



## Establishment of background pollution levels and spatial analysis of moss data on a regional scale



Pablo Giráldez <sup>a,\*</sup>, Rosa M. Crujeiras <sup>b</sup>, J. Ángel Fernández <sup>a</sup>, Jesús R. Aboal <sup>a</sup>

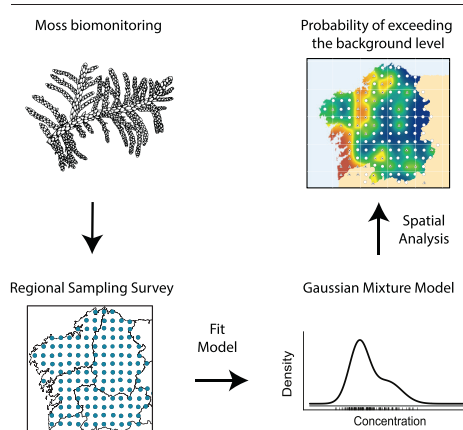
<sup>a</sup> CRETUS, Ecology Area, Department of Functional Biology, Faculty of Biology, Universidade de Santiago de Compostela, 15782 Santiago de Compostela, Spain

<sup>b</sup> CITMAGA, Universidade de Santiago de Compostela, 15782 Santiago de Compostela, Spain

### HIGHLIGHTS

- New method for modeling background level distribution based on Gaussian mixtures.
- Assignment of sampling stations (SS) to background or above-background distribution.
- Qualitative treatment of moss SSs to obtain above-background probability maps.

### GRAPHICAL ABSTRACT



### ARTICLE INFO

Editor: Elena Paoletti

#### Keywords:

Atmospheric pollution  
Biomonitoring  
Bryophytes  
Heavy metals  
Mercury  
Normal mixture model

### ABSTRACT

The determination of background pollution levels is fundamental for the interpretation of the results obtained from environmental biomonitoring. In this paper we propose a new probabilistic method, based on a Gaussian mixture model, for determining the distribution of regional background levels of different pollutants. The distribution of the reference level is used to categorize the observations as “background” or “above-background” and spatial statistical techniques are then applied to determine the probability of the background level being exceeded. To exemplify its use, we applied the method to concentrations of five potentially toxic elements (Cd, Cu, Hg, Pb and Zn) measured in the moss *Pseudoscleropodium purum*. The proposed method was applied to data resulting from sampling at ca. 150 sampling stations in a regular grid (15 × 15 km) in Galicia (NW Spain). Sampling was carried out in June in 2000 and 2002, and in March and September in 2004, 2006, 2008 and 2014. The proposed method yielded consistent results for all of the different sampling surveys, and the pollution levels were found to be closely related to the sources of pollution identified in the study region. In short, although not an optimal solution, the proposed method seems to be suitable and realistic for the qualitative assessment of regional pollution.

### 1. Introduction

Anthropogenic pollution currently represents one of the main threats facing humanity in the 21st century (Mishra et al., 2019; Thompson and

Darwish, 2019; Manisalidis et al., 2020). In this context, assessment of environmental pollution, whether atmospheric (e.g., Osborne et al., 2021), marine (e.g., Tornero and Ribera d'Alcalà, 2014; Ausili et al., 2020), fluvial (e.g., Giri, 2021) or edaphic (e.g., Huang et al., 2020; Cai et al., 2021), has

\* Corresponding author.

E-mail address: [pablo.giraldez.suarez@usc.es](mailto:pablo.giraldez.suarez@usc.es) (P. Giráldez).

become imperative. A multitude of tools and techniques that enable monitoring (e.g., Shellaiah and Sun, 2021; Yang et al., 2021) and biomonitoring (e.g., Hoang et al., 2021; González-Rubio et al., 2021) of environmental pollution have been developed. However, the concentrations of pollutants determined by these techniques include both anthropogenic and natural inputs (Hopke et al., 2020; Giráldez et al., 2021), thus hampering the interpretation of the results. In order to address this issue, researchers often try to establish background levels (Dung et al., 2013) when applying these techniques.

Although the type of background level (also referred to as baseline level) is rarely specified, the term can be defined in different ways. For example, Carballeira et al. (1997) identified 4 types of background levels: pre-industrial level (before any human activity), natural level (representative of an area with human activity but which is well preserved), standard level (reference value for a region) and pre-operational or zero state level (before a particular activity, regardless of the state of conservation). The natural background level seems the most appropriate reference level for monitoring changes in the environment at the regional level as it fluctuates with the environmental conditions, thus allowing observation of changes in regional pollution levels and identification of trends in the levels. In previous studies, different empirical and statistical methods have been used to determine background levels (Dung et al., 2013; Birch, 2017). Regardless of the system used, once the background level has been estimated, a contamination threshold is generally established from a representative value.

The moss monitoring technique has been used for more than 40 years to biomonitor air pollution (Rühling and Tyler, 1968), but its use is mainly restricted to Europe (Boquete et al., 2017). Unlike other biomonitoring methods (e.g., use of macroinvertebrates for water quality assessment), the moss technique has not been included in the European legislation. One possible reason for the lack of success in implementing this technique is that the data obtained (potentially toxic element -PTE- concentrations in moss) have traditionally been interpreted quantitatively and as integrative observations of the concentrations or levels of PTEs in the air. However, the lack of significant correlations between pollutant concentration in moss and concentrations of many PTEs in bulk deposition (Aboal et al., 2010; Boquete et al., 2015), as well as the high spatial (e.g., Fernández et al., 2007; Varela et al., 2014) and temporal variability (Boquete et al., 2011) in moss concentrations, among other aspects (see Boquete et al., 2017), make this approach inappropriate. Researchers familiar with the moss monitoring technique are aware that false negative results can occur (i.e., high atmospheric levels of pollutants and low concentrations in moss). However, as false positives do not occur, the method is valid for differentiating contaminated sites from sites characterised by low levels of pollution. Therefore, a qualitative approach that enables use of moss to determine whether or not pollution is present in concentrations exceeding the background level would be preferable (Boquete et al., 2015).

Empirical methods, such as sampling pristine or unpolluted areas (e.g., Steinnes et al., 1997) or selecting the lowest concentrations of pollutants (e.g., Yan et al., 2016), have been widely used to estimate the natural background levels of PTE concentrations in mosses. The result obtained (representative value of the background level) is usually used together with the concentrations of the pollutants to calculate pollution indices or enrichment factors (e.g., Fernández and Carballeira, 2001). Although a few studies consider the associated variability, that may be different for each pollutant and each sampling (e.g., Loppi et al., 2021), most of them only determine a mean value. Threshold values that differentiate between contaminated and clean samples cannot be established using a common empirical factor and they must be determined using a probabilistic method, i.e., by considering the distribution of the data and not only mean values.

In this study we propose a new statistical method for determining natural background levels of pollution at the regional scale (as the distribution of a subset of samples). This probabilistic method will enable researchers to determine whether the concentrations measured at a sampling site (SS) exceed the previously established background level. In addition, the

method enables qualitative analysis of the results in order to predict whether the pollution levels at any point within a sampling region exceed the natural background level. This method is based on the hypothesis that at a regional scale the natural background level of a PTE has a Gaussian distribution. Although there is not yet literature to support this hypothesis, empirical evidence seems to strongly support it. In previous studies (Boquete et al., 2011; Varela et al., 2011), we observed that PTE concentrations are normally distributed (Table A.1) in uncontaminated SS, either at different subsamples within the SS (Fig. 1D, E) or at a same SS over time (Fig. 1G, H). This implies that the natural background levels will also be normally distributed at the regional scale, as they are simply the sum of the sources of variability corresponding to uncontaminated SSs. This is confirmed by the fact that the differences in PTE concentrations in unpolluted SSs separated by short distances (ca. 1 km; see: Fernández et al., 2007 for more details) are normally distributed around a mean value of 0, which reinforces the hypothesis of normality. The same applies to metabolically regulated nutrients in moss, which tend to be normally distributed, at both regional (e.g., K in Fig. 1C and N in Varela et al., 2013) and local scales (i.e., in the same SS, Fig. 1F). This similarity in the distributions may indicate that PTEs, when present at low concentrations, are also metabolically regulated (although not as strongly as nutrients). Nonetheless, and although the underlying biochemical mechanisms have not yet been identified, the natural background levels are clearly normally distributed at the regional scale.

