the course of the main river mainly starting from the Middle Valdarno where natural processes began to mix with anthropic ones and near to the mouth due to the marine ingression and mixing processes.

After the calculus of Aitchison's distance, the data were ranked for its increasing values and the perturbation difference between subsequent compositions $x_{i+1} - x_i$ was determined. This step is important to highlight the different behavior of single chemical components inside the composition. In this way, similar Aitchison's distance values could be attributed to different processes depending on the association of variables in the perturbation. The results, as unclosed perturbation factors, are reported in Fig. 4 (Egozcue and Pawlowsky-Glahn 2011) while the histogram of Aitchison's distance values, with the associated kernel density estimation curve,
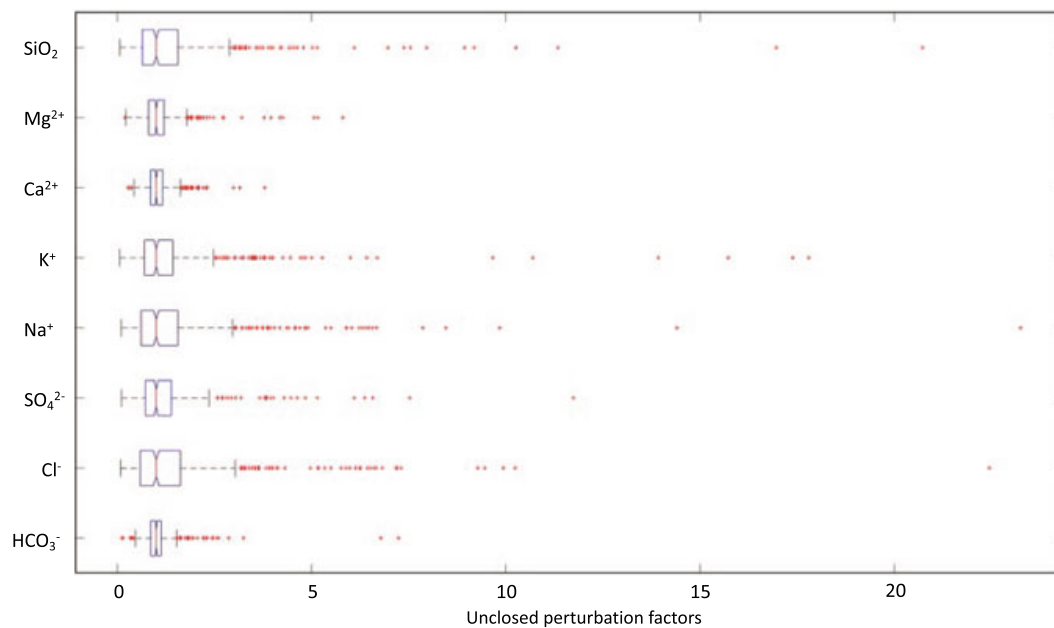
**Fig. 4** Comparative box plots for the unclosed perturbation factors for each chemical species for data ranked by considering the increasing values of Aitchison's distance from the spring water. Perturbation difference represents the difference $x_{i+1} - x_i$
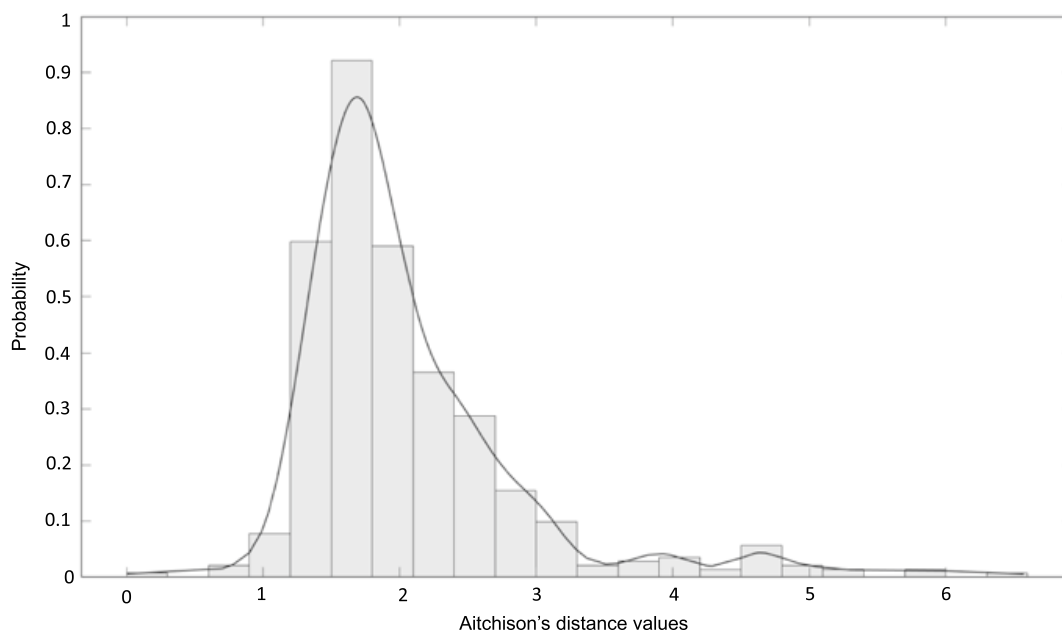


