

# Ceramic investigation. How to perform statistical analyses.

## Abstract

The aim of this article is to summarize and organize the statistical methodologies used for the statistical analysis towards ceramic investigation and in particular study of ceramic provenance. An update and review of all related methodologies is provided during the presentation of a typical statistical analysis. The presentation is given in a step-by-step process and emphasis is on interpretation of the intermediate and final results. The analysis attempts to cover the following:

- What issues to examine in a preliminary analysis
- Data transformation
- Cluster Analysis
- Clustering assessment
- Data dimension reduction methods as part of a clustering visualization and assessment
- Outliers and small groups
- Mixed-mode analysis
- Cluster characterization and discriminating factors
- Classification

## 1. Introduction

This paper is a part of a larger project, a series of papers, aiming to provide a guide to a researcher in the study of ceramic materials, ideally in every single stage of the study. Statistical methods in general are used throughout the process of an archaeological survey, from the initial stage of planning the survey and sampling until the stage of data collection. However, with the term statistical analysis in the title of this paper we focus on the last part of the survey, when the data have already been collected. In particular, we refer to the specific problem of analyzing the data obtained from an archaeological survey aiming to answer questions with respect to provenance, as a part of reconstruction of the past. Most of the statistical methodologies used for this problem fall within the branch of multivariate statistical analysis. This is due to the fact that the data collected for a provenance study in Archaeology consist of a multivariate data matrix, where the columns, i.e. the variables of the problem are in general correlated and a univariate study would be inadequate. Methods from multivariate statistical analysis would be in particular, data reduction methods, clustering and classification.

Quantitative methods in Archaeology, in general, have a long history. It is widely accepted by the community of Archeology scientists that statistical theory can be a valuable tool, sometimes necessary, in order to plan and organize the survey, handle the volume of data

1  
2  
3  
4 collected, verify or clarify scientific hypotheses and make statistical inference based on the data  
5 at hand. First publications of statistical methods applied to problems in Archaeology appear the  
6 decade of 1960s (Binford, 1964) and more specific publications to provenance and multivariate  
7 statistical techniques appear in the mid of 1970s (Bieber et al., 1976). Baxter (2008) provides a  
8 thorough review on use of Mathematics and Statistics in the last fifty years in Archaeometry.  
9

10  
11 The paper with title ‘The awful truth about statistics in Archaeology’ by D.H. Thomas in 1978  
12 is an indicator of the great expansion of statistical methodology used in Archaeology over the  
13 recent years before the date of its publication. The author although admits that undoubtedly  
14 statistical theory, when correctly used, can assist the researcher to efficiently and subjectively  
15 derive results on the archaeological questions, he discusses the misuses and in some cases abuses  
16 of Statistical techniques. The paper is organized in sections with titles ‘the good’, ‘the bad’ and  
17 ‘the ugly’ presenting in each section from harmless to more serious mistakes in statistical  
18 analysis applied to Archaeology. It is mentioned that the instruction of his editor was to ‘shake  
19 things up in a pleasant way’ and the author definitely accomplished this purpose.  
20  
21

22  
23 Since then, Archaeological scientists have gain more experience in using quantitative tools,  
24 they are aware of the methodologies appropriate for their analysis, they are more educated and  
25 they have access to statistical/computer packages that can assist towards the implementation of  
26 such methodologies. This article does not attempt to present some new statistical methodologies  
27 –although a review on various choices in each case is given- but to provide some guidance with  
28 respect to the sequence of the steps an analysis and the implementation. The instruction of my  
29 editor was to present ‘a tutorial approach –solving problem oriented’ and my contribution to  
30 have ‘an educational character’. I hope this is accomplished. The paper although technical has  
31 kept Mathematics to a minimum and emphasis is on ‘how’, ‘why’ and ‘when’ we use each  
32 method.  
33  
34

35  
36 In particular, with respect to its content the article is organized as follows. In section 2, the  
37 issues we examine during a preliminary analysis are presented. Questions such as why a  
38 preliminary analysis is important and how we can use any conclusions made at this stage in the  
39 subsequent steps of the analysis, are answered. Section 3, is the main part of the paper,  
40 reviewing various approaches for clustering, illustrated in simulated data. Section 4 lists the tools  
41 and the steps the researcher can follow if he wishes to conduct an assessment of the clustering  
42 result(s). Within this section, as part of the clustering assessment, methods for data dimension  
43 reduction are also presented. Moreover, a presentation and proposed solution of special  
44 problems one has to deal with in cluster analysis, such as existence of outliers is given. In section  
45 6 a brief presentation of the classification problem and possible approaches is listed. Finally, a  
46 summary of the steps of the analysis is listed in the last section.  
47  
48