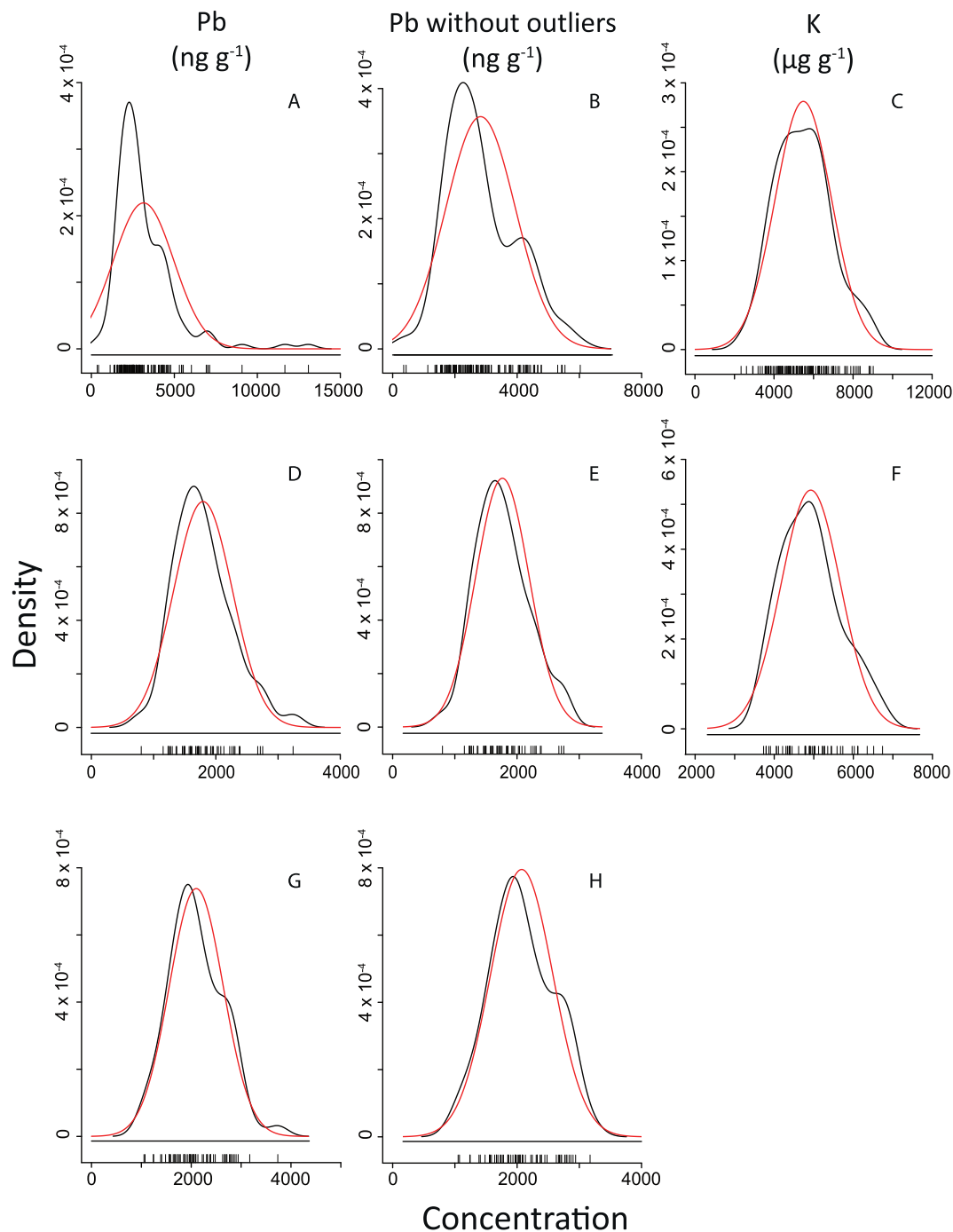
For regional scale sampling, some polluted SS in which pollutant concentrations exceed the natural background level are likely to be present and, therefore the distribution will be skewed to the right (Fig. 1A, B). The regional distribution is, therefore, a mixed distribution in which the concentration of the natural background level can be interpreted as “background noise” with a normal distribution and the concentrations exceeding that background level, as pollution “signals” or “flags”. Thus, a Gaussian model could be fitted where the first one, i.e., the one with the lowest mean, would be the natural background level. A lognormal distribution model (Ott, 1995) could be fitted to all observations at the regional level, even including extremely high polluted SSs. However, this is not appropriate, as the concentrations of contaminated and uncontaminated observations have different origins. Moreover, this approach would not enable the two groups to be differentiated, which is essential for assessing contamination at the regional level.

The idea of using a Gaussian mixture model to describe the distribution of pollutants at regional level and thus obtain the natural background level is not new, and has previously been applied (Carballeira et al., 2002) using the NORMSEP software (Tomlinson, 1971). However, in this case, the method was only used to calculate modal value and not to estimate the mixed distributions to enable each observation to be assigned a probability of belonging to each component. In the present study, we provide a more accurate and updated version of the method, including advanced statistical techniques. Furthermore, the novel aspect of the present study is the use of the background level to assign classes to the SSs and the subsequently qualitative processing of the results, which, as previously mentioned, is the most appropriate way of dealing with PTE concentrations in moss.

## 2. Material and methods

### 2.1. Sampling and processing

Samples of the moss *Pseudoscleropodium purum* (Hedw.) M. Fleisch were collected following the recommendations of Fernández et al. (2015). The samples were collected in June of 2000 and 2002 and in March and September of 2004, 2006, 2008 and 2014, from ca. 150 sampling sites (SSs) in Galicia (NW Spain; Fig. 2A, B). The SSs are located at the nodes of a quasi-regular  $15 \times 15$  km grid covering the entire region (for more details, see Fernández et al., 2005). Briefly, whenever possible, the samples were collected at distances of at least 300 m from main roads, 100 m from secondary roads, 4 km from industries and 3 km from cities. The samples were collected in open areas, or where this was not possible, in forest

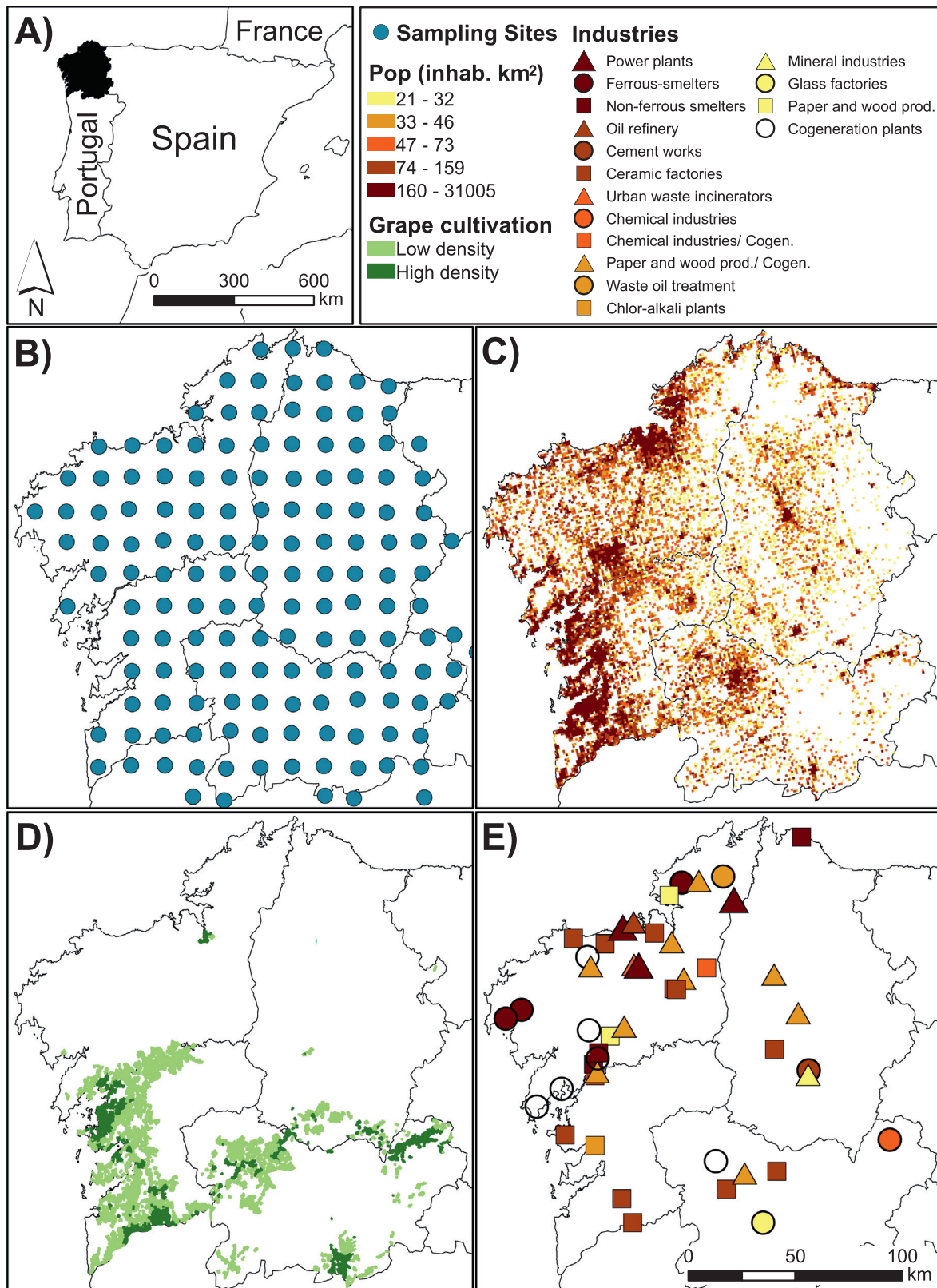


**Fig. 1.** Plots of nonparametric estimates of density (black) and normal distributions (red) corresponding to sample means and variances. A, B and C denote the 2000 sampling in Galicia ( $n = 132$ ) (data from Aboal et al., 2006). D, E and F are the result of sampling of spatial variability conducted in an unpolluted sampling site ( $n = 50$ ) (data from Varela et al., 2011). G and H correspond to a study of temporal variability conducted in another unpolluted sampling site ( $n = 62$ ) (data from Boquete et al., 2011). The first column includes the results for Pb with all the data, while the second column includes the results without the outliers. The third column includes the K results.

