

1 Combining raw and compositional data to determine the spatial
2 patterns of Potentially Toxic Elements in soils

3

4 C. Boente^a, M.T.D. Albuquerque^b, A. Fernández-Braña^a, S. Gerassis^c, C. Sierra^d,
5 J.R. Gallego^a

6

7 ^a *INDUROT and Environmental Technology, Biotechnology, and Geochemistry Group,*
8 *Universidad de Oviedo, Campus de Mieres, 33600 Mieres (Asturias), Spain*

9 ^b *Instituto Politécnico de Castelo Branco, 6001-909 Castelo Branco, Portugal and CERENA/FEUP*
10 *Research Center, Portugal*

11 ^c *Department of Natural Resources and Environmental Engineering, Univ. of Vigo, Lagoas*
12 *Marcosende, 36310 Vigo, Spain*

13 ^d *Departamento de Transportes, Tecnología de Procesos y Proyectos, Universidad de Cantabria,*
14 *Campus de Torrelavega, 39300 Torrelavega (Cantabria), Spain*

15

16 **Abstract**

17 When considering complex scenarios involving a multiset of attributes, such as in
18 environmental characterization, a clearer picture of reality can be achieved through the
19 dimensional reduction of data.

20 In this context, maps facilitate the visualization of spatial patterns of contaminant
21 distribution and the identification of enriched areas. A set, of 15 Potentially Toxic
22 Elements (PTEs) – (As, Ba, Cd, Co, Cr, Cu, Hg, Mo, Ni, Pb, Sb, Se, Tl, V, and Zn), was
23 measured in soil, collected in the municipality of Langreo (80 Km²), in Asturias, northern
24 Spain, a paradigmatic industrial area.

25 With the aim to explore PTE dissemination trends and to define clusters of relative
26 enrichment, the mechanisms through which these contaminants are spatially distributed
27 were examined.

28 Relative enrichment (RE) is introduced here to refer to the proportion of elements present
29 in a given context. Indeed, a novel approach is provided for research into PTE fate. This

30 method involves studying the variability of PTE proportions throughout the study area,
31 thereby allowing the identification of dissemination trends.

32 Transformations to open closed data are widely used for this purpose. As compositions
33 are shown along with their spatial locations, spatial patterns have an indubitable interest.
34 In this study, the Centered Log-ratio transformation (*clr*) was used, followed by its back-
35 transformation, to build a set of compositional data that, combined with raw data, allowed
36 to establish the sources of the PTEs and trends of spatial dissemination.

37 Based on the obtained findings was possible to conclude that the Langreo area is deeply
38 affected by its industrial and mining legacy. The city centre is highly enriched in Pb and
39 Hg and As shows enrichment in a northwesterly direction. Overall, the multivariate
40 geochemical approach presented facilitates the identification and quantification of
41 anthropogenic impacts and consequent adequate monitoring measures required to
42 safeguard the health of local communities.

43

44 **Keywords:** Soil Pollution, PTEs, Compositional Data, Ordinary Kriging, Local G-
45 clustering, Relative Enrichment.

46

47 **1. Introduction**

48 Environmental characterization involves complex scenarios in which a multiset of
49 attributes must be considered. A dimensional reduction of data is pivotal to gain a clear
50 picture of reality (Moen and Ale, 1998). Maps are useful to visualize pollutant
51 concentrations, as well as to determine zones of contaminant enrichment, whether
52 natural or caused by anthropogenic activity. In this context, Potentially Toxic Elements
53 (PTEs) are increasingly affecting soils all over the world, thus posing a threat to both
54 public health and the environment (McIlwaine et al., 2016). The presence of these
55 elements in soils can be explained by many factors (Alloway, 1990), the growth of

56 urbanization and resulting increase in industrial activities being among the most
57 important (Biasioli et al., 2006). Given that high concentrations of PTEs can endanger
58 human and environmental health, it is of utmost importance to characterize their spatial
59 distribution, determine their source, and screen for enrichment trends (Fayiga and Saha,
60 2016; Li et al., 2014; Boente et al., 2017; Cachada et al., 2013).

61 The area of Langreo (Asturias, NW Spain) (Fig. 1) is one of the regions in the Iberian
62 Peninsula most marked by industrialization (Gallego et al., 2016). Coal mining and
63 industries devoted to energy, metallurgy, pharmacology, and fertilizers, among others,
64 have been operating in this region for decades, leaving a lasting imprint on the
65 environment (Martínez et al., 2014; Megido et al., 2017). In this regard, great amounts
66 of PTEs have been identified in soils from former industrial plots in this area (Boente et
67 al., 2016; Gallego et al., 2016).

68 A comparative study of a set of 15 chemical elements was performed, analyzed in soils
69 gathered in the Langreo area (80 Km²), a paradigmatic industrial area as described
70 above. In this sort of studies, the distribution of PTEs cannot be studied by merely
71 considering the total concentrations (raw data), especially when the concentration of
72 chemical elements in almost all datasets is compositional (Pawłowsky-Glahn., 1989;
73 Filzmoser et al 2009), where attributes vary together with all the others. In this context,
74 transformations that open closed data are widely used and, as compositions are
75 recorded along with their spatial locations, spatial patterns are of interest (Pawłowsky-
76 Glahn., 1989). The contributions of Pawłowsky-Glahn to regionalized compositions
77 (Pawłowsky-Glahn, 1989; Pawłowsky-Glahn and Burger, 1992; Pawłowsky-Glahn et al.,
78 1995) and their applications are widely applied (Odeh et al., 2003; Lark and Bishop,
79 2007). In this context, multiple log-ratio transformations are commonly used, the most
80 common being the additive log-ratio transformation (alr), the centered log-ratio
81 transformation (clr) (e.g. Aitchison, 1986), and the isometric log-ratio transformation (ilr)
82 (Egozcue et al., 2003). In this study, the clr transformation and its back-transformation

83 were performed through CoDaPack v2.02.21 software to create a set of compositional
84 data that provides information about the comparative magnitudes of their constituents.
85 This compositional dataset was used to map patterns of RE, thereby allowing to identify
86 spatial dissemination trends for PTEs.

87 In summary, the main goal of this study was to test a methodology that, by means of
88 combining raw and compositional data, has the capacity to identify spatial patterns, areas
89 of pollution risk and anthropogenic or natural sources of PTEs. All the evidence provided
90 is supported by uni- and multi-variate statistical analysis, together with ordinary kriging
91 and Local G clustering for the area of Langreo. Finally, core strengths and weaknesses
92 are extrapolated to make this methodology useful and applicable to studies of a similar
93 nature.