**Fig. 5** Histogram and kernel density estimation of Aitchison's distance values from the spring composition considered as a pristine water
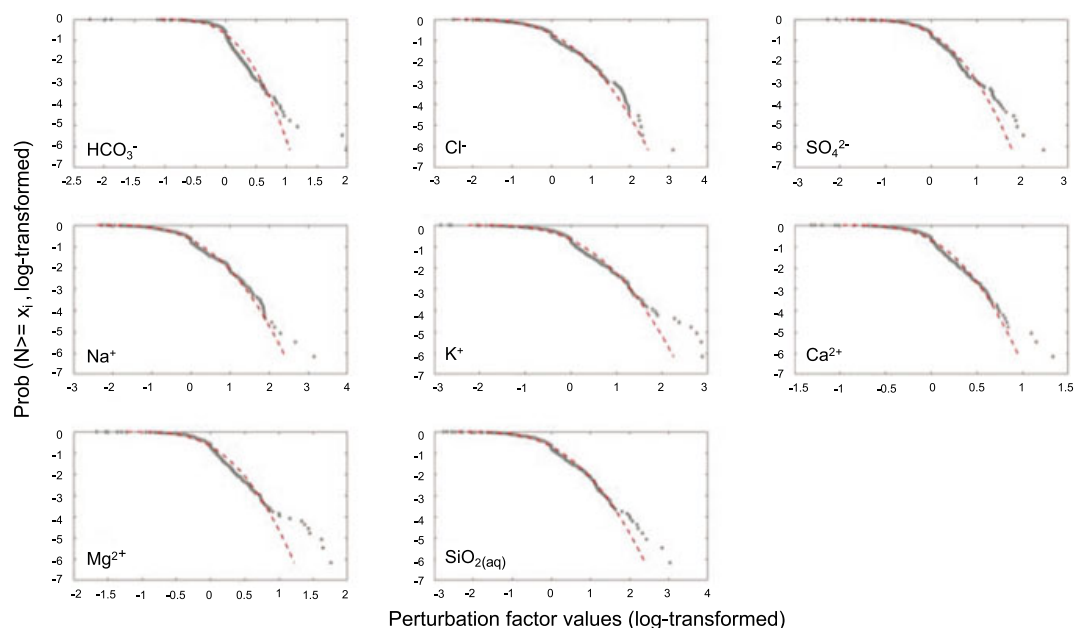
**Fig. 6** Complementary cumulative distribution function for unclosed perturbation factors. The continuous line is related to the log-normal model

is shown in Fig. 5. The distance presents a positive asymmetrical distribution with rare high values in the right tail, indicating that most of the data are characterized by a distance lower than 3 with respect to the pristine condition. High values are related to the mouth (marine ingression) and to peculiar pollution conditions (e.g., Usciana and Chiana Channels) or to the presence of lenses of $CaSO_4$ and $NaCl$ in clays in the valleys of Elsa and Era Tributaries (Fig. 2). The unclosed perturbation factors of Fig. 4 provide information about the contribution of each variable to the compositional changes. $SiO_{2(aq)}$, $Cl^-$, $Na^+$, $K^+$, and $SO_4^{2-}$ show a more scattered contribution while $HCO_3^-$, $Ca^{2+}$, and $Mg^{2+}$ exhibit a more stable signal. Hence, this procedure appears to be able to divide the variables into two different groups, one more sensible to environmental changes, perhaps related to an intermittent spatial behavior, and the other more resilient. The result indicates that monitoring plans for a long temporal range, with the aim of intercepting important changes in the catchment, also related to climatic variations, should be based on the sub-composition related to the cycle of carbonates ($HCO_3^-$, $Ca^{2+}$, and $Mg^{2+}$). To investigate in more detail the nature of the contribution of the variables to the compositional changes, the complementary cumulative distribution function for the unclosed perturbation factors was determined and plotted (Mitzenmacher 2004; Egozcue and Pawlowsky-Glahn 2011). Results are reported in Fig. 6 with the continuous line related to the description of the patterns by the log-normal model. The comparison between experimental data and the log-normal distribution indicates that in no case it is possible to adopt it to describe the behavior of the chemical species. On the other hand, some portions of the curves are clearly linear highlighting the possible presence of more than one power law or a multifractal behavior (Seely and Macklem 2012).