clearings, as far as possible from trees. A minimum of 30 subsamples of similar weight were collected in each SS over an area of between 201 m<sup>2</sup> (circular area,  $\varnothing > 16$  m) and 2500 m<sup>2</sup> (50 × 50 m). The subsamples were combined to form a single composite sample for each SS. In the laboratory, moss samples were cleaned to remove remains of adhered material, and the green parts were cut from the shoots. The material was dried in an oven at 45 °C and homogenized in a metal-free ultracentrifuge mill (Retsch ZM200, Retsch GmbH). Finally, the homogenized material was stored in glass vials for analytical determination.

## 2.2. Chemical analysis

One aliquot of the moss samples was digested in HNO<sub>3</sub> (Hiperpur) and H<sub>2</sub>O<sub>2</sub>, in a microwave oven (Ethos-1, Milestone). An undigested aliquot of each moss sample was used to determine the Hg content, in a mercury analyzer (Milestone DMA80). The concentrations of Cd, Cu, Pb and Zn in the digested samples were determined by inductively coupled plasma mass spectrometry (Agilent 7700 ×). The determinations were carried out at the Research Support Services Unit and the Ecology Unit, University of



**Fig. 2.** Maps showing (A) the location of Galicia in Spain, (B) the location of sampling sites in and around Galicia, (C) the population density (inhabitants per square kilometres) (Source: Galician Institute of Statistics, year 2018), (D) the area under grape cultivation in Galicia, according to vineyard density (source: Institute of Territorial Studies of the Regional Ministry of the Environment, Territory and Infrastructures, Xunta de Galicia, year 2000), and (E) location of the main industries labelled according to their main productive sector (E) (source: E-PRTR emissions inventory, modified according to the text).

Santiago de Compostela. The concentrations of almost all elements were above the limits of quantification (LOQ) of the analytical technique, except for Cu in one sample and Cd in 7 samples. In these cases, values equal to half of the corresponding LOQ were used. To check the analytical quality, one sample of certified reference material (M2 and M3, *Pleurozium schreberi*, Steinnnes et al., 1997; Poplar leaves, GBW07604, National Institute of Metrology, China) was analysed every ten samples. In addition, one of every 10 samples was also remeasured (duplicate analysis). Analytical blanks were also analysed (one every ten samples) to test for any possible contamination. The analytical quality of the process was satisfactory; the overall percentage error was between 3% and 10% in most cases (except for except for the 2002 sampling campaign: 14–18%, Cd in 2004S: 21% and Hg in 2014: 12%). The percentage recovery of the reference materials was between ca. 70 and 128% in most cases, with the exception of Cu and Zn in 2014 (69% and 66%, respectively), and Pb in 2000 (65%).

### 2.3. Statistical analysis

We present here a method for determining whether the concentration of a certain element in a SS can be assigned to the background distribution or whether it corresponds to an “above-background” sample. This term includes those SSs that are clearly contaminated and exceed the background contamination level. The assignment is based on fitting a Gaussian mixture model to the observed data, as described below. Finally, an indicator kriging method was used to create a global prediction map for the region from the dichotomized observations (i.e., each SS is categorised as background or “above-background”).

As a starting point in the procedure, we fitted a Gaussian mixture model to the data sample, after removing outliers (clearly contaminated points exceeding the upper limit of outliers). The Gaussian mixture density model with  $k$  components can be written as follows:

$$f(y) = \sum_{i=1}^k \pi_i \varphi(y; \mu_i, \sigma_i) \tag{1}$$

where  $\varphi$  denotes the normal density with mean  $\mu_i$  and standard deviation  $\sigma_i$  and the weights  $\pi_i$  are positive values summing up to 1. Given that it is not known to which mixture component a certain observation should be assigned, parameter estimation in (1) is usually done by considering the Expectation-Maximization (EM) algorithm introduced by Dempster et al. (1977). However, as a first step in trying to fit a mixture distribution to a data sample, a decision must be made about how many components  $k$  to include in (1). This is done by a likelihood ratio test which compares the likelihood of a fit with  $k = k_0$  components with another fit with  $k = k_0 + 1$ .

Specifically, a likelihood ratio test with a test statistic given by Eq. (2) is applied:

$$-2 \log \lambda = 2 \left\{ \log L(\hat{\theta}_1) - \log L(\hat{\theta}_0) \right\}, \tag{2}$$

where  $L$  is the likelihood function of the corresponding mixture with  $k_0$  or  $k_0 + 1$  components and  $\hat{\theta}_0$  and  $\hat{\theta}_1$  are the estimated parameters in the mixture model. Note that the collection of parameters includes  $k$  mean values,  $k$  standard deviations and  $(k-1)$  weights.

The likelihood ratio test is applied sequentially ( $k_0 = 1$  vs  $k_1 = 2$ ,  $k_0 = 2$  vs  $k_1 = 3$ ...) until the null hypothesis is no longer rejected, i.e., until the  $p$ -value obtained in the test is greater than the established significance level (in this case 5%). Although for a likelihood ratio statistic, an asymptotic chi-square distribution can be used to obtain the  $p$ -values, in practice for mixtures of distributions, the classical regularity conditions are lost, which makes application of a bootstrap algorithm necessary. Hence, the test calibration is done using a bootstrap method, as shown in appendix (A.1). We used the boot.comp function in the R mixtools package (Benaglia et al., 2009) to carry out this test.

Once the value of the number of components  $k$  has been determined, estimation of the parameters in the mixture model (1) was done using the EM

algorithm. The normalmixEM function of the R mixtools package (Benaglia et al., 2009) was used in this step.

Finally, the goodness-of-fit of the proposed model was checked using a smoothing-based test. The parametric estimator obtained with the EM algorithm was compared using a nonparametric approach given by a kernel density estimator. Under the null hypothesis that the distribution of the sample data will be described by the mixture model with the previously determined number of components  $k$ , we obtained estimated values of the parameters (means and standard deviations of each component and weights).  $f_n$  was used to denote the parametric estimator under the null hypothesis, and a kernel density estimator, namely  $f_{n,h}$  was considered under the (general) alternative. A suitable test statistic has been proposed by Pavia (2015):

$$T_{n,h} = \int_{-\infty}^{\infty} |f_{n,h}(y) - f_n(y)| dy \tag{3}$$

Computation of the test statistic in Eq. (3) requires numerical integration methods and calibration. In practice, as the test considers a nonparametric estimator of the density and this implies that convergence to the asymptotic distribution occurs slowly, calibration in practice must be conducted using bootstrap methods. The algorithm considered for calibration is detailed in appendix (A.2). This test was carried out using the dgeometric.test function of the R GoFKernel package (Pavia, 2015).

Finally, once the model was obtained and validated, the observations that were assigned to the first component, according to the probabilities obtained by the EM algorithm, were classified as “background” while the rest of the observations (including outliers removed prior to the normal mixture fitting) were classified as “above-background”.

Once the SSs were classified, each location was assigned a label indicating whether it was background or not. Binary kriging can be used to obtain predictions regarding the probability of a SS being “above-background”. This tool provides a probability map that may enable identification of “above-background” areas. For implementation of the tool, computation of a valid (parametric) variogram function that provides an acceptable fit for the data dependence structure is required. On the other hand, some hypothesis (stationarity, isotropy) should be checked on the sample. For both the computation of a valid variogram function and assessment of the hypothesis, a nonparametric variogram is required. A robust empirical variogram (in a transformed scale) is constructed on the basis of the square root of the absolute value of the differences between the data:

$$\hat{\gamma}^*(h) = \frac{1}{|N(h)|} \sum_{(i,j) \in N(h)} \sqrt{|y_i - y_j|}, \tag{4}$$

where for each pair of observations  $(i,j)$  at a distance  $h$  a set of pairs of locations with similar distance is computed,  $N(h) = \{(i,j) : h_{ij} \in b(h)\}$ , where  $b(h)$  is an interval containing  $h$  and  $|N(h)|$  denotes the cardinal of  $N(h)$ .