## 49 **2. Preliminary analysis. Data manipulation and transformation**

50  
51 The data collection should normally be a part of a more general process which includes the  
52 questions of the research. These questions, the purpose of the research project, are set at an  
53  
54  
55

1  
2  
3  
4 early stage of this project and the data collection is adjusted so that the data at hand will include  
5 information sufficient for answering the questions. Moreover, the type of the data collected  
6 needs to be in accordance with the methods of data manipulation which are going to be used.  
7 Alternatively, the methods of analysis need to be adapted to the type of data that are meaningful  
8 for the purpose of the analysis.  
9

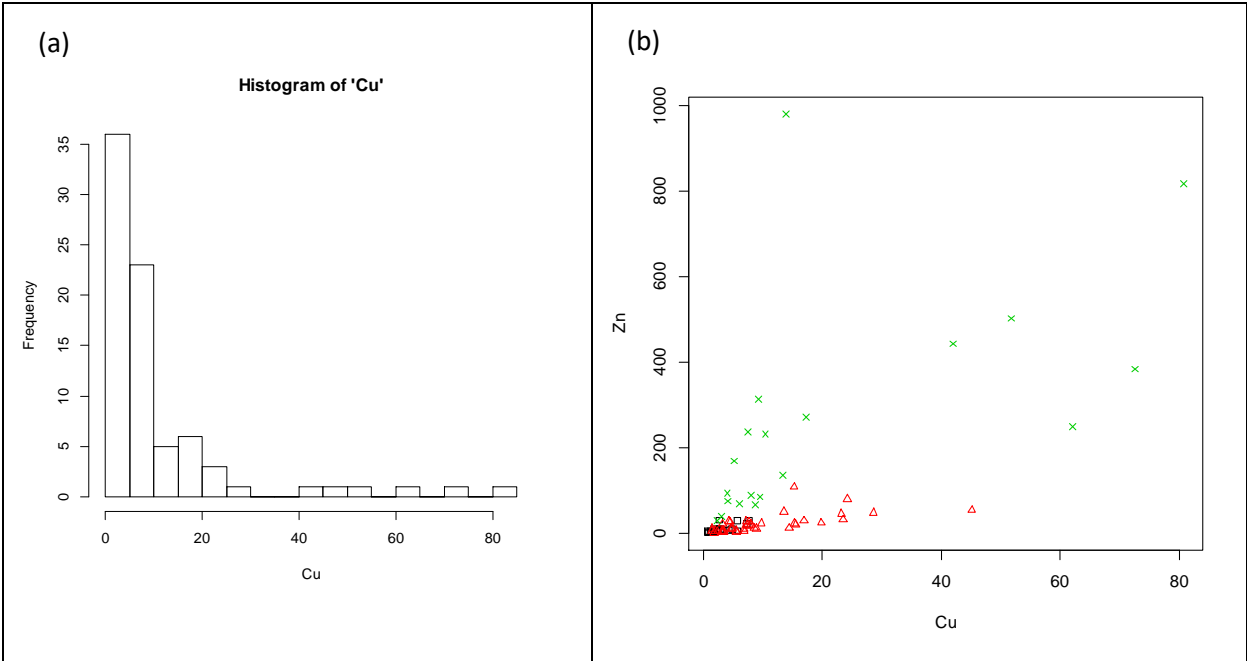
10  
11 For example a key question in Archaeology is provenance. Are we able to identify which  
12 artifacts have common provenance and which differentiate from others? Is it possible to  
13 determine the factors which differentiate the distinct groups and learn from this process either  
14 with respect to the characterization of a source, or the technology used towards ceramic  
15 production? If provenance is the question of the analysis, data carrying information about the  
16 characteristics of an origin are necessary, such that the chemical composition. This will lead to a  
17 number of continuous variables in statistical terms, i.e. quantitative variables which can be  
18 measured on a scale and take infinite number of values. Those continuous variables will be the  
19 input data for a cluster analysis for example. Apart from chemical composition, other types of  
20 variables may be available, such as data collected from typology or microscopic study of the  
21 specimens. These variables will be categorical, nominal or ordinal. For example, variables taking  
22 values 'yes' or 'no' to the question if a certain mineral is present or not to a specimen, or variables  
23 taking values 1, 2, 3 and 4 where level 1 corresponds to the 'no presence' outcome, 2 to the  
24 outcome 'few', 3 for 'moderate' and 4 for 'plenty'. If we assume that this type of data also carry  
25 information about the provenance, methods of clustering which accommodate categorical  
26 variables would be appropriate to pursue or methods solely dedicated for categorical data.  
27