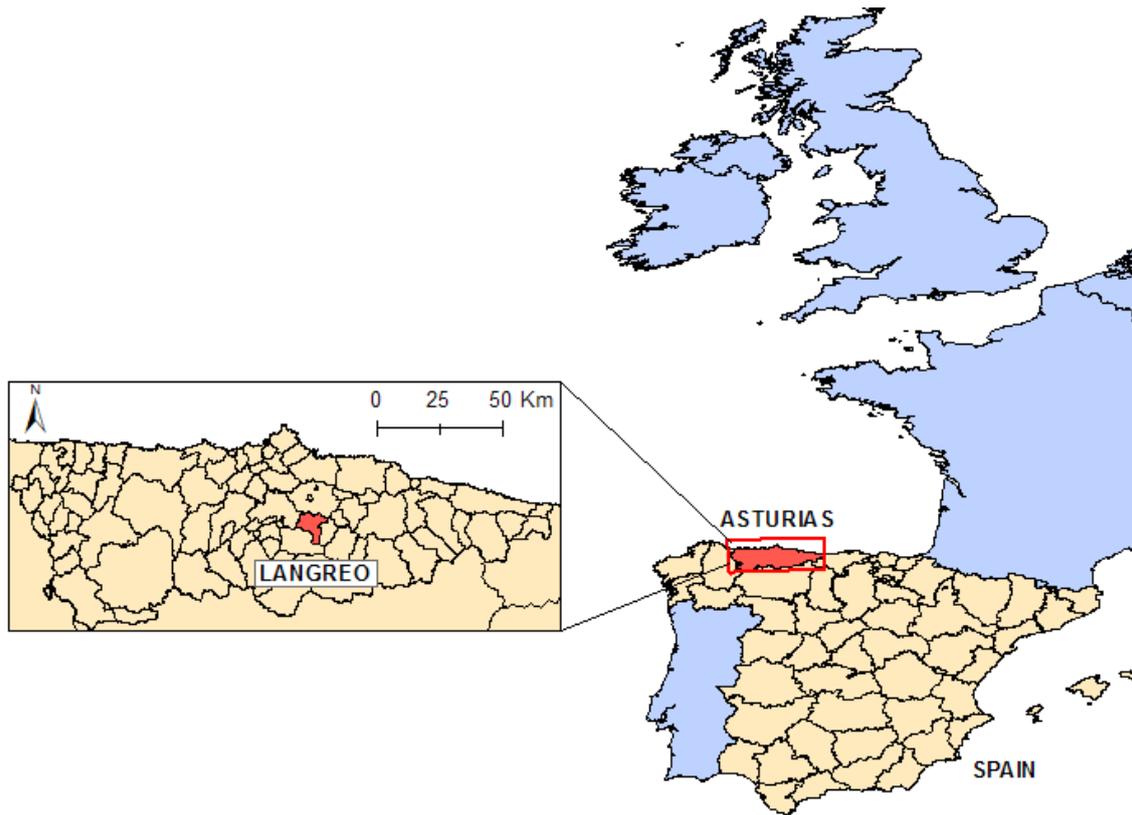
94

95 **2. Materials and Methods**

96 *2.1. Study area*

97 Covering 80 km², the municipality of Langreo (Asturias, NW Spain, Fig. 1) has a history
98 of mining and industrial activity that dates back to the 1850s (Martínez et al., 2014). This
99 activity left behind a legacy of polluted sites, making this zone one of the most
100 contaminated areas in northern Spain (Gallego et al., 2016) and thus an ideal site in
101 which to test the method presented in this study.

102 The region lies along the Nalón River, which is the longest and the most voluminous in
103 Asturias. Altitudes in the area vary from 200 m (location of the urban areas and industry)
104 to 900 m (rural environments, forests), with the presence of steep mountains. This
105 geography gives rise to an enclosed area that facilitates the accumulation of PTEs by
106 atmospheric deposition.



107

108 Fig. 1. Location of the study area in the municipality of Langreo in Asturias, Spain.

109

110 2.2. Data collection and chemical analyses

111 Samples were collected using a stratified systematic sampling method at random
112 distances to obtain a representative set of data on the total variability of PTE content and
113 site diversity (natural or anthropic environments, geomorphology, land uses, etc.). To
114 this end, 10 equidistant transects, 250 m wide and each one 1000 m apart, were
115 distributed perpendicular to the Nalón River (Fig. 2). A total of 150 samples were
116 collected, the number *per* transect being determined proportionally to its length. The
117 sample location within each transect was selected at random (Fig. 2).

118 Each sample composed of five increases taken from each vertex of a 1-m edge square
119 and its central point from the top 20-25 cm of the soil, using an Edelman Auger.
120 Afterwards, samples were passed through a 2-cm mesh screen *in situ* to remove large
121 material such as organic matter, rocks, and gravel. The samples were then dried in an

122 oven at 35°C to prevent the evaporation of volatile compounds, and finally quartered by
123 means of a Jones riffle splitter for soil homogenization and representativeness.

124 These fractions were ground in an RS100 Resch mill at 400 RPM for 40 s. Then, 1-g
125 representative sub-samples were sent to the ISO 9002-accredited Bureau Veritas
126 Laboratories (Vancouver, Canada) and subjected to 1:1:1 “aqua regia” digestion. The
127 total concentrations of the elements Ag, Al, As, Au, B, Ba, Bi, Ca, Cd, Co, Cr, Cu, Fe,
128 Ga, Hg, K, La, Mg, Mn, Mo, Na, Ni, P, Pb, S, Sb, Sc, Se, Sr, Te, Th, Ti, Tl, U, V, W and
129 Zn in the digested material were determined by Inductively Coupled Plasma-Optical
130 Emission Spectroscopy (ICP-OES).

131 A subset of the analyzed elements corresponding to PTEs was used for this study. This
132 subset was chosen because it represented a set of typical contaminants (heavy
133 metal(loid)s) found in environmental studies in Asturias (Albuquerque et al., 2017;
134 Boente et al., 2016; Gallego et al., 2015), in addition the Risk Based Soil Screening
135 Levels (RBSSLs) for these contaminants are available for this region of Spain (BOPA,
136 2014). Furthermore, the dispersal of the concentrations of these contaminants never
137 exceeded three orders of magnitude and thus provided readable proportions. Therefore,
138 of the original list of 36 elements, the following 15 were examined (PTE group): As, Ba,
139 Cd, Co, Cr, Cu, Hg, Mo, Ni, Pb, Sb, Se, Tl, V and Zn.

140

141

142

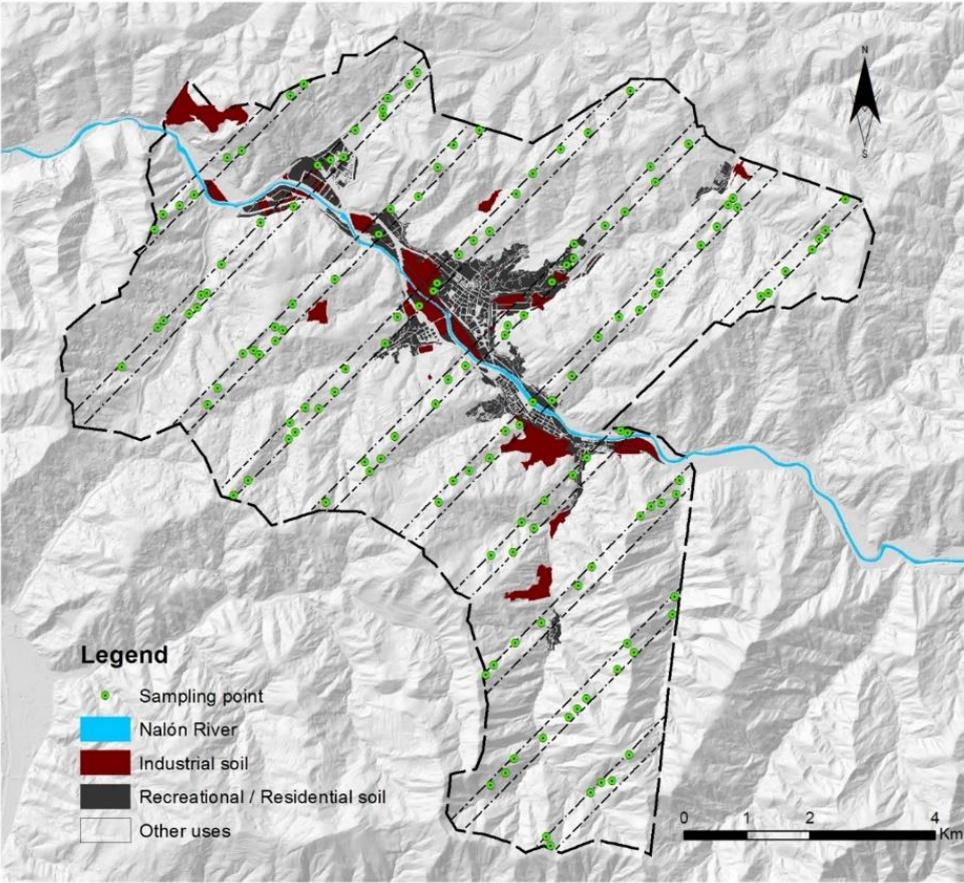
143

144

145

146

147
148
149
150
151
152
153
154
155
156
157
158



159 Fig. 2. Sampling design and land use categories in the study area.

160

161 *2.3. Data transformation – compositional data and the closure problem*

162 In geochemistry, compositional data is obtained by transforming each original raw
163 concentration (i.e. mg/kg of an element in a sample) into proportions of a whole whose
164 elements sum one or 100% (Pawlowsky-Glahn and Egozcue, 2006). However, the
165 unfeasibility of analyzing all the elements in a given soil hinders the consideration of
166 proportions. Indeed, this issue has been heavily debated and is referred to by
167 researchers as the closure problem (Filzmoser et al., 2009b). In environmental science
168 studies, it is generally accepted that the elements analyzed make up the entirety of the
169 soil on the condition that a suitable number of such elements is included in the study
170 (Campbell et al., 2009; Reimann et al., 2012). Moreover, other authors work with

171 subcompositions, defined as a subset of components of parts of a composition
172 (Pawlowsky-Glahn and Buccianti, 2011). Subcompositions are feasible when they
173 respect the principles of compositional data (Greenacre and Lewi, 2009), including the
174 subcompositional coherence principle (Aitchison, 1986).