In practice, from a realization of a spatial process, one can estimate this variogram on the transformed scale, obtaining a point cloud, which in turn can be smoothed by the usual techniques used in regression (e.g., kernel or spline methods). From the smoothed version, the spatial dependence of the sample is analysed by means of a hypothesis test, so that if the null hypothesis of independence is rejected, the spatial analysis is continued; otherwise, it is not continued. If the data are independent, then the variogram should be flat, and this is the idea on which the contrast proposed by Diblasi and Bowman (1997) and later extended to assess stationarity and isotropy by Bowman and Crujeiras (2013) is based. To perform the independence test, we used the sm.variogram function of the R package sm (Bowman and Azzalini, 2021).

Once all the tests were performed, we fitted a parametric semivariogram (exponential in this case) similar to the empirical semivariogram. For this purpose, we used the fit.variogram function of the gstat R package (Pebesma, 2004). When it was not possible to achieve convergence of the parametric semivariogram, we used the autofitVariogram function, in the automap package (Hiemstra et al., 2009).

Binary kriging predictions were obtained from the parametric semivariogram by using indicator kriging (approximated by ordinary kriging) and represented by probability maps. These probabilities were estimated with the krigé function included in the R package gstat (Pebesma, 2004).

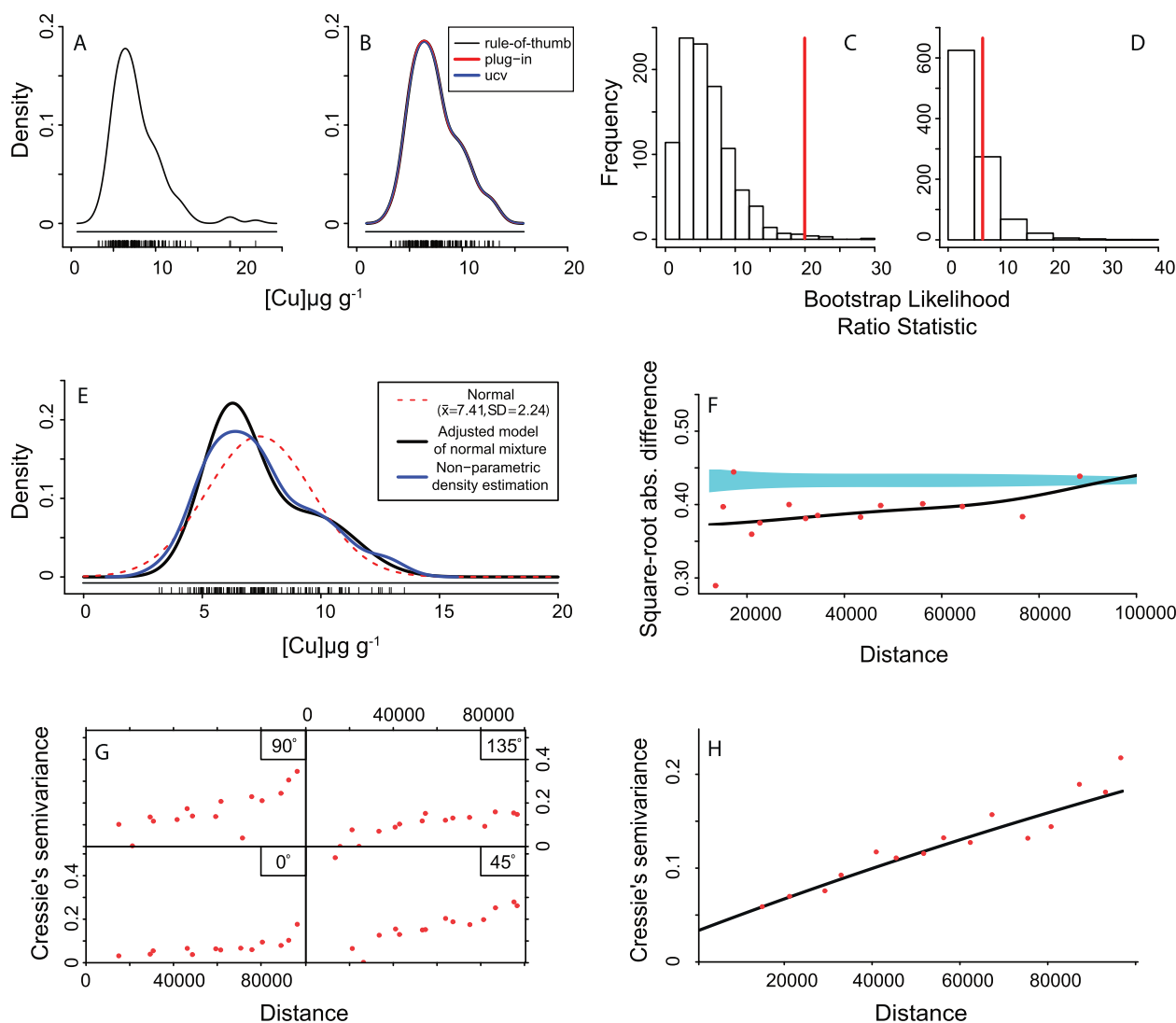
### 3. Results

The proposed method was applied to the data for all elements determined in the sampling surveys considered. To exemplify the process, two cases were chosen, one in which the method was applied without any problems and the other in which various problems emerged. The first case involves Cu data from September 2006 (Cu06S, Fig. 3) and the second case, Cd data from September 2004 (Cd04S, Fig. 4). Both figures show the results for each of the steps of the method, except the goodness-of-fit test, which does not yield a graphical result. Kernel estimation of the density is shown in the upper left corner of each figure (Figs. 3A, B and 4A, B), with all the data and without those data that exceed the upper limit of outlier detection. Kernel estimation require the selection of a smoothing parameter,

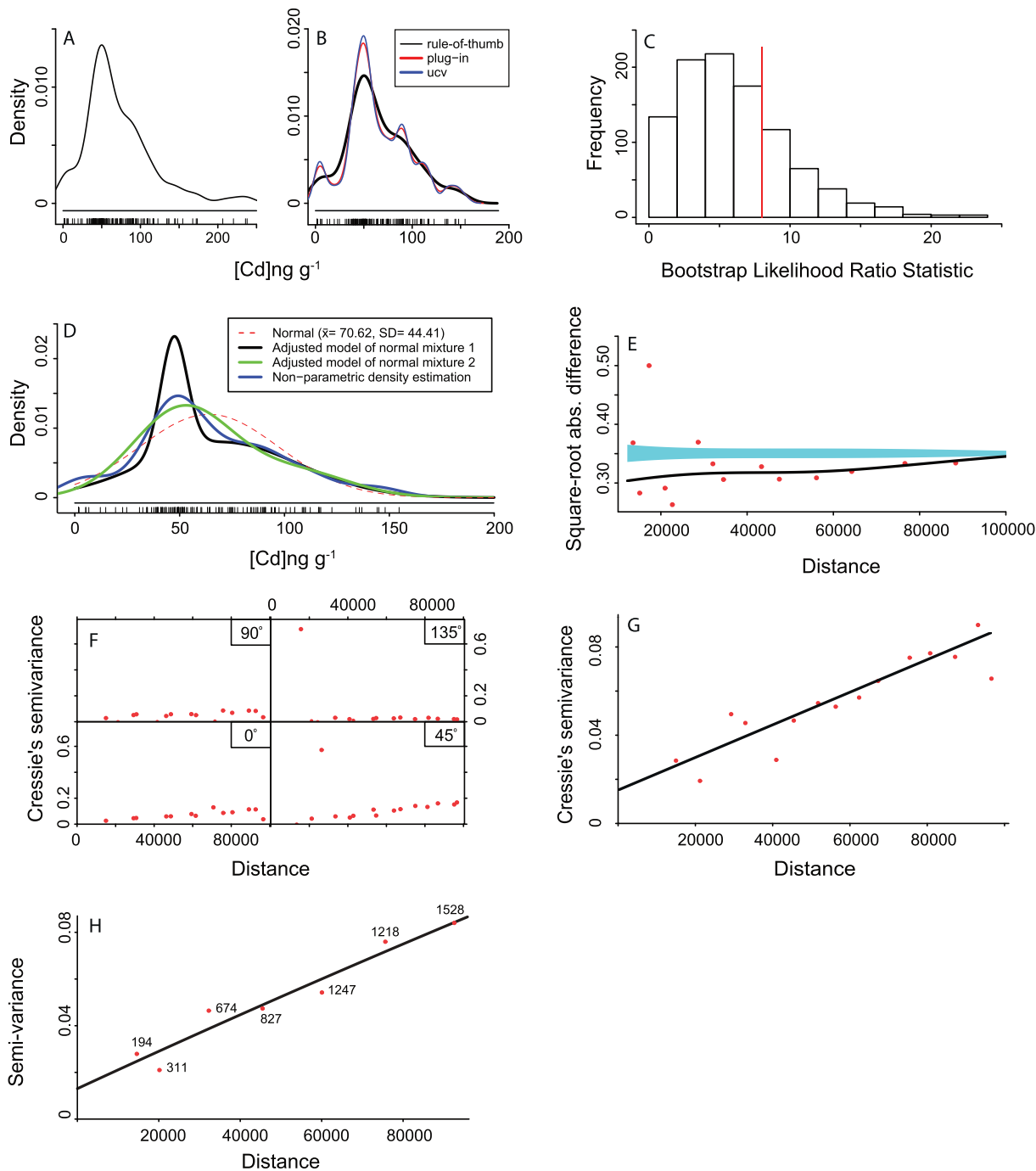
that controls the regularity of the estimated curve (Wand and Jones, 1994). For the whole data set, a rule-of-thumb bandwidth is considered, whereas three different bandwidths are used for the reduced dataset (namely rule-of-thumb, plug-in and ucv windows).