28  
29 The first stage of the analysis is to have the questions clearly stated. Moving to the numerical  
30 variables the preliminary analysis can include univariate summary statistics for each of the  
31 variables and univariate or bivariate plots for pairs. The preliminary analysis in general can give  
32 an idea of the data at hand. For example, range of values for a variable, shape of the distribution,  
33 existence of specimens with extreme values and possibly an obvious pattern, e.g a bi-modal  
34 distribution for one or more variables may suggest to the existence of two distinct groups. The  
35 preliminary analysis will not be sufficient to answer questions of the research project, but is  
36 essential and has an auxiliary role in the proceed analysis.  
37

38  
39 More analytically, for continuous variables a preliminary analysis can include summary statistics,  
40 i.e. mean, variance, standard deviation, median, first and second quantile, minimum and  
41 maximum value. Graphs, such as histograms for univariate variables will give evidence for  
42 symmetry or not, existence of a long tail due to extreme high values on this variable for instance,  
43 bi-modality etc. Moreover, it will point out a specimen with extreme value on that particular  
44 variable, a fact that needs to explore further. For example, if a statistical method which assumes  
45 normality is going to be used, a transformation in the log-scale will be suggested before the  
46 analysis in the case of a long right tail at the variable distribution.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 Scatter plots for pairs of variables will also provide similar information, but moreover will  
5 add information about the relation between the two variables. Ellipsoidal shape of a scatter plot  
6 instead of spherical will indicate significant level of correlation. An obvious pattern in the data  
7 with respect to distinct groups may also be apparent. However, as it is the case for any simple or  
8 complex statistical analysis with multivariate data, the difficulty is the number of variables  
9 available, usually they are too many, and the contradict information that they may contain.  
10 Moreover, by studying the variables one by one the part of information due to the correlations  
11 among the variables is left out. A bivariate analysis is definitely superior that a univariate study,  
12 but the number of possible pairs is even higher and a selection of informative pair or pairs is  
13 essential.  
14

15  
16 As an example, we use simulated data for giving measurements on five chemical elements.  
17 The data set of total size 60 is comprised from three groups, where each group is generated using  
18 the lognormal distribution and parameters, i.e. means and covariance matrices differ among the  
19 three groups. Figure 1(a) plots a histogram of one of the element measurement. A long right tail  
20 is apparent and log transformation could transform the data producing a more symmetric plot  
21 allowing a symmetric model (e.g. normal distribution) to be assumed for the statistical  
22 methodology. Figure 1(b) is a scatter plot in two dimensions, plotting the measurements of two  
23 elements and using different symbols and colors for the three groups. From figures 1(a) and 1(b)  
24 we draw immediate results that the data presents heterogeneity, in particular high values in  
25 some of the variables and different type of correlation for at least two subsets of the data as  
26 indicated from the two selected variables in Fig. 1(b). However, this preliminary analysis does not  
27 form an analysis or enough evidence to draw any conclusions on certain structure regarding  
28 distinct groups in data.  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



**Figure 1.** (a) A Histogram of univariate measurements of one selected element. (b) A scatter plot of bivariate data set. Color and symbol shape correspond to the groups.

### Data Transformation

Data transformation plays an important role in subsequent analysis. If provenance is the problem of interest, a method of multivariate statistics within the framework of clustering will be applied and any result of cluster analysis will be visualized through principal component analysis (PCA) or any other statistical multivariate technique aiming to data reduction. Performance of either clustering or data reduction techniques will depend on data transformation. For example, for PCA in particular, it is a fundamental result that if the data variable variances differ, the first components of a PCA analysis will be dominated from the variable(s) with large variance. This is against the PCA performance since the first components instead of explaining the majority of data heterogeneity will solely explain the variability that data present at one or two variables, the variable with dominant variances.

For implementing cluster analysis, as seen in Section 3 analytically, both factors of data transformation and measure of distance will be essential for the analysis. Transformation of data primarily holds for continuous data in Archaeology, i.e. data resulting from the chemical composition of specimens. This analysis will provide us with a large number of variables, each one corresponding to one chemical element or their oxides. The measurements are either counts i.e. absolute measures or percentages out of 100% of weight. These data where the measurements for the set of variables sum to 100% for each sample are called fully compositional. Compositional data in general are called either the fully compositional or data that can be considered as a subset (with respect to variables) of fully compositional data.

1  
2  
3  
4 Transformations usually implemented in analysis of Archaeometric data are (i) standardization  
5 or scale transformation in zero mean and variance one, (ii) log-transformation (iii) log and scale,  
6 (iv) log-ratio transformation for compositional data.