175 The most frequently used log-ratio transform functions (*alr*, *clr* and *ilr*) have both
176 advantages and disadvantages, which are widely discussed in the literature. The *clr*
177 transformation is the prevailing function in geochemical studies as it uses the geometric
178 mean as normalizer parameter and it was chosen for the purposes of the present study
179 **(introduzir citações).**

180 The centred log-ratio transformation (*clr*) equation was adapted from (Aitchison, 1986):

$$clr(x) = \ln \left(\frac{C_j}{\sqrt[D]{\prod_{j=1}^D C_j}} \right) \quad (1)$$

181

182 where C_j is the concentration of pollutant j and D is the number of parts into which the
183 composition is divided (in this case, the number of pollutants considered).

184 The back-transformation equation is computed as:

$$\overline{clr}(x) = \frac{e^{clr(x)}}{\sum_{j=1}^D e^{clr(x)}} \quad (2)$$

185

186 **Where.....This** equation allows representation of the *clr*-transformed data as
187 compositional data (proportions). This means that the sum of all the elements after back-
188 transformation is equal to 1. The *clr* transformation and the calculation of its back-
189 transformation was performed using CoDaPack v2.02.21 software
190 (<http://www.compositionaldata.com/codapack.php>).

191

192

193 2.4. Spatial modeling

194 The spatial characterization of PTE distribution was performed with the following two
195 complementary objectives in mind. First, to define spatial clusters of PTE concentration.
196 To accomplish this, the raw dataset was used, allowing to interpret contamination
197 outbreaks and therefore locate the main sources of PTEs. Second, to define Relative
198 Enrichment (RE) spots along the study area and thus, evaluate trends of enrichment.
199 Indeed, rather than simply looking at PTE content enrichment, it was possible to
200 approach the study of PTE's fate, by examining the changes in their proportions
201 throughout the study area and therefore. The compositional dataset was used to tackle
202 this issue, and spatial clusters of RE were computed.

203 A four-step methodology was adopted as follows:

- 204 • Principal Components Analysis (PCA) for reducing dimensionality and for
205 evaluating variable association was performed. PCA is one of the most important
206 multivariate statistical methods and it is widely used for data preprocessing and
207 dimension reduction (raw and compositional data). The aim of PCA is to reduce
208 the dimensionality of data while simultaneously preserving the within variability
209 structure (variance–covariance) (e.g. Zuo et al., 2016). The analysis starts with p
210 random attributes X_1, X_2, \dots, X_p , where no assumption of multivariate normality is
211 required. The axes of the constant ellipsoids correspond to the new synthesis
212 variables, the principal components. The XlStat 2013.1.01 software
213 (<https://www.xlstat.com/en/>) was used for computational purposes.
- 214 • Selected attributes were subjected to a structural analysis, and experimental
215 variograms were computed for both raw and compositional data. The variogram
216 is a vector function used to calculate the spatial variation structure of regionalized
217 variables (Matheron, 1971; Journel and Huijbregts, 1978; Gringarten and
218 Deutsch, 2001).

- 219 • Spatial prediction through Ordinary Kriging (OK), aiming to predict the values for
220 the variables at any arbitrary spatial location within the study region, was
221 performed. The raw dataset was used to infer the concentration and PTE origin,
222 as the compositional dataset was used for dissemination trend detection and
223 local RE evaluation. Of note, geostatistics are a reference approach for the
224 characterization of environmental hazards in contexts in which the information
225 available is scarce (citações). The primary application of geostatistics is to
226 estimate and map environmental attributes in unsampled areas where Kriging is
227 a generic name for a set of generalized least-squares regression algorithms.
228 Ordinary Kriging (OK) accounts for local fluctuations of the mean by limiting the
229 field of stationarity of the mean to the local neighborhood (Goovaerts 1997). For
230 the computation, the Space-Stat Software V. 4.0.18, Biomedware was used
231 (Albuquerque et al., 2014) (Fig. 6).
- 232 • Finally, Local G clustering was performed. This technique allows measurement
233 of the degree of association that results from the concentration of weighted points
234 (or region represented by a weighted point) and all other weighted points included
235 within a radius of distance from the original and defining clusters of high (high-
236 ring) and low (low-ring) significance. For computation, the SpaceStat V. 4.0-18.
237 software (<https://www.biomedware.com/>) was used.

238

239 3. Results and discussion

240 3.1. Descriptive statistics

241 Descriptive statistics for raw and *clr*-transformed data were computed (Table 1). The raw
242 data revealed considerable variability for some elements, which was of concern for As,
243 Cd, Cu, Pb, Sb and Zn, whose maximum values surpassed the RBSSLs (BOPA, 2014).
244 The 5% trimmed mean allowed to conclude that extreme values were concentrated
245 mainly in the upper 2.5% intervals, as the remaining 97.5% can be approximated by the

246 normal distribution. Once the *clr*-transformed data were applied, the associated standard
 247 deviation was clearly reduced and the mean, median and 5% trimmed mean tended to
 248 be similar. Indeed, the *clr* data showed a normal distribution as a result of diminishing
 249 the weight of outliers. This diminished weight enhanced the prediction of data proportions
 250 after the back-transformation of *clr* data, and compositional data were obtained.

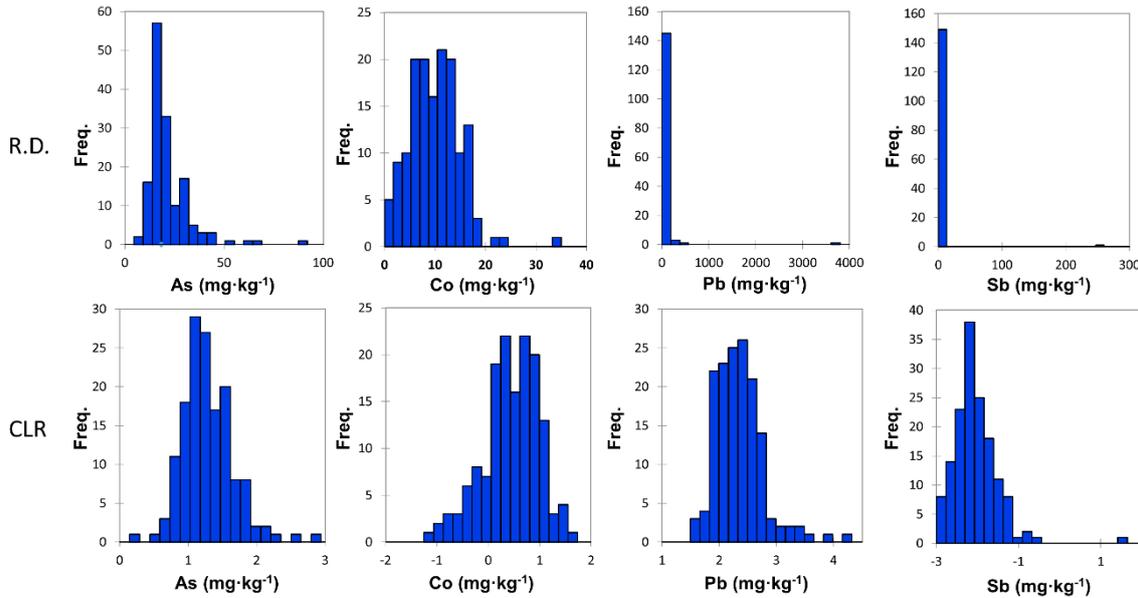
251
 252 Table 1. Descriptive statistics for 15 PTEs: Range, Mean, Median, Standard Deviation (SD), and
 253 Trimmed Mean (T.Mean 5%) are expressed in mg·kg⁻¹, Relative Standard Deviation (RSD) is
 254 expressed in %.