The result of the test to determine the number of components is shown in the upper right corner of Figs. 3 and 4 (Figs. 3C, D and 4C). Two components were selected for Cu06S: in the first test, the null hypothesis ( $k = 1$ ) was rejected ( $pval = 0.005$ ) but in the second test the null hypothesis ( $k = 2$ ) was not rejected ( $pval = 0.231$ ). However, for Cd04S, the null hypothesis ( $k = 1$ ) was not rejected ( $pval = 0.232$ ), and therefore a single component (i.e., a Gaussian model) was selected.

The fitted Gaussian mixture model is plotted together with the nonparametric estimate of the outlier-free sample density and the density of the normal with the mean and standard deviation of the sample, in Figs. 3E and 4D. The fitted model for Cu06S was not rejected in the goodness-of-fit test ( $Tn = 0.11654$ ,  $pval = 0.8119$ ), while the normal model for Cd04S was rejected ( $Tn = 0.22885$ ,  $p-value < 2.2e - 16$ ), and therefore a Gaussian mixture model with two components was fitted using the EM algorithm.



**Fig. 3.** Results of different steps involved in applying the method to the Cu concentrations obtained in the September 2006 sampling survey. (A and B) Kernel density estimates with and without outliers, respectively [(B) the kernel density estimated with different smoothing values]. (C and D) Histograms with the bootstrap result of the contrast on the number of components in the Gaussian mixture model. The value of the observed statistic, corresponding  $p$  values of 0.003 (C) and 0.229 (D) are shown in red. (E) The nonparametric estimate of the density with a two-normal mixture model fitted to the data. The normal distribution with the sample mean and standard deviation is also shown. (F) Result of the contrast regarding spatial independence. The confidence band for spatial independence is indicated in blue. (G) Cressie directional empirical semivariograms. (H) Parametric fit (black line) for the Cressie empirical semivariogram.



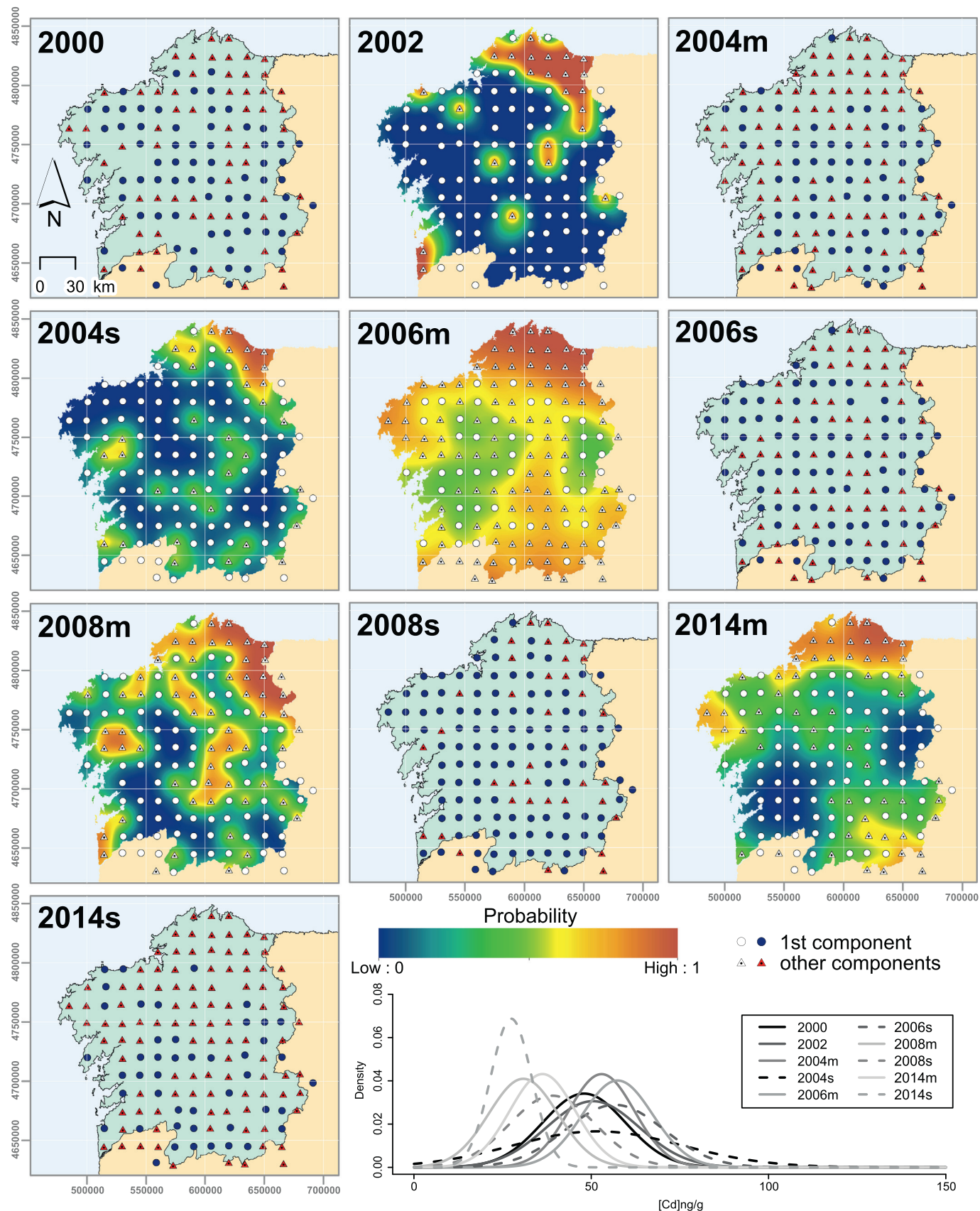
**Fig. 4.** Results of different steps involved in the applying the method to the Cd concentrations obtained in the September 2004 sampling survey. (A and B) Nonparametric density estimates with and without outliers, respectively [(B) the kernel density estimated with different smoothing values]. (C) Histogram with the bootstrap result of the contrast on the number of components in the Gaussian mixture model. The value of the observed statistic, which corresponds to a p value of 0.262, is indicated in red. (D) Nonparametric estimate of the density with the fit of two different mixture models of two normal distributions. The normal distribution, with the sample mean and standard deviation, is also shown. (E) Result of the contrast on spatial independence. The confidence band for spatial independence is shown in blue. (F) Cressie directional empirical semivariograms. (G) Parametric fit (black line) for the Cressie empirical semivariogram. In this case the fit does not converge, so another fit was made, as shown in H.

This algorithm yielded two possible models: model 1 and model 2 (Fig. 4D). Model 2 was selected for the rest of the Cd04S analysis, because it provided the best fit to the estimated density for the data.