The consequence is that an interaction-dominant dynamics is no longer exhaustive and the typical multiplicative processes of a log-normal model would be associated with interdependent feedback transactions covering different time or space scales (van Rooij et al. 2013). Thus, compositional changes that enrich riverine water in solutes starting from its pristine condition are apparently described by features that are typical of a complex system such as the occurrence of multifractality. This property characterizes dynamical systems in which energy dissipation can no longer be neglected. However, in this general framework, considering the unclosed perturbation factors of Fig. 4, species such as $HCO_3^-$, $Ca^{2+}$, and $Mg^{2+}$ appear to be more efficient in stabilizing fluctuations. In fact, they reveal the presence of a sub-system characterized by homogeneity and connectivity which is particularly resistant to change but also subjected to critical transitions (Scheffer et al. 2012; Sauro Graziano et al. 2020). The cause–effect relationship able to govern this situation could be related at the catchment scale to: (i) the spatial diffusion of the lithologies from which these chemical species migrate; (ii) their homogeneity and connections; and (iii) the efficiency of geochemical processes, affecting the way in which a system, with potential local alternative states, can respond to changing conditions. The behavior of all the other chemical species, $SiO_{2(aq)}$, $Cl^-$, $Na^+$, $K^+$, and $SO_4^{2-}$, could instead indicate: (i) an heterogeneous distribution of the lithologies from which they came from; (ii) incomplete connectivity with presence of modularity; and (iii) less diffused (intermittency) or less efficient geochemical processes, all features that lead to a higher adaptive capacity and favor gradual changes. The values of the saturation indices, calculated for some fundamental minerals such as calcite, dolomite, and quartz, appear stationary when cases are ranked for increasing Aitchison's distance as well as box plots equilibrated around the median values, with the exception of some anomalous data (Zhu and Anderson 2002). Thus, chemical equilibria, able to limit the concentration of the ions in water, have a buffering effect in the compositional changes as registered by the perturbation difference from a pristine water. The result indicates this compositional tool to be very useful as an explorative aid in water geochemistry to check for different evolutive geochemical hypothesis. Therefore it is important to stress that Aitchison's distance is a compositional distance between different compositions and only the association with the perturbation difference will help us to discriminate different situations characterized by similar distance values. In our case study, the distance values are not associated with an anomalous behavior of the perturbation factors (e.g., presence of trends and pluri-modality), indicating that similar processes appear to have affected the evolutive path of the water starting from the pristine reference.

## 4  Conclusions and Future Developments

The research on complexity theory and nonlinear dynamics involves concepts such as dissipative structures and fractal patterning, as well as instability, resilience, adaptive cycles, and uncertainty. Such promising concepts are still developing but CoDA

appears to offer the tools to investigate chemical compositions as a whole and to probe the system dynamics based on the nature of the interactions of the components under the forcing effects of different environmental drivers. In this framework, the perturbation operator seems to present powerful features able to discover the resiliency of chemical variables versus intermittency and instability. This opens important paths to reveal warning signals for relevant changes, also related to climatic variations. The results obtained for the river chemistry of the Arno catchment are very comforting. Further investigation would be performed to link perturbation differences to non-compositional variables, or environmental drivers, characterizing the catchment (e.g. slope, structure of the drainage network, erodibility, runoff, and discharge). The joint analysis would help to point out the nature of the interactions and feedback mechanisms able to govern the development of chemical reactions, thus determining the chemistry of the waters (Gozzi et al. 2019a, b). The proposed approach gives a potential basis to expand this research since Aitchison's distance and compositional changes from a pristine water synthesize the complexity of the system from the chemical point of view. Other non-compositional variables can be then called into question and interrelated with the previous ones.

# References

J. Aitchison, *The Statistical Analysis of Compositional Data*. Monographs on Statistics and Applied Probability (Chapman & Hall Ltd., London (UK), 1986). (Reprinted in 2003 with additional material by The Blackburn Press), 416 p

J. Aitchison, The statistical analysis of compositional data (with discussion). J. R. Stat. Soc. Ser. B (Stat. Methodol.) **44**, 139–177 (1982)

C. Arrighi, M. Masi, R. Iannelli, Flood risk assessment of environmental pollution hotspots. Environ. Model. Softw. **100**, 1–10 (2018)

E.K. Berner, R.A. Berner, *Global Environmental: Water, Air and Geochemical Cycles* (Prentice-Hall, Upper Saddle River NJ, 1996), 376 p