PTE	Raw Data						Clr-Transformed Data				
	Range	Mean	Median	SD	RSD	T.Mean 5%	Mean	Median	SD	RSD	T.Mean 5%
As	6.4 - 91.1	21.8	18.5	10.9	49.8	21.0	21.9	20.9	6.3	28.9	21.7
Ba	11.0 - 1747.1	107.9	66.9	168.7	156.3	90.2	79.2	74.6	16.7	21.1	78.3
Cd	0.02 - 26.9	0.6	0.3	2.2	382.6	0.4	0.4	0.3	0.1	19.2	0.3
Co	1.1 - 34.0	10.0	9.8	5.0	49.8	9.9	9.4	10.2	11.4	121.7	9.4
Cr (III)	5.7 - 69.0	18.9	18.6	6.7	35.6	18.5	19.6	20.1	4.8	24.5	19.6
Cu	3.0 - 2022.2	39.0	22.7	163.6	419.2	24.6	24.4	24.2	7.3	29.7	24.1
Hg	0.1 - 2.6	0.4	0.3	0.4	95.5	0.4	0.3	0.3	0.1	21.3	0.3
Mo	0.4 - 4.6	1.0	0.9	0.6	53.6	1.0	1.0	1.0	0.2	16.0	1.0
Ni	1.4 - 52.8	18.3	16.5	9.1	49.7	18.0	17.5	17.5	7.2	41.1	17.5
Pb	10.5 - 3729.5	91.6	52.2	302.7	330.6	64.0	62.8	60.7	11.1	17.7	61.8
Sb	0.3 - 256.6	2.5	0.6	20.8	821.8	0.8	0.8	0.7	0.2	26.6	0.8
Se	0.1 - 1.9	0.9	0.8	0.4	45.1	0.8	0.8	0.9	0.3	30.3	0.8
Tl	0.0 - 0.5	0.2	0.2	0.1	33.8	0.2	0.2	0.2	0.0	10.6	0.2
V	7.0-56.0	27.9	27.0	6.9	24.8	27.8	29.6	29.8	6.3	21.2	29.8
Zn	16.9-2161.0	136.2	107.2	179.4	131.7	120.8	119.8	120.8	11.8	9.9	120.1

255
 256 On the basis of comparison of the histograms (Fig. 3) of the raw and compositional
 257 datasets, it is possible to reason that: a) when considering the raw dataset, asymmetric
 258 distributions are found for almost all the PTEs, and these distributions are biased mainly
 259 by the presence of outliers; b) the *clr*-transformed dataset shows an important feature as
 260 it allows the assumption of normality. Therefore, it was possible to conclude that the *clr*-
 261 transformed dataset and the compositional dataset (after *clr* back-transform) have two
 262 principal advantages, namely, they allow work with proportions and at the same time,
 263 improves data normalization.

264 Of note were the anomalous As, Cd, Cu, Pb, Sb and Zn concentrations, which greatly
 265 exceeded the RBSSLs (BOPA, 2014) (Table 1). These elements are classic fingerprints
 266 of heavy industrial activity. However, the presence of Ba, Co, Cr, Hg, Mo, Ni, Se, Tl and
 267 V did not constitute an immediate risk for human health or the environment.

268



269

270 Fig. 3. As, Co, Pb and Sb histograms for raw data (R.D.) and *clr*-transformed data (CLR).

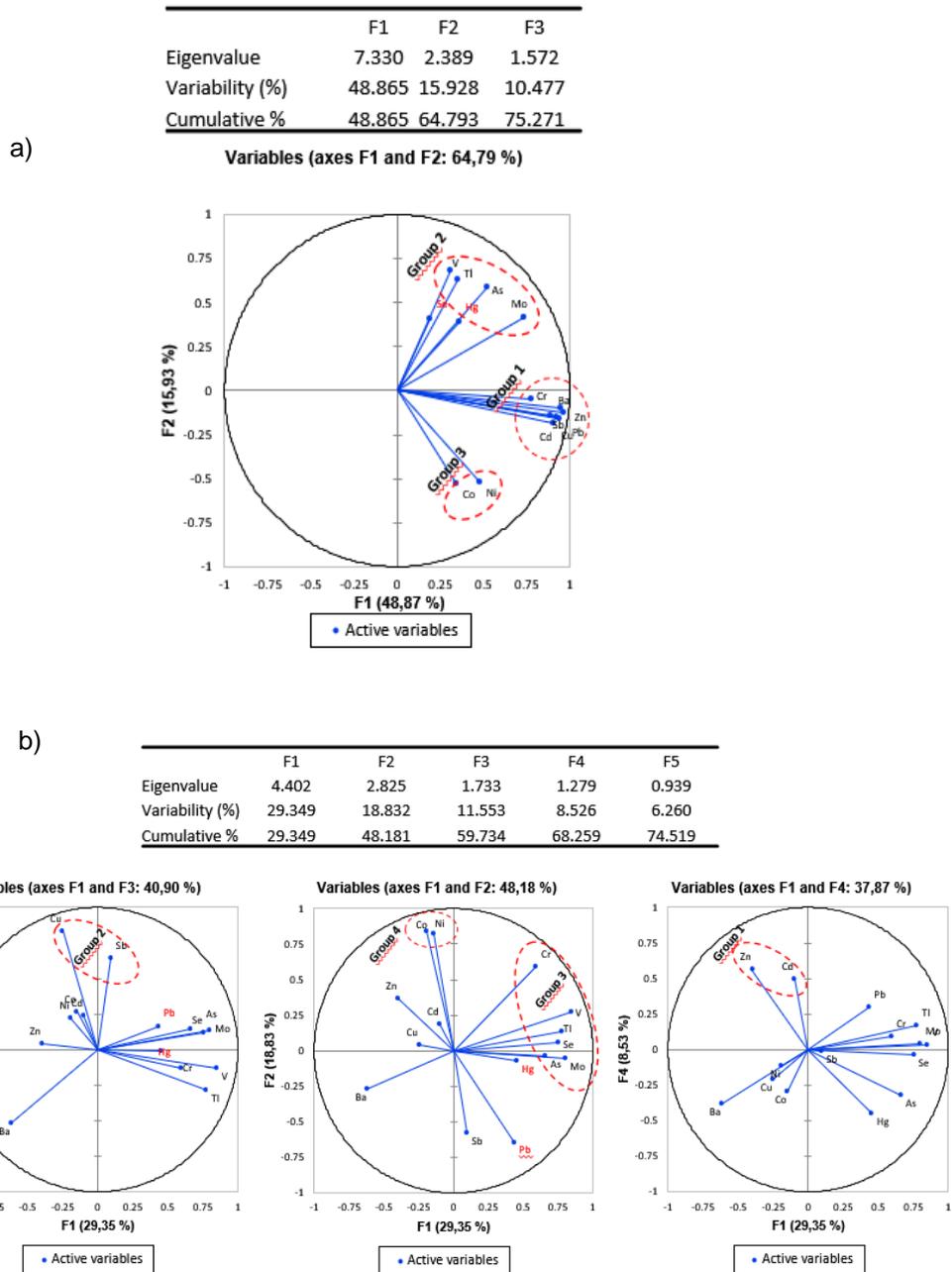
271

272 3.2. Multivariate statistics – Principal Components Analysis

273 When running the raw dataset, PCA results revealed three groups (Fig. 4 a)): a) the first
 274 formed by Ba, Cd, Cr, Cu, Pb, Sb and Zn—a typical association of heavy metals; b) the
 275 second composed by As, Mo, Tl and V; and c) the third representing Co and Ni. Finally,
 276 Hg and Se showed independent behaviors, thereby possibly indicating different sources.
 277 On the other hand, when considering the compositional dataset, slight differences in the
 278 results were observed (Fig. 4 b). The first-mentioned group (Ba, Cd, Cr, Cu, Pb, Sb and
 279 Zn) was split in two: a) the first comprising Cd and Zn; b) the second Cu and Sb.
 280 Furthermore, two more groups were identified, c) the third comprising As, V, Tl, Mo, Se
 281 and Cr; and d) the fourth Ni and Co. Mercury (Hg) and Pb were found to be independent.

282 The PCA's results lead to conclude that the compositional dataset provides a fuller
 283 recognition of relevant contaminant associations. When setting a dependence on weight
 284 between elements, those which increase or decrease proportionally tend to be
 285 associated.

286



287

288

289 Fig. 4. a) PCA - Raw dataset; b) PCA - Compositional data.

290 3.3. Spatial modeling – geostatistical approach

291 At this point, As, Cu, Hg, Pb, and Zn were chosen for spatial modeling purposes as they
292 are core PTEs in contamination forecasts and representative of the most important
293 groups identified (Fig. 4).

294 The spatial stochastic patterns of the **five chosen** PTEs were constructed following a
295 three-step geostatistical modeling method.