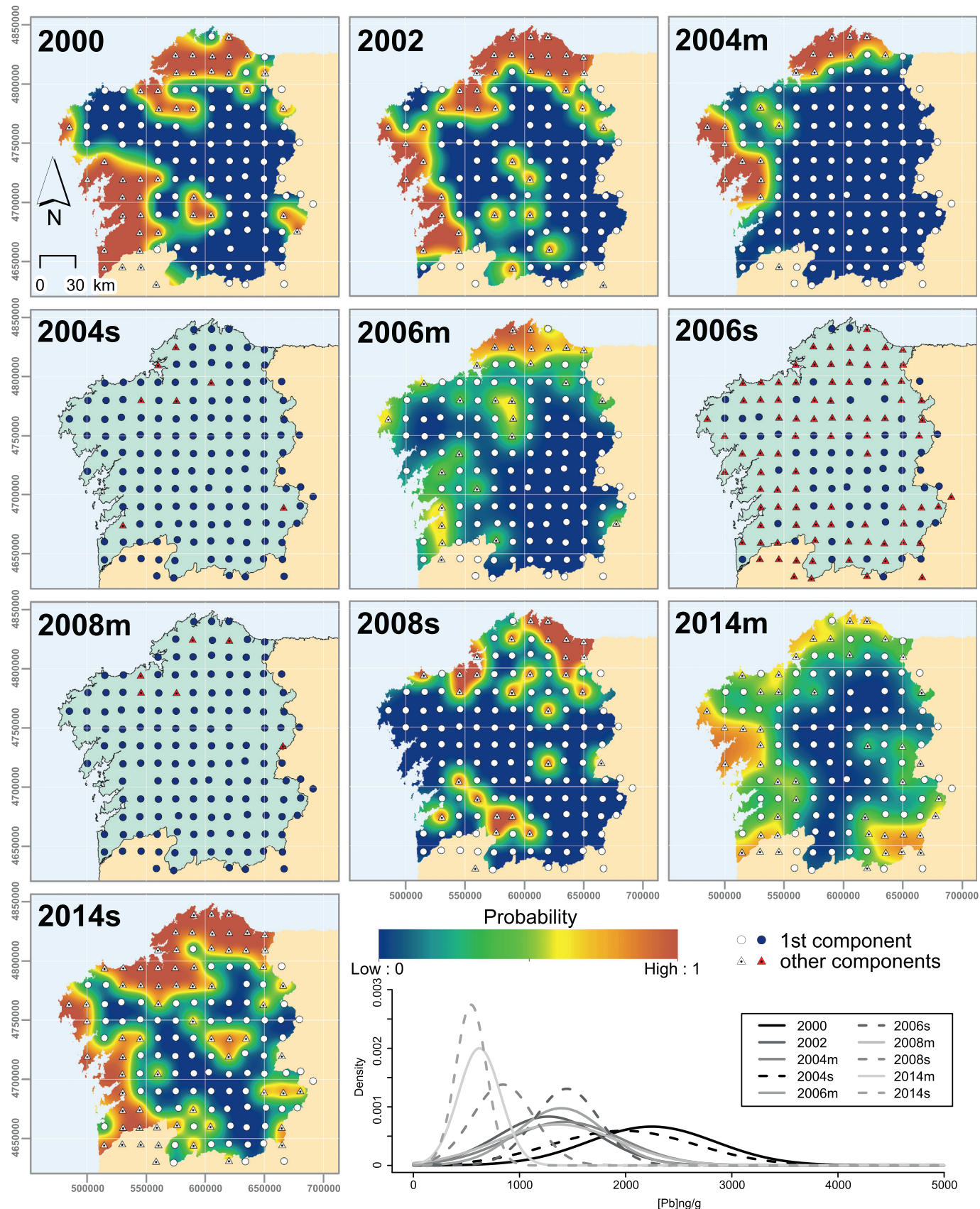
The results of the spatial independence tests are represented in Figs. 3F and 4E. In both cases spatial independence was rejected (pval of 0.030 and 0.037, respectively) and the spatial analysis was thus complete. The resulting Cressie's empirical directional semivariograms (Figs. 3G and 4F) revealed anisotropy in the data. The exponential semivariograms fitted to

the Cressie's empirical semivariograms are shown in Figs. 3H and 4F. In the case of Cd04S, the fit with the fit.variogram function did not converge, and the autofitVariogram function was used; the fit is shown in Fig. 4H.

As the final result of the application of the method, the SSs were classified as "1st component" or "other component" and in those cases in which spatial independence was rejected, probability maps, including the probability of any point on the map belonging to "other component" (Cd and Pb in Figs. 5 and 6, respectively; and Cu, Hg and Zn in Figs. A.1, A.2 and



**Fig. 5.** Maps resulting from the different *Pseudoscleropodium purum* sampling surveys in Galicia. Two surveys were conducted every year (March -m- and September -s-), except in 2000 and 2002 (June only). Sampling sites (SSs) are classified according to whether the Cd concentration ( $\text{ng g}^{-1}$ ) in moss is associated with the 1st component of a Gaussian mixture model fitted to the regional distribution of the element concentrations (circles) or whether it is associated with the other components (triangles). For years when spatial dependence in the classification of SSs was detected, estimates of the probability of not belonging to the 1st component obtained by indicator kriging are plotted. The density of the fitted 1st-component normal distribution for each sampling survey is shown in the graph in the lower right-hand side of the diagram.



**Fig. 6.** Maps resulting from the different *Pseudoscleropodium purum* sampling surveys carried out in Galicia. Two surveys were conducted every year (March -m- and September -s-), except in 2000 and 2002 (June only). Sampling sites (SSs) are classified according to whether the Pb concentration ( $\text{ng g}^{-1}$ ) in moss is associated with the 1st component of a Gaussian mixture model fitted to the regional distribution of the element concentrations (circles) or whether it is associated with the other components (triangles). For years when spatial dependence in the classification of SSs was detected, estimates of the probability of not belonging to the 1st component, as determined by indicator kriging, are plotted. The density of the fitted 1st-component normal distribution for each sampling survey is shown in the graph in the lower right-hand side of the diagram.

A.3, respectively), were constructed. In addition to the maps, these figures also include a plot of the first component of the model fitted to each campaign for each element. These plots clearly show that the background level decreased over time.

#### 4. Discussion

The study findings show that the proposed method allows the regional distribution of pollutant concentrations in moss to be modelled by a Gaussian mixture model. In other words, the method can be used to determine the distribution of the natural background level and also whether the concentrations of elements measured at a SS fall within this distribution. Concentrations exceeding the distribution range indicate that the SS is clearly affected by a local source of pollution that increases the concentration above the background level for the region. In addition to the method of establishing background levels, another innovative contribution of this work is the qualitative treatment of the SS according to the background level. This treatment allows maps of the probability of exceeding the background level to be created, by means of indicator kriging, while avoiding the problems derived from quantitative treatment of the data.

The maps shown in Figs. A.1–A.3 and especially in Figs. 4 and 5 provide a clear, consistent view of which zones in the region include the greatest numbers of “above-background” SSs. The consistency is apparent as areas with a higher probability of exceeding the background level are repeated in different surveys and coincide with the location of pollution hotspots in the region. These hotspots are associated with the most populated area, the western zone, with more than 2 million inhabitants (ca. 77% of the region's inhabitants, Fig. 2C). Different types of pollution are directly associated with the number of inhabitants and are caused by combustion of fossil fuels (by vehicles or in the domestic environment for heating, cooking, etc.), generation of waste and garbage, and creation of debris derived from the demolition and/or construction of buildings. In addition, the presence of agglomerations of people generally implies good communication routes and manpower, leading to the creation of industrial parks on the outskirts of cities, which are a source of industrial pollution (Fig. 2E). In fact, most of the large regional industries that do not fall into the groups “intensive livestock and aquaculture” or “animal and vegetable products of the food and beverage industry” (E-PRTR emissions inventory; [www.prtr-es.es](http://www.prtr-es.es)) are located in the western zone, around the “Atlantic axis”, which crosses the western region from north to south, linking 5 of the 7 Galician cities (Fig. 2C and E). This axis, and the industries established around it, account for a large part of the areas of Galicia identified in this study as being likely to be categorised as “above-background”, mainly for Pb and Cu. It also accounts for many isolated “above-background” SSs, which coincide with point sources of contamination.

In addition to the industries located in industrial parks in large cities or in nearby, well communicated regions, there are a few highly polluting industries in the north of Galicia, including a large thermal power plant (1400 MW) and aluminium and alumina factories. These industries, together with other small industries in this area account for a large proportion of the atmospheric emissions of Cd, Cu, Hg, Pb and Zn emitted by point sources (Fig. 2E). Emissions from these industries lead to the SSs in this region being “above-background”, as consistently observed in the different sampling surveys. Indeed, in 2019, these companies accounted for 57% of Zn emissions from point sources in Galicia, 28% of Pb, 40% of Cu, 30% of Cd and 49% of Hg (calculated from E-PRTR emissions inventory).

The other point sources associated with 2 population centres in remote parts of the region or with other population centres and other diffuse pollution sources may explain the presence of “above-background” SSs outside the western and northern zones. These diffuse sources include the use and exploitation of soils, specifically agriculture, and the type of crop, e.g., the use of copper-containing treatments in viticulture (Fig. 2D) can influence the pattern of “above-background” SSs. The viticulture case can explain the presence of Cu-contaminated SSs in the south and remote inland areas of the region, where more than 87% (calculated from data on soil type according to use and exploitation 2000–2010 from the Ministry of

Agriculture, Fisheries and Food; [www.mapa.gob.es](http://www.mapa.gob.es)) of the land in Galicia dedicated to viticulture is located (Figs. 2D and A.1).

In cases where spatial independence is rejected, the resulting probability maps are easy to interpret. However, even when spatial independence is not rejected, the maps contain useful information. Although in these cases the probability of the presence of “above-background” SSs for the whole territory cannot be predicted, the maps provide an idea of which sites are “above background” and, therefore, which areas may be affected by sources of pollution. Moreover, temporal variations in background levels can also be described by means of this analysis. In the case study, an obvious decrease in concentrations over time was detected for almost all the elements considered (Table A.2).