C.M. Bethke, *Geochemical and Biogeochemical Reaction Modeling* (Cambridge University Press, 2008), 543 p

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

E. Dinelli, G. Corteccci, F. Lucchini, E. Zantedeschi, Source of major and trace elements in the stream sediments of the Arno river catchment (northern Tuscany, Italy). Geochem. J. **39**, 531–545 (2005)

J.J. Egozcue, C. Barceló-Vidal, J.A. Martín-Fernández, E. Jarauta-Bragulat, J.L. Díaz-Barrero, G. Mateu-Figueras, Elements of simplicial linear algebra and geometry. In (Pawlowsky-Glahn & Buccianti, 2011), pp. 141–157, 378 p (2011)

J.J. Egozcue, V. Pawlowsky-Glahn, Basic concepts and procedure. In (Pawlowsky-Glahn & Buccianti, 2011), pp. 12–28, 378 p (2011)

J.J. Egozcue, V. Pawlowsky-Glahn, G. Mateu-Figueras, C. Barceló-Vidal, Isometric logratio transformations for compositional data analysis. Math. Geol. **35**, 279–300 (2003). https://doi.org/10.1023/A:3A1023818214614

C. Gozzi, R.S. Graziano, A. Buccianti, Statistical methods for the geochemical characterisation of surface waters: the case study of the Tiber River basin (Central Italy). Comput. Geosci. **131**, 80–88 (2019a)

C. Gozzi, R.S. Graziano, F. Frondini, A. Buccianti, Innovative monitoring tools for the complex spatial dynamics of river chemistry: case study for the Alpine region. Environ. Earth Sci. **77**(16), 579 (2019b)

R.S. Graziano, C. Gozzi, A. Buccianti, Is compositional data analysis a theory able to discover complex dynamics in aqueous geochemical systems?. J. Geochem. Explor. **211**, 106465, 1-9 (2020)

ISPRA Ambiente, Italian hydrogeological complexes [Download MAIS]. http://www.sinanet.isprambiente.it/it/sia-ispra/download-mais/complessi-idrogeologici/view (24 September, 2019) (2019)

A. Kleidon, Life, hierarchy, and the thermodynamic machinery of planet Earth. Phys. Life Rev. **7**, 424–460 (2010)

W. Krzanowski, *Principles of Multivariate Analysis (Second Edition)*. Oxford Statistical Science Series 23 (Oxford, UK, 2000), 563 p

M. Mitzenmacher, A brief history of generative models for power law and lognormal distributions. Internet Math. **1**(2), 226–251 (2004)

B. Nisi, A. Buccianti, O. Vaselli, G. Perini, F. Tassi, A. Minissale, G. Montegrossi, Hydrochemistry and strontium isotopes in the Arno River Basin (Tuscany, Italy): constraints on natural controls by statistical modeling. J. Hydrol. **1–4**, 166–183 (2008)

V. Pawlowsky-Glahn, A. Buccianti (eds.), *Compositional Data Analysis: Theory and Applications* (Wiley, 2011), 378 p

V. Pawlowsky-Glahn, J.J. Egozcue, R. Tolosana-Delgado, *Modeling and Analysis of Compositional Data*. Statistics in Practice (Wiley, Chichester, UK, 2015), 247 p

V. Pawlowsky-Glahn, J.J. Egozcue, Geometric approach to statistical analysis on the simplex. Stoch. Environ. Res. Risk Assess. **15**(5), 384–398 (2001)

X. Sanchez-Vila, M. Dentz, L.D. Donado, Transport-controlled reaction rates under local non-equilibrium conditions. Geophys. Res. Lett. **34**, 1–5 (2007)

M. Scheffer, S.R. Carpenter, T.M. Lenton, J. Bascompte, W. Brock, V. Dakos, J. van de Koppel, I.A. van de Leemput, S.A. Levin, E.H. van Nes, M. Pascual, J. Vandermeer, Anticipating critical transitions. Science **338**, 344–348 (2012)

A.J.E. Seely, P. Macklem, Fractal variability: an emergent property of complex dissipative systems. Chaos **22**, 013108-1–013108-7 (2012)

S.L. Shvartsev, Self-organizing abiogenic dissipative structures in the geologic history of the Earth. Earth Sci. Front. **16**(6), 257–275 (2009). https://doi.org/10.1007/BF02066299. https://doi.org/10.1016/S1872-5791(08)60114-1

M.M.J.W. van Rooij, B. Nash, S. Rajaraman, J.G. Holden, A fractal approach to dynamic inference and distribution analysis. Front. Physiol. **4**, 1–16 (2013)

C. Zhu, G. Anderson, *Environmental Applications of Geochemical Modeling* (Cambridge University Press, 2002), 284 p, 378 p