296

297 3.3.1 Structural analysis and experimental variograms

298 The selected variables were subjected to a structural analysis, and experimental
299 variograms were computed. The variogram is a vector function used to calculate the
300 spatial variability of regionalized variables defined by the following equation (Matheron,
301 1971; Journel and Huijbregts, 1978):

$$\gamma(h) = \frac{1}{2N(h)} \sum_{2N(h)}^{N(h)} [Z(x_i) - Z(x_i + h)]^2 \quad (3)$$

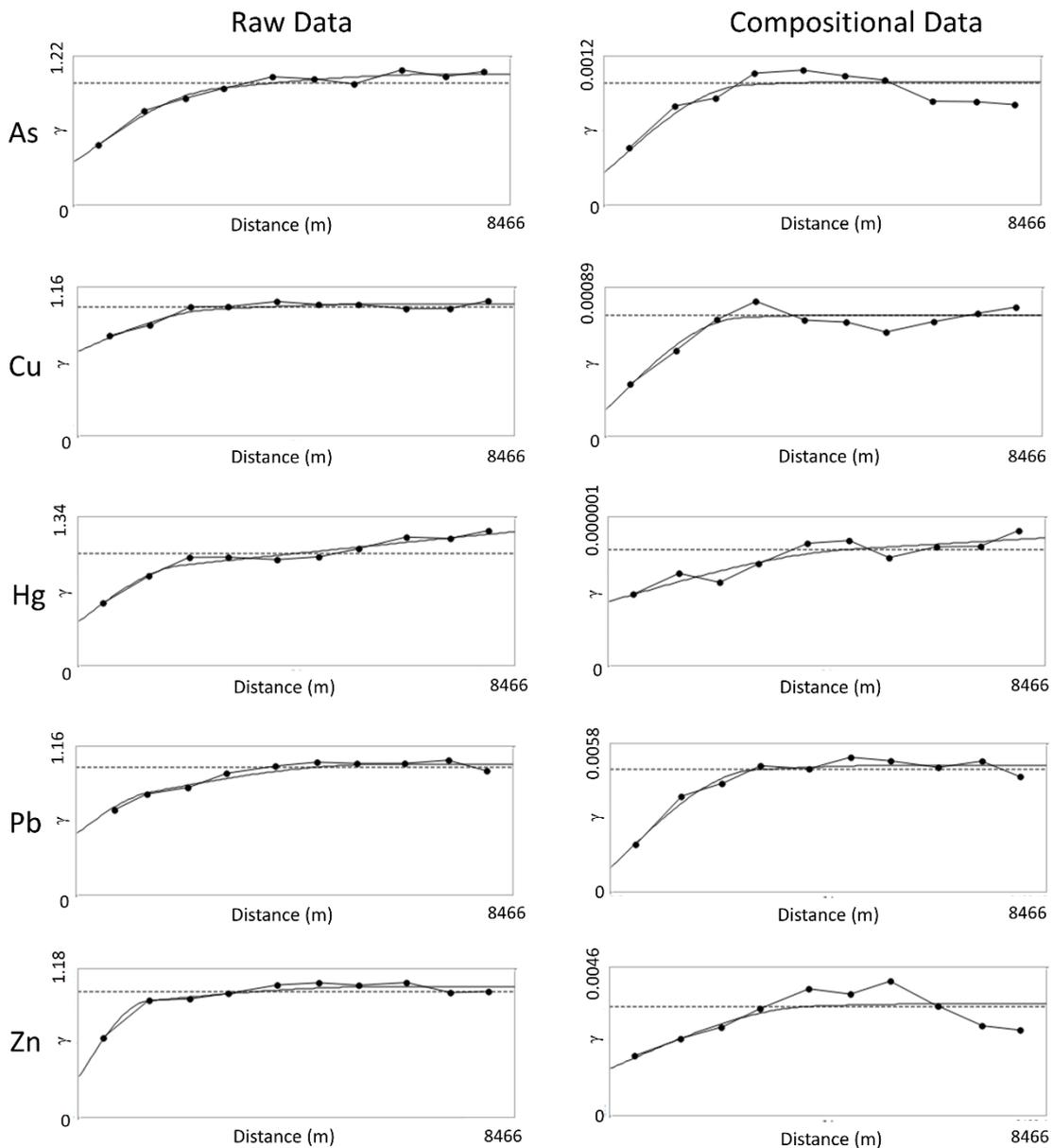
302

303 Its argument is h (distance), where $Z(x_i)$ and $Z(x_i+h)$ are the numerical values of the
304 observed variable at points x_i and x_i+h . The number of pairs forming for a h distance is
305 $N(h)$. Thus, it is the median value of the square of the differences between all pairs of
306 points in the geometric field spaced at a h distance. The graphic study of the variograms
307 obtained provides an overview of the spatial structure of the variable. One of the
308 parameters that provide such information is the nugget effect (C_0), which shows the
309 behavior at the origin. The other two parameters are the sill (C_1) and the amplitude (a)
310 which define the inertia used in the interpolation process and the influence radius of the
311 variable, respectively (Table 2).

312 The experimental variograms $\gamma(h)$ were then fitted to a theoretical model, $\hat{\gamma}(h)$ (Isaaks and
313 Srivastava 1989). The adjusted parameters for the five PTEs of the theoretical

314 variograms (raw and compositional datasets) (Fig. 5) allowed to observe that the
 315 isotropic variograms obtained generally showed a better fit for the compositional dataset.
 316 Indeed, the attributes showed a nugget effect below 40% of the total variance of all the
 317 attributes (Table 2). The error associated with the interpolation **procedure (OK)**, is
 318 therefore minimized when using the compositional dataset.

319



320

321 Fig. 5. Isotropic experimental variograms and fitted models for the raw and compositional
 322 datasets.

323

324 Table 2. Experimental variogram parameters for the raw and compositional datasets: a (m) is the
 325 amplitude; C_0 represents the value of the nugget effect; C_1 and C_2 , the value of the sill of the first
 326 and the second spherical structure respectively, and $C_0(\%Var)$ and $C_1+C_2 (\%Var)$ the mutual
 327 variances weighing for nugget and sill respectively.

	Parameters	As	Cu	Hg	Pb	Zn
Raw Data	A	2738	2575	1997	1376	1327
	C_0	0.356	0.664	0.401	0.488	0.330
	C_1	0.465	0.256	0.411	0.201	0.544
	C_2	0.260	0.110	1.17	0.339	0.172
	$C_0(\%Var)$	33	64	20	47	32
	$C_1+C_2 (\%Var)$	67	36	80	53	68
Comp. Data	A	2700	2569	4758	2808	3903
	C_0	$2.77 \cdot 10^{-4}$	$1.63 \cdot 10^{-4}$	$5.93 \cdot 10^{-7}$	$9.90 \cdot 10^{-4}$	$1.47 \cdot 10^{-3}$
	C_1	$6.45 \cdot 10^{-4}$	$4.83 \cdot 10^{-4}$	$3.50 \cdot 10^{-7}$	$3.52 \cdot 10^{-3}$	$1.58 \cdot 10^{-3}$
	C_2	$1.14 \cdot 10^{-4}$	$8.11 \cdot 10^{-5}$	$6.90 \cdot 10^{-7}$	$5.35 \cdot 10^{-4}$	$4.18 \cdot 10^{-4}$
	$C_0(\%Var)$	27	22	36	20	42
	$C_1+C_2 (\%Var)$	73	78	64	80	58