An example of the direct application of the method is given for Cu in the September 2006 sampling survey (Fig. 3), and the findings can be summarised as follows: i) the windows scarcely affected the nonparametric density estimate; ii) the selection of the number of components was clear; iii) the fitted model density scarcely differed from the nonparametric density estimate; iv) spatial independence was clearly rejected; v) a minor deviation from isotropy was detected; and vi) the parametric semivariogram fitted perfectly (there was convergence) with the fit.variogram function. Nevertheless, for the Cd determined in the September 2004 survey (Fig. 4), some difficulties arose in applying the method: i) it can be seen the window selection effect in the nonparametric estimation of the density; ii) despite selection of 1 component, the fit derived from this selection is rejected in the goodness-of-fit test; iii) when attempting a two-component fit the EM algorithm fits two different models; and iv) when fitting the parametric semivariogram with the selected model there was no convergence with the fit.variogram function, and we therefore used the autofitVariogram function. We must also add the possible problems of the power of the spatial independence test, as may have occurred with Cu in September 2004, as both the layout of the SS and the empirical semivariogram seem to have spatial structure, although the result yielded a *p*-value greater than 0.05. This issue can be explained by examining the (overlapping) variability bands around the semivariogram.

Although all of the problems are unlikely to occur in a single sampling survey, the appearance of any of these problems can complicate the application of the method by requiring decisions to be made by the researcher. One of the most difficult decisions, which can have a major impact on the results, is the choice of the number of components when there is variability in the result of the number of components contrast. The choice of the model is also difficult when for a given number of components there are several possible models. In these cases, it is advisable to fit all of the models (with the different components selected) and then to analyse the estimated parameters, the log-likelihood of the model and its density curve before applying goodness-of-fit tests and comparing the test statistics, which refer to the integral of the difference between the density of the fitted model and the nonparametric estimate of the sample density. If a simple decision regarding the most appropriate model still cannot be reached, the distribution on the map of the SS classified according to whether they belong to the background level must be examined and compared against data already available for the region.

Many of the used statistical tools are non-parametric or based on bootstrap procedures, and their robustness therefore largely depends on the number of samples considered. Therefore, the application of this method requires a substantial number of SSs. As an example, for independent observations in a density estimation context, one should not proceed with less than 30 or 50 observations. The number of sample observations should be larger in the presence of spatial dependence. If the number of SSs too low a few SSs with slightly higher concentrations (probably due to chance) than the others may distort the normal distribution and force a two-component fit.

Although the main points of the proposed method include modeling the distribution as a Gaussian mixture and classification of SS, we must mention the problems that can arise during the spatial analysis. These include the power of the test for spatial independence (discussed above) and the lack of stationarity and isotropy in most of the sampling surveys.

The (nonparametric) testing procedures are flexible tools that enable assessment of the usual hypothesis for kriging interpolation. However, the calibration may not be very accurate, as occurs the independence test, because it is based on nonparametric smoothing procedures. However, if the patterns of non-stationarity or anisotropy are not very strong, the results of the indicator kriging can still be interpreted, although further inferential analysis must be conducted with caution. Although stationarity and isotropy are basic conditions for kriging interpolation, researchers do not usually assess the plausibility of these characteristics for their data samples.

In short, the proposed method represents an innovative approach that enables probabilistic and realistic interpretation (regarding pollution hotspots) of the levels of pollution in a region. Moreover, although the method was used here to analyse PTE pollution, the method could be used for any other type of pollutant, as long as the regional distribution is similar to that described here.

## 5. Conclusion

The proposed method allows the distribution of the background level of pollutants in different sampling sites to be determined and the observations to then be classified according to whether they belong to this distribution (background) or not (above-background). This helps to avoid the problems derived from quantitative analysis of the data and provides consistent results. Some problems may arise when applying the method. However, many of these problems can be solved by the researcher's informed judgment, while others require a more detailed study of the statistical approach.

## CRedit authorship contribution statement

Pablo Giráldez: Methodology, Software, Formal analysis, Data Curation, Writing - Original Draft, Visualization.

Rosa M. Crujeiras: Methodology, Software, Formal analysis, Writing - Review & Editing, Supervision, Project administration.

Jesús R. Aboal: Conceptualization, Validation, Investigation, Resources, Writing - Review & Editing, Visualization, Supervision, Project administration, Funding acquisition.

J. Ángel Fernández: Conceptualization, Validation, Investigation, Resources, Writing - Review & Editing, Supervision, Project administration, Funding acquisition.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

J. R. Aboal, J. A. Fernández and P. Giráldez belong to the Grupo de Referencia Competitiva GRC GI-1252/GPC2020-23 (ED431C 2020/19) which is co-funded by ERDF (EU). Authors would like to thank RIAIDT-USC for the use of analytical facilities. P. Giráldez is grateful to the Spanish Ministerio de Ciencia, Innovación y Universidades for a grant awarded within the Programa de Formación de Profesorado Universitario (FPU 2018 [grant number FPU18/04134]). Research of R. M. Crujeiras has been supported by MINECO (Grant PID2020-116587GB-I00), and by Xunta de Galicia (Grupos de Referencia Competitiva ED431C 2021/24), all of them through the ERDF.

## Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.scitotenv.2022.156182>.