7  
8 In particular, standardization will be suggested, among other cases, for a PCA visualization as  
9 mentioned before, log-transformation will be appropriate and improve the performance of  
10 either a data reduction/visualization technique or a clustering method if the data present  
11 asymmetry. It is proposed from the literature (see for example Bieber et al., 1976 and Bishop and  
12 Neff, 1989) logarithms to be taken with base 10. This transformation will result to a set of variables  
13 that will have nearly equal variances (a characteristic that standardization also guarantees) and  
14 variables measured in percentages or ppm are transformed to measurements which weight  
15 almost the same. Towards the implementation of log transformation special care is required for  
16 measurements that are zero.

17  
18 Zero measurement naturally results in compositional data when the measured value is below  
19 the detection threshold. A possible treatment in this case would be to substitute zero with the  
20 threshold or with a value  $\alpha$  smaller than the threshold (Beardah et al, 2003) or substitute zero  
21 with  $\alpha/(p - 1)$  with  $\alpha$  a small number and  $p$  the number of variables (Aitchison, 1986). Other  
22 proposal is to adjust the remaining variable measurements when a zero is replaced by  $\alpha$ . More  
23 specifically if  $x_{ij}$  is the measurement of  $i$  sample for the  $j$  variable,  $x_{ij}$  is replaced with  $(x_{ij} -$   
24  $\alpha x_{ij})/100$  (see Pawlowsky-Glahn, 2002). Another approach is to impute those values based on  
25 the remaining measurements (see Palarea-Albaladejo et al. 2015). Most of the  
26 statistical/computational packages include imputation techniques.

27  
28 The log-ratio transformation is aiming to take into account the condition that compositional  
29 data have that they add up to a constant number. It is proposed by Aitchison (1986) and adopted  
30 by Buxeda I Garrigos (1999), Martin – Fernandez et al. (2015) and Pawlowsky-Glahn and Buccianti  
31 (2011) amongst others. If  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  is the vector of measurements of artefact  $i$  for  
32 the set of  $p$  variables, the log-ratio transformation of  $x_{ij}$ , ( $j = 1, 2, \dots, p$ ) is  $\log(x_{ij}/g(\mathbf{x}_i))$  where  
33  $g(\mathbf{x}_i) = (x_{i1} x_{i2} \dots x_{ip})^{1/p}$  (centered log-ratio transformation). The merit of the log-ratio  
34 transformation is that the condition which holds for the data have is taken into account.  
35 Transforming the data in log-ratio scale such way the distances among vectors are expressed in  
36 a geometry that represents better their relation. There is however a discussion in the literature  
37 for choosing log-ratio and simply standardize the data for compositional data. When using log-  
38 ratio transformation an element with small absolute values, but relatively high variance is  
39 promoted in comparison to other elements with higher absolute presence and more relative to  
40 the question of the analysis, e.g. provenance. See Baxter (2001) and Baxter et al (2006) for some  
41 examples.

42  
43 Nevertheless, if a structure is apparent in the data, both transformation will suggest this  
44 same structure and differences caused by practical problems as mentioned or presence of  
45 outliers will need further examination. It is important at this stage to mention that a PCA plot, in  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 either transformation will not be a method to conclude about a certain structure in the data,  
5 even if this it is apparent. A clustering technique is the appropriate method to answer this  
6 question and PCA will only act complimentary to a clustering technique confirming the validity of  
7 a clustering result. Other methods for assessing a clustering technique are also available and a  
8 common practice is to reckon in all them. Moreover, a two dimensional PCA plot, based on the  
9 first two principal components will not count for the 100% of the data heterogeneity and  
10 therefore it can serve as a data visualization method, as the closest view of the data we can have  
11 in two dimensions, but yet it is not the complete information. Some theoretical and practical  
12 considerations on PCA are presented in section 3.  
13  
14  
15  
16  
17

18 Another practice would be to take ratios of variables instead of raw measurements in  
19 variables as an attempt to explain dilution and alleviate its effect. This dilution may be the result  
20 of different proportion of temper added to the paste towards the production of ceramic. The  
21 paste source may be the same, but different proportion of temper may result in different  
22 composition. One way to cancel out this effect if to work with ratios, i.e. instead of measurements  
23  $(x_{i1}, x_{i2}, \dots, x_{ip})$  for the  $i$  artefact in elements  $p$  elements, the ratios  $(x_{i1}/x_{ik}, x_{i2}/x_{ik}, \dots, x_{ip}/$   
24  $x_{ik})$  can be used for the analysis, where  $x_{ik}$  is the measurement at a chosen element  $k$ . For a  
25 detailed discussion see Baxter (2001).  
26  
27  
28  
29  
30  
31  
32