328

329 3.3.2 Spatial prediction: Ordinary kriging

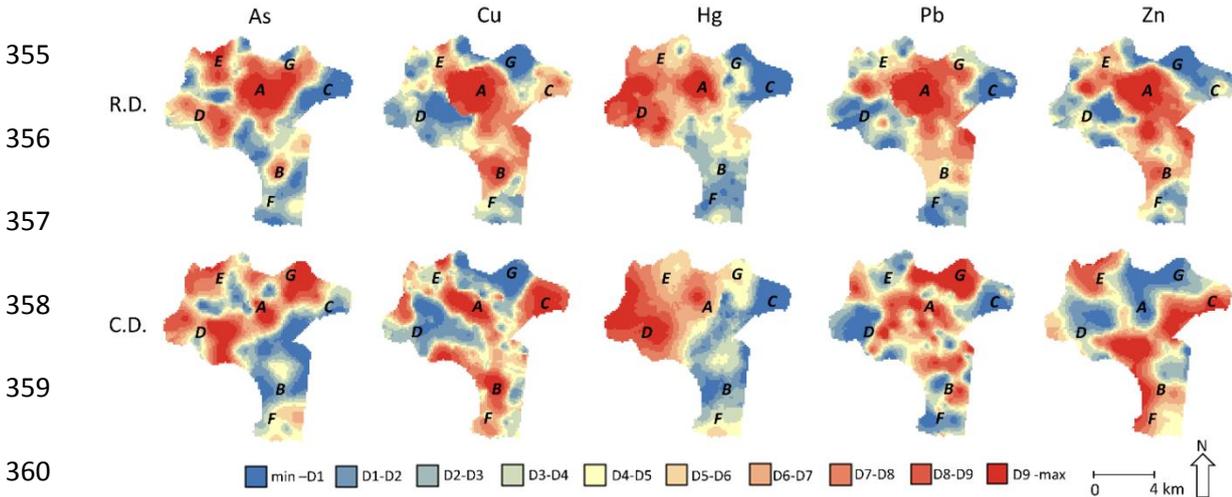
330 Analysis of the outputs obtained (Fig. 6) revealed evident contrasts between the raw
 331 and the compositional dataset representations. Care must be taken when interpreting
 332 representations as they reflect distinct data. In this regard, the raw dataset mapping
 333 shows the estimated picture of PTE concentration distribution, thus indicating possible
 334 sources of these contaminants. In contrast, the compositional dataset mapping shows
 335 the spatial variability of PTE proportion, thus reflecting PTE's Relative Enrichment, and
 336 providing crucial information about the fate of these compounds within the study area.
 337 To facilitate the understanding of the results, the study area was divided into various
 338 zones of interest (Fig. 6) and interpreted as follows:

339 a) Considering the maps of the raw dataset (Fig. 6 -R.D.), high concentrations for all
 340 PTEs (Zn, Hg, As, Pb and Cu) in the central zone (zone A) can be observed, which
 341 coincides with the city of Langreo (Fig. 6). Moreover, Cu and Zn showed notable
 342 presence in the southern area (zone B), where the mining industry (coal mines and
 343 processing) were located (Fig. 1). The Cu map shows a north-eastern red-colored site
 344 (zone C) coinciding with a former coal-mining area. On the other hand, high

345 concentrations of Hg and As were observed in the western (zone D) and northern (zone
346 E) areas, which may be explained by the proximity to a derelict Hg mine (El Terronal site)
347 whose impact has been widely discussed (e.g. Gallego et al., 2015, González-Fernández
348 et al., 2018);

349 b) Concerning the compositional dataset (Fig. 6-C.D.), **Relative Enrichment** in Cu, Pb
350 and Zn was identified towards south (zone F) and northeast (zone C) of the area (Fig.
351 6), where the corresponding distribution was at its lowest level when using the raw data.
352 Cu, Pb and Zn showed a significant distribution throughout the area and therefore a
353 marked RE.

354



361

362 Fig. 6. Ordinary **kriging (OK) results**. Raw data (R.D) and compositional data (CD) respectively.
363 Scale is expressed in deciles (D_i) of mg·kg⁻¹ (R.D). and of % (C.D).

364

365

366

367

368

369

370 3.3.3 Spatial prediction: Local G clustering

371 To reinforce the findings of the previous section, a Local G clustering was conducted to
372 assess the level of association resulting from the concentration of weighted points (or
373 region represented by a weighted point) and all other weighted points included within a
374 radius from the original point. In this regard, a given zone was subdivided into n regions,
375 $l = 1, 2, \dots, n$, where each neighborhood is distinguished with a point whose Cartesian
376 coordinates are known. Each i has a value x (a weight) taken from a variable X
377 associated with it. The variable holds a natural origin and it is positive. The $G(i)$ statistic
378 developed below allows the testing of hypotheses concerning the spatial concentration
379 of the sum of x values associated with the j points within d of the i^{th} point. The following
380 statistic is obtained:

$$G_i(d) = \frac{\sum_{j=1}^n W_{ij}(d)X_j}{\sum_j^n x_j} \quad (4)$$

381

382 where W_{ij} is a symmetric one/zero spatial weight matrix with a value of 1 for all links
383 defined as being within distance d of a given i ; all other links are zero, including the link
384 of point i to itself. The numerator is the sum of all x_j within d of i but not including x_i . The
385 denominator is the sum of all x_j , excluding x_i (Getis and Ord, 1992).

386 The maps obtained (Fig. 7) provide a faster and more intuitive way to verify whether the
387 problematic zones detected previously are indeed of concern. Thus, red areas (high ring)
388 show the sites with the greatest accumulation of the PTEs, while the blue areas (low
389 ring) represent zones with low accumulation (Fig. 7). The highest accumulation of PTEs,
390 when considering the raw data clusters, was in the city center (high ring-zone A). The
391 soils in this area were clearly affected by PTE deposition, presumably due to heavy
392 industry and/or the transport of pollutants. However, examination of the significance of
393 the spatial clusters obtained using the compositional data shows several differences.
394 The central high ring (high significance) is now smaller, showing that the areas with the

395 highest concentration of these PTEs (Zn, Hg, As, Pb and Cu) do not totally overlap with
 396 the corresponding higher proportions and indicating that PTE transport and RE occurs
 397 in a westerly and southerly direction.

398

399

400

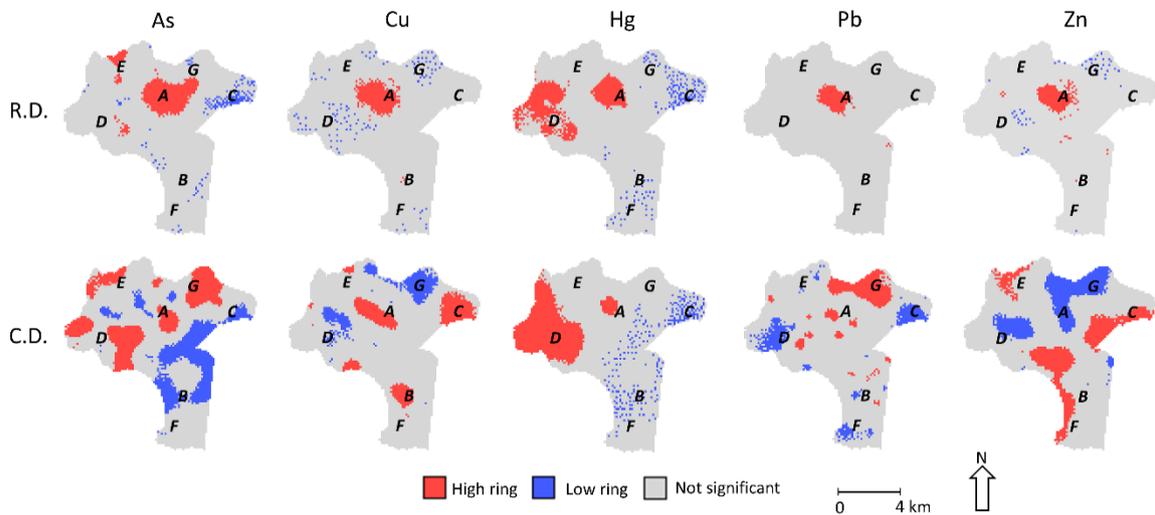
401

402

403

404

405



406 Fig. 7. Local G clusters. Raw data (R.D.) and compositional data (C.D.) respectively.

407

408 4. Conclusions

409 The degree of PTE contamination in the soil of an industrial area can be characterized
 410 using two datasets, namely raw and compositional (*clr*-transformed followed by the back-
 411 transformation function). To exemplify the complementary attributes of these two types
 412 of dataset, 150 soil samples were collected, and 36 elements were analyzed in Langreo
 413 (80 km²), a paradigmatic example of an industrial area affected by heavy metal and
 414 metalloids contamination. Univariate statistics allowed recognition of redundant
 415 information and the identification of outliers. The space of analysis was then reduced for
 416 both datasets by building the synthesis variables held by PCA. Five PTEs, namely Zn,
 417 Hg, As, Pb and Cu, were retained for spatial modeling due to their significance in the
 418 contamination forecast. Ordinary Kriging (OK) and Local G clustering allowed the
 419 construction of hazard maps, which facilitate the evaluation of the probable origin of
 420 PTEs (raw data) and their possible Relative Enrichment (compositional data).