## References

- Aboal, J.R., Fernández, J.A., Boquete, T., Carballeira, A., 2010. Is it possible to estimate atmospheric deposition of heavy metals by analysis of terrestrial mosses? *Sci. Total Environ.* 408 (24), 6291–6297. <https://doi.org/10.1016/j.scitotenv.2010.09.013> Elsevier.
- Aboal, J.R., Real, C., Fernández, J.A., Carballeira, A., 2006. Mapping the results of extensive surveys: The case of atmospheric biomonitoring and terrestrial mosses. *Sci. Total Environ.* 356, 256–274. [10.1016/j.scitotenv.2005.04.025](https://doi.org/10.1016/j.scitotenv.2005.04.025).
- Ausili, A., Bergamin, L., Romano, E., 2020. Environmental status of Italian coastal marine areas affected by long history of contamination. *Front. Environ. Sci.* 8. <https://doi.org/10.3389/fenvs.2020.00034> Frontiers Media S.A.
- Benaglia, T., Chauveau, D., Hunter, D.R., Young, D.S., 2009. Mixtools: an R package for analyzing finite mixture models. *J. Stat. Softw.* 32 (6), 1–29. <https://doi.org/10.18637/jss.v032.i06>.
- Birch, G.F., 2017. Determination of sediment metal background concentrations and enrichment in marine environments – a critical review. *Sci. Total Environ.* 580, 813–831. <https://doi.org/10.1016/j.scitotenv.2016.12.028> Elsevier B.V.
- Boquete, M.T., Aboal, J.R., Carballeira, A., Fernández, J.A., 2017. Do mosses exist outside of Europe? A biomonitoring reflection. *Sci. Total Environ.* 593–594, 567–570. <https://doi.org/10.1016/j.scitotenv.2017.03.196>.
- Boquete, M.T., Fernández, J.A., Aboal, J.R., Carballeira, A., 2011. Analysis of temporal variability in the concentrations of some elements in the terrestrial moss *Pseudoscleropodium purum*. *Environ. Exp. Bot.* 72 (2), 210–216. <https://doi.org/10.1016/j.envexpbot.2011.03.002>.
- Boquete, M.T., Fernández, J.A., Carballeira, A., Aboal, J.R., 2015. Relationship between trace metal concentrations in the terrestrial moss *Pseudoscleropodium purum* and in bulk deposition. *Environ. Pollut.* 201, 1–9. <https://doi.org/10.1016/j.envpol.2015.02.028>.
- Bowman, A.W., Crujeiras, R.M., 2013. Inference for variograms. *Comput. Stat. Data Anal.* 66, 19–31. <https://doi.org/10.1016/j.csda.2013.02.027>.
- Bowman, A.W., Azzalini, A., 2021. R package `sm`: nonparametric smoothing methods (version 2.2-5.7). URL <http://www.stats.gla.ac.uk/~adrian/sm>.
- Cai, X., Duan, Z., Wang, J., 2021. Status assessment, spatial distribution and health risk of heavy metals in agricultural soils around mining-impacted communities in China. *Pol. J. Environ. Stud.* 30 (2), 993–1002. <https://doi.org/10.15244/pjoes/124742> HARD Publishing Company.
- Carballeira, A., Couto, J.A., Fernández, J.A., 2002. Estimation of background levels of various elements in terrestrial mosses from Galicia (NW Spain). *Water Air Soil Pollut.* 133 (1–4), 235–252. <https://doi.org/10.1023/A:1012928518633>.
- Carballeira, A., Carral, E., Puente, X.M., Villares, R., 1997. In: de Copostela, U. de S. (Ed.), *Estado de conservación de la costa de Galicia: Nutrientes y metales pesados en sedimentos y organismo intermareales*.
- Dempster, A.P., Laird, N.M., Rubin, D.B., 1977. Maximum likelihood from incomplete data via the EM algorithm. *J. R. Stat. Soc. Ser. B Methodol.* 39 (1), 1–22. <https://doi.org/10.1111/J.2517-6161.1977.tb01600.x>.
- Diblati, A., Bowman, A.W., 1997. Testing for constant variance in a linear model. *Statist. Probab. Lett.* 33 (1), 95–103. [https://doi.org/10.1016/S0167-7152\(96\)00115-0](https://doi.org/10.1016/S0167-7152(96)00115-0).
- Dung, T.T.T., Cappuyns, V., Swennen, R., Phung, N.K., 2013. From geochemical background determination to pollution assessment of heavy metals in sediments and soils. *Rev. Environ. Sci. Biotechnol.* 12 (4), 335–353. <https://doi.org/10.1007/s11157-013-9315-1> Springer.
- Fernández, J.A., Aboal, J.R., Real, C., Carballeira, A., 2007. A new moss biomonitoring method for detecting sources of small scale pollution. *Atmos. Environ.* 41 (10), 2098–2110. <https://doi.org/10.1016/j.atmosenv.2006.10.072>.
- Fernández, J.A., Boquete, M.T., Carballeira, A., Aboal, J.R. (Eds.), 2015. *A Critical Review of Protocols for Moss Biomonitoring of Atmospheric Deposition: Sampling and Sample Preparation*. 517, p. 132.
- Fernández, J.A., Carballeira, A., 2001. A comparison of indigenous mosses and topsoils for use in monitoring atmospheric heavy metal deposition in Galicia (northwest Spain). *Environ. Pollut.* 114 (3), 431–441. [https://doi.org/10.1016/S0269-7491\(00\)00229-3](https://doi.org/10.1016/S0269-7491(00)00229-3).
- Fernández, J.A., Real, C., Couto, J.A., Aboal, J.R., Carballeira, A., 2005. The effect of sampling design on extensive bryomonitoring surveys of air pollution. *Sci. Total Environ.* 337 (1–3), 11–21. <https://doi.org/10.1016/j.scitotenv.2004.07.011>.
- Giráldez, P., Varela, Z., Aboal, J.R., Fernández, J.A., 2021. Testing different methods of estimating edaphic inputs in moss biomonitoring. *Sci. Total Environ.* 778, 146332. <https://doi.org/10.1016/j.scitotenv.2021.146332>.
- Giri, S., 2021. Water quality prospective in Twenty First Century: status of water quality in major river basins, contemporary strategies and impediments: a review. *Environ. Pollut.* 271, 116332. <https://doi.org/10.1016/j.envpol.2020.116332> Elsevier Ltd.
- González-Rubio, S., Ballesteros-Gómez, A., Asimakopoulos, A.G., Jaspers, V.L.B., 2021. A review on contaminants of emerging concern in European raptors (2002–2020). *Sci. Total Environ.* 760. <https://doi.org/10.1016/j.scitotenv.2020.143337> Elsevier B.V.
- Hiemstra, P.H., Pebesma, E.J., Twenhöfel, C.J.W., Heuvelink, G.B.M., 2009. Real-time automatic interpolation of ambient gamma dose rates from the dutch radioactivity monitoring network. *Comput. Geosci.* 35 (8), 1711–1721. <https://doi.org/10.1016/j.cageo.2008.10.011>.
- Hoang, A.Q., Tu, M.B., Takahashi, S., Kunisue, T., Tanabe, S., 2021. Snakes as bionitors of environmental pollution: a review on organic contaminants. *Sci. Total Environ.* 770. <https://doi.org/10.1016/j.scitotenv.2020.144672> Elsevier B.V.
- Hopke, P.K., Dai, Q., Li, L., Feng, Y., 2020. Global review of recent source apportionments for airborne particulate matter. *Sci. Total Environ.* 740, 140091. <https://doi.org/10.1016/j.scitotenv.2020.140091> Elsevier B.V.
- Huang, J.H., Shetaya, W.H., Osterwalder, S., 2020. Determination of (Bio)-available mercury in soils: a review. *Environ. Pollut.* 263, 114323. <https://doi.org/10.1016/j.envpol.2020.114323> Elsevier B.V.

- Loppi, S., Kosonen, Z., Meier, M., 2021. Estimating background values of potentially toxic elements accumulated in moss: a case study from Switzerland. *Atmosphere* 12 (2), 177. <https://doi.org/10.3390/ATMOS12020177> 2021, Vol. 12, Page 177.
- Manisalidis, I., Stavropoulou, E., Stavropoulos, A., Bezirtzoglou, E., 2020. Environmental and health impacts of air pollution: a review. *Front. Public Health* 8, 14. <https://doi.org/10.3389/fpubh.2020.00014> Frontiers Media S.A.
- Mishra, S., Das, A.P., Charan, Rath C., 2019. Marine microfiber pollution: a review on present status and future challenges. *Mar. Pollut. Bull.* 140, 188–197. <https://doi.org/10.1016/j.marpolbul.2019.01.039> Elsevier Ltd.
- Osborne, S., Uche, O., Mitsakou, C., Exley, K., Dimitroulopoulou, S., 2021. Air quality around schools: Part I - a comprehensive literature review across high-income countries. *Environ. Res.* 196, 110817. <https://doi.org/10.1016/j.envres.2021.110817> Academic Press Inc.
- Ott, W.R., 1995. *Environmental Statistics and Data Analysis*. CRC Press.
- Pavia, J.M., 2015. Testing goodness-of-fit with the kernel density estimator: GoFKernel. *J. Stat. Softw.* 66 (1), 1–27. <https://doi.org/10.18637/jss.v066.c01>.
- Pebesma, E.J., 2004. Multivariable geostatistics in S: the gstat package. *Comput. Geosci.* 30 (7), 683–691. <https://doi.org/10.1016/J.CAGEO.2004.03.012>.
- Rühling, Å., Tyler, G., 1968. An ecological approach to the Lead problem. *Bot. Notiser* 121 (3), 321–342.
- Shellaiah, M., Sun, K.W., 2021. Inorganic-diverse nanostructured materials for volatile organic compound sensing. *Sensors (Switzerland)* 21 (2), 1–61. <https://doi.org/10.3390/s21020633> MDPI AG.
- Steinnes, E., Rühling, Å., Lippo, H., Mäkinen, A., 1997. Reference materials for large-scale metal deposition surveys. *Accred. Qual. Assur.* 2 (5), 243–249. <https://doi.org/10.1007/s007690050141>.
- Thompson, L.A., Darwish, W.S., 2019. Environmental chemical contaminants in food: review of a global problem. *J. Toxicol.* 2019. <https://doi.org/10.1155/2019/2345283> Hindawi Limited.
- Tomlinson, P.K., 1971. NORMSEP: normal distribution separation. In: *FAO Fish Tech (Ed.), Computer Programs for Fish Stock Assessment*. FAO.
- Tornero, V., Ribera d'Alcalà, M., 2014. Contamination by hazardous substances in the Gulf of Naples and nearby coastal areas: a review of sources, environmental levels and potential impacts in the MSFD perspective. *Sci. Total Environ.* 466–467, 820–840. <https://doi.org/10.1016/j.scitotenv.2013.06.106> Elsevier.
- Varela, Z., Aboal, J.R., Carballeira, A., Real, C., Fernández, J.A., 2014. Use of a moss biomonitoring method to compile emission inventories for small-scale industries. *J. Hazard. Mater.* 275, 72–78. <https://doi.org/10.1016/J.JHAZMAT.2014.04.061>.
- Varela, Z., Carballeira, A., Fernández, J.A., Aboal, J.R., 2013. On the use of epigeic mosses to biomonitor atmospheric deposition of nitrogen. *Arch. Environ. Contam. Toxicol.* 64 (4), 562–572. <https://doi.org/10.1007/S00244-012-9866-0/FIGURES/4>.
- Varela, Zulema, Fernández, J.A., Aboal, J.R., Real, C., Carballeira, A., 2011. Determination of the optimal size of area to be sampled by use of the moss biomonitoring technique. *J. Atmos. Chem.* 65 (1), 37–48. <https://doi.org/10.1007/S10874-010-9180-Z/TABLES/3>.
- Wand, M.P., Jones, M.C., 1994. Kernel smoothing. *Kernel Smoothing*. Chapman and Hall/CRC <https://doi.org/10.1201/B14876>.
- Yan, Y., Zhang, Q., Wang, G.G., Fang, Y.M., 2016. Atmospheric deposition of heavy metals in Wuxi, China: estimation based on native moss analysis. *Environ. Monit. Assess.* 188 (6), 1–8. <https://doi.org/10.1007/s10661-016-5315-2>.
- Yang, S., Chen, Z., Cheng, Y., Liu, T., Yin, Lihong, Pu, Y., Liang, G., 2021. Environmental toxicology wars: organ-on-a-chip for assessing the toxicity of environmental pollutants. *Environ. Pollut.* 268, 115861. <https://doi.org/10.1016/j.envpol.2020.115861> Elsevier B.V.