421 Regarding the Langreo area, it is extensively affected by its industrial and mining history.
422 The following observations support this conclusion: 1. The city centre is highly enriched
423 in PTEs, which can be explained by heavy industry and pollutant transport, Pb being the
424 main contaminant; 2. The spatial distribution of Cu indicates a strong association with
425 coal mining and processing; and 3. Hg and As show enrichment in a northwesterly
426 direction, which is linked to natural mineralization and former Hg mining and metallurgy.
427 Future work would require an exhaustive study of covariates to shed light on PTE
428 dynamics and to clarify the main sources of PTEs, as well as their RE throughout the
429 study area.

430 The information gathered provides a basis for delimiting the polluted zones and the
431 sources of pollutants, thus facilitating the development of specific air and soil monitoring
432 activities, urban planning, and environmental policies.

433

434 **Acknowledgements**

435 C. Boente obtained a grant from the “Formación del Profesorado Universitario” program,
436 financed by the “Ministerio de Educación, Cultura y Deporte de España”.

437 M.T.D Albuquerque acknowledges a scholarship 567 SFRH/BSAB/127907/2016 from
438 the Foundation for Science and Technology (Portugal).

439

440 **References**

441 Aitchison, J., 1986. The Statistical Analysis of Compositional Data. *J. R. Stat. Soc.* 44,
442 139–177. doi:10.1007/978-94-009-4109-0.

443 Albuquerque, M.T.D., Gerassis, S., Sierra, C., Taboada, J., Martín, J.E., Antunes,
444 I.M.H.R., Gallego, J.R., 2017. Developing a new Bayesian Risk Index for risk
445 evaluation of soil contamination. *Sci. Total Environ.* 603–604, 167–177.

446 doi:10.1016/j.scitotenv.2017.06.068.

447 Alloway, B.J., 1990. The origins of heavy metals in soils, in: *Heavy Metals in Soils*. p.
448 339.

449 Antunes, I.M.H.R., Albuquerque, M.T.D., 2013. Using indicator kriging for the evaluation
450 of arsenic potential contamination in an abandoned mining area (Portugal). *Sci.*
451 *Total Environ.* 442, 545–552. doi:10.1016/j.scitotenv.2012.10.010.

452 Biasioli, M., Barberis, R., Ajmone-Marsan, F., 2006. The influence of a large city on some
453 soil properties and metals content. *Sci. Total Environ.* 356, 154–164.
454 doi:10.1016/j.scitotenv.2005.04.033

455 Bishop, T.F.A., McBratney, A.B., 2001. A comparison of prediction methods for the
456 creation of field-extent soil property maps. *Geoderma* 103, 149–160.
457 doi:10.1016/S0016-7061(01)00074-X.

458 Boente, C., Matanzas, N., García-González, N., Rodríguez-Valdés, E., Gallego, J.R.,
459 2017. Trace elements of concern affecting urban agriculture in industrialized areas:
460 A multivariate approach. *Chemosphere* 183, 546–556.
461 doi:10.1016/j.chemosphere.2017.05.129.

462 Boente, C., Sierra, C., Rodríguez-Valdés, E., Menéndez-Aguado, J.M., Gallego, J.R.,
463 2016. Soil washing optimization by means of attributive analysis: Case study for the
464 removal of potentially toxic elements from soil contaminated with pyrite ash. *J.*
465 *Clean. Prod.* doi:10.1016/j.jclepro.2016.11.007.

466 BOPA, Boletín Oficial del Principado de Asturias, 91, April 21, 2014. Generic Reference
467 Levels for Heavy Metals in Soils from the Principality of Asturias, Spain,
468 <https://sede.asturias.es/porta/site/Asturias/menuitem.1003733838db7342ebc4e191100000f7/?vgnextoid=d7d79d16b61ee010VgnVCM1000000100007fRCRD&fecha=21/04/2014&refArticulo=2014-06617&i18n.http.lang=es> (Accessed December
470 2017).
471

472 Bucciatti, A., Grunsky, E., 2014. Compositional data analysis in geochemistry: Are we
473 sure to see what really occurs during natural processes? *J. Geochemical Explor.*
474 141, 1–5. doi:10.1016/j.gexplo.2014.03.022.

475 Cachada, A., Dias, A.C., Pato, P., Mieiro, C., Rocha-Santos, T., Pereira, M.E., Da Silva,
476 E.F., Duarte, A.C., 2013. Major inputs and mobility of potentially toxic elements
477 contamination in urban areas. *Environ. Monit. Assess.* 185, 279–294.
478 doi:10.1007/s10661-012-2553-9.

479 Campbell, G.P., Curran, J.M., Miskelly, G.M., Coulson, S., Yaxley, G.M., Grunsky, E.C.,
480 Cox, S.C., 2009. Compositional data analysis for elemental data in forensic science.
481 *Forensic Sci. Int.* 188, 81–90. doi:10.1016/j.forsciint.2009.03.018.

482 Candeias, C., da Silva, E.F., Salgueiro, A.R., Pereira, H.G., Reis, A.P., Patinha, C.,
483 Matos, J.X., Ávila, P.H., 2011. The use of multivariate statistical analysis of
484 geochemical data for assessing the spatial distribution of soil contamination by
485 potentially toxic elements in the Aljustrel mining area (Iberian Pyrite Belt, Portugal).
486 *Environ. Earth Sci.* 62, 1461–1479. doi:10.1007/s12665-010-0631-2.

487 Dalla Libera, N., Fabbri, P., Mason, L., Piccinini, L., Pola, M., 2017. Geostatistics as a
488 tool to improve the natural background level definition: An application in
489 groundwater. *Sci. Total Environ.* 598, 330–340.
490 doi:10.1016/j.scitotenv.2017.04.018.

491 Egozcue, J.J., Pawlowsky-Glahn, V., Mateu-Figueras, G., Barceló-Vidal, C., 2003.
492 Isometric Logratio Transformations for Compositional Data Analysis. *Math. Geol.*
493 35, 279–300. doi:10.1023/A:1023818214614.

494 Fayiga, A.O., Saha, U.K., 2016. Soil pollution at outdoor shooting ranges: Health effects,
495 bioavailability and best management practices. *Environ. Pollut.* 216, 135–145.
496 doi:10.1016/j.envpol.2016.05.062.

497 Filzmoser, P., Hron, K., Reimann, C., 2009a. Principal component analysis for

498 compositional data with outliers, in: *Environmetrics*. pp. 621–632.
499 doi:10.1002/env.966.

500 Filzmoser, P., Hron, K., Reimann, C., 2009b. Univariate statistical analysis of
501 environmental (compositional) data: Problems and possibilities. *Sci. Total Environ.*
502 407, 6100–6108. doi:10.1016/j.scitotenv.2009.08.008.

503 Gallego, J.R., Esquinas, N., Rodríguez-Valdés, E., Menéndez-Aguado, J.M., Sierra, C.,
504 2015. Comprehensive waste characterization and organic pollution co-occurrence
505 in a Hg and As mining and metallurgy brownfield. *J. Hazard. Mater.* 300, 561–571.
506 doi:10.1016/j.jhazmat.2015.07.029.

507 Gallego, J.R., Rodríguez-Valdés, E., Esquinas, N., Fernández-Braña, A., Afif, E., 2016.
508 Insights into a 20-ha multi-contaminated brownfield megasite: An environmental
509 forensics approach. *Sci. Total Environ.* 563–564, 683–692.
510 doi:10.1016/j.scitotenv.2015.09.153.

511 Getis, A., Ord, J.K., 1992. The Analysis of Spatial Association by Use of Distance
512 Statistics. *Geogr. Anal.* 24, 189–206. doi:10.1111/j.1538-4632.1992.tb00261.x.

513 González-Fernández, B., Rodríguez-Valdés, E., Boente, C., Menéndez-Casares, E.,
514 Fernández-Braña, A., Gallego, J.R., 2018. Long-term ongoing impact of arsenic
515 contamination on the environmental compartments of a former mining-metallurgy
516 area. *Sci. Total Environ.* 610–611, 820–830. doi:10.1016/j.scitotenv.2017.08.135.

517 Goovaerts, P., 1997. *Geostatistics for Natural Resources Evaluation*. University Press,
518 New York: Oxford.

519 Goovaerts, P., 1999. *Geostatistics in soil science: State-of-the-art and perspectives*.
520 *Geoderma* 89, 1–45. doi:10.1016/S0016-7061(98)00078-0.

521 Greenacre, M., Lewi, P., 2009. Distributional equivalence and subcompositional
522 coherence in the analysis of compositional data, contingency tables and ratio-scale
523 measurements. *J. Classif.* 26, 29–54. doi:10.1007/s00357-009-9027-y.

524 Gringarten, E., Deutsch, C. V., 2001. Teacher's Aide Variogram Interpretation and
525 Modeling. *Math. Geol.* 33, 507–534. doi:10.1023/a:1011093014141.

526 Guagliardi, I., Cicchella, D., De Rosa, R., 2012. A geostatistical approach to assess
527 concentration and spatial distribution of heavy metals in urban soils. *Water. Air. Soil*
528 *Pollut.* 223, 5983–5998. doi:10.1007/s11270-012-1333-z.

529 Isaaks, E.H., Srivastava, R.M., 1989. An introduction to applied geostatistics. University
530 Press, New York: Oxford. pp. 278–322.

531 Journel, A.G., Huijbregts, C.J., 1978. *Mining Geostatistics*. Academic Press, San Diego.

532 Lark, R.M., Bishop, T.F., 2007. Cokriging particle size fractions of the soil. *Eur. J. Soil*
533 *Sci.* 58, 763–774. doi:10.1111/j.1365-2389.2006.00866.x.

534 Li, Z., Ma, Z., van der Kuip, T.J., Yuan, Z., Huang, L., 2014. A review of soil heavy metal
535 pollution from mines in China: Pollution and health risk assessment. *Sci. Total*
536 *Environ.* doi:10.1016/j.scitotenv.2013.08.090.

537 Lu, A., Wang, J., Qin, X., Wang, K., Han, P., Zhang, S., 2012. Multivariate and
538 geostatistical analyses of the spatial distribution and origin of heavy metals in the
539 agricultural soils in Shunyi, Beijing, China. *Sci. Total Environ.* 425, 66–74.
540 doi:10.1016/j.scitotenv.2012.03.003.

541 Martínez, J., Saavedra, Á., García-Nieto, P.J., Piñeiro, J.I., Iglesias, C., Taboada, J.,
542 Sancho, J., Pastor, J., 2014. Air quality parameters outliers detection using
543 functional data analysis in the Langreo urban area (Northern Spain). *Appl. Math.*
544 *Comput.* 241, 1–10. doi:10.1016/j.amc.2014.05.004.

545 Mateu-Figueras, G., Pawlowsky-Glahn, V., 2008. A critical approach to probability laws
546 in geochemistry, in: *Progress in Geomathematics*. pp. 39–52. doi:10.1007/978-3-
547 540-69496-0_4.

548 Matheron, G., 1963. Principles of geostatistics. *Econ. Geol.* 58, 1246–1266.
549 doi:10.2113/gsecongeo.58.8.1246.

550 McIlwaine, R., Doherty, R., Cox, S.F., Cave, M., 2016. The relationship between
551 historical development and potentially toxic element concentrations in urban soils.
552 Environ. Pollut. 220, 1036–1049. doi:10.1016/j.envpol.2016.11.040.

553 McKinley, J.M., Hron, K., Grunsky, E.C., Reimann, C., de Caritat, P., Filzmoser, P., van
554 den Boogaart, K.G., Tolosana-Delgado, R., 2016. The single component
555 geochemical map: Fact or fiction? J. Geochemical Explor. 162, 16–28.
556 doi:10.1016/j.gexplo.2015.12.005.

557 Megido, L., Suárez-Peña, B., Negral, L., Castrillón, L., Fernández-Nava, Y., 2017.
558 Suburban air quality: Human health hazard assessment of potentially toxic
559 elements in PM10. Chemosphere 177, 284–291.
560 doi:10.1016/j.chemosphere.2017.03.009.

561 Moen, J., Ale, B.J.M., 1998. Risk maps and communication. J. Hazard. Mater. 61, 271–
562 278.

563 Odeh, I.O.A., Todd, A.J., Triantafyllis, J., 2003. Spatial prediction of soil particle size
564 fractions as compositional data. Soil Sci. 168, 501–515.

565 Pawlowsky, V., 1989. Cokriging of regionalized compositions. Math. Geol. 21, 513–521.
566 doi:10.1007/BF00894666.

567 Pawlowsky, V., Burger, H., 1992. Spatial Structure-analysis of Regionalized
568 Compositions. Math. Geol. 24, 675–691. doi:10.1007/BF00894233.

569 Pawlowsky, V., Olea, R.A., 2004. Geostatistical Analysis of Compositional Data,
570 Chapman & Hall, Ltd., London, UK.

571 Pawlowsky, V., Olea, R.A., Davis, J., 1995. Estimation of regionalize compositions: A
572 comparison of three methods. Math. Geol. 27, 105-127.

573 Pawlowsky-Glahn, V., Egozcue, J.J., 2006. Compositional data and their analysis: an
574 introduction. Geol. Soc. London, Spec. Publ. 264, 1–10.
575 doi:10.1144/GSL.SP.2006.264.01.01.

576 Pawlowsky-Glahn, V., Egozcue, J.J., 2001. Geometric approach to statistical analysis
577 on the simplex. *Stoch. Environ. Res. Risk Assess.* 15, 384–398.
578 doi:10.1007/s004770100077.

579 Pawlowsky-Glahn, V., Buccianti, A., 2011. *Compositional Data Analysis: Theory and*
580 *Applications*. Wiley.

581 Reimann, C., Filzmoser, P., Fabian, K., Hron, K., Birke, M., Demetriades, A., Dinelli, E.,
582 Ladenberger, A., Albanese, S., Andersson, M., Arnoldussen, A., Baritz, R., Batista,
583 M.J., Bel-lan, A., Cicchella, D., De Vivo, B., De Vos, W., Duris, M., Dusza-Dobek,
584 A., Eggen, O.A., Eklund, M., Ernsten, V., Finne, T.E., Flight, D., Forrester, S.,
585 Fuchs, M., Fugedi, U., Gilucis, A., Gosar, M., Gregorauskiene, V., Gulan, A.,
586 Halamic, J., Haslinger, E., Hayoz, P., Hobiger, G., Hoffmann, R., Hoogewerff, J.,
587 Hrvatovic, H., Husnjak, S., Janik, L., Johnson, C.C., Jordan, G., Kirby, J., Kivisilla,
588 J., Klos, V., Krone, F., Kwecko, P., Kuti, L., Lima, A., Locutura, J., Lucivjansky, P.,
589 Mackovych, D., Malyuk, B.I., Maquil, R., McLaughlin, M.J., Meuli, R.G., Miosic, N.,
590 Mol, G., Négrel, P., O'Connor, P., Oorts, K., Ottesen, R.T., Pasiieczna, A., Petersell,
591 V., Pfeleiderer, S., Ponavic, M., Prazeres, C., Rauch, U., Salpeteur, Schedl, A.,
592 Scheib, A., Schoeters, I., Sefcik, P., Sellersjö, E., Skopljak, F., Slaninka, I., Šorša,
593 A., Srvkota, R., Stafilov, T., Tarvainen, T., Trendavilov, V., Valera, P.,
594 Verougstraete, V., Vidojevic, D., Zissimos, A.M., Zomeni, Z., 2012. The concept of
595 compositional data analysis in practice - Total major element concentrations in
596 agricultural and grazing land soils of Europe. *Sci. Total Environ.* 426, 196–210.
597 doi:10.1016/j.scitotenv.2012.02.032.

598 Tepanosyan, G., Maghakyan, N., Sahakyan, L., Saghatelyan, A., 2017. Heavy metals
599 pollution levels and children health risk assessment of Yerevan kindergartens soils.
600 *Ecotoxicol. Environ. Saf.* 142, 257–265. doi:10.1016/j.ecoenv.2017.04.013.

601 Venkatramanan, S., Chung, S.Y., Kim, T.H., Kim, B.W., Selvam, S., 2016. Geostatistical
602 techniques to evaluate groundwater contamination and its sources in Miryang City,

603 Korea. Environ. Earth Sci. 75. doi:10.1007/s12665-016-5813-0.

604 Zuo, R., Carranza, E.J.M., Wang, J., 2016. Spatial analysis and visualization of
605 exploration geochemical data. Earth-Science Rev. 158, 9–18.
606 doi:10.1016/j.earscirev.2016.04.006.

607