### 33 **3. Statistical methods for provenance**

34 Cluster Analysis (CA) is the most widely used method of multivariate statistical analysis in  
35 Archaeology (Baxter, 2008). It is a term to include any statistical method, now-days any machine  
36 learning as well, seeking for similarities among observations, based on a number of variables, and  
37 identify groups consisted of observations that have common characteristics. Those groups are  
38 called clusters. CA is applied to Archaeometric data in Archaeology aiming to identify artefacts  
39 that share common characteristics in composition and therefore make inference about these  
40 with respect to provenance, technology and draw subsequent conclusions on economic and  
41 social relation of past societies.  
42  
43  
44  
45

46 A plethora of clustering techniques is available and this is in practice one of the reasons the  
47 task it can be challenging. CA is typically used when after a preliminary analysis of the data, we  
48 suspect that there is heterogeneity within the dataset and it cannot be assumed as a sample  
49 coming from a unique population. This can be apparent either from univariate or bivariate plots  
50 (Fig. 1), summary statistics or plotting the data using any data reduction technique (e.g. PCA) and  
51 detect some pattern in the data, e.g. bi-modality or different type of relation among observations  
52 (e.g. Fig 1b). Another charactering of CA that make the task challenging, is that no prior  
53 knowledge is of group membership is assumed. Any CA method is classified within the  
54 unsupervised statistical learning techniques. With 'unsupervised' we refer to the fact that there  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 is no prior information available to either assist towards the group formation or the  
5 assessment/comparison of clustering techniques.  
6

7  
8 Moreover, the range of statistical methods available for CA varies with respect to the  
9 statistical assumptions they make, the complexity of the methodology and the computational  
10 demand. Most importantly, they may also vary in special characteristics, such as the effect of  
11 presence of outliers in clustering procedure and the properties they impose in resulting groups,  
12 e.g. some method impose spherical shape into clusters or they assume equal size groups.  
13 Therefore, for a scientist of Archaeology, or any other area, some knowledge of the theory behind  
14 each technique is recommended and in any case an exhaustive exploration of clustering  
15 techniques would be necessary before drawing any conclusions. Some literature in CA among  
16 others, are: Everitt et al (2011), Baxter (2015) and Papageorgiou (2018). For a practical guide to  
17 implementation CA see Kassambara (2017).  
18  
19  
20  
21

22 The statistical clustering techniques are classified into four categories with respect to their  
23 approach towards the problem. Namely, the categories are: Hierarchical, optimization or  
24 partitioning, model based and density based. From those four categories hierarchical and  
25 partitioning methods are most often used in Archaeology mainly because of the ease in  
26 implementation and lack of statistical hypothesis they assume. Most of the statistical  
27 methodologies may handle both continuous and categorical data, but it is more common in  
28 Archaeology to use the chemical composition (continuous) data only for obtaining the groups  
29 and verify or compare those with the information available from categorical data, such as  
30 mineralogical data. The main reason why this is the case is that mineralogical data are less  
31 frequently recorded in a manner that invites quantitative analysis and they are considered semi-  
32 quantitative. If quantitative discrete data are available, an analysis to the combined data, mixed-  
33 mode as they called, is appropriate. A discussion on mixed-mode data analysis follow the  
34 presentation of the clustering techniques and methodologies with the relative references will be  
35 provided at this stage.  
36  
37  
38  
39  
40  
41  
42  
43  
44

### 45 **Hierarchical Clustering**

46 Hierarchical clustering is an approach of clustering that is algorithmic and is based on distances  
47 or equivalently similarities that can be calculated among observations based on the set of  
48 variables that contribute to the analysis. There are two 'symmetric' ways to implement  
49 hierarchical clustering. The agglomerative, according which the algorithm initiates from the  
50 situation that each of the  $n$  observations form a separate cluster and algorithm proceeds by  
51 merging observations until all observations are located in one single group and the divisible  
52 hierarchical clustering where the starting scenario is the inverse, i.e. all  $n$  observations form a  
53 single group and the algorithm proceeds by repetitive partitions until all observations are  
54 separated. In both cases, the whole procedure, from the one extreme situation until the other,  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



is completed and we decide about the number of clusters and the observations' cluster membership by inspection.

Two factors determine the hierarchical method. (a) The distance measure and (b) the linkage method. Distance or dissimilarity measure for continuous variables, can be any distance metric. The distances between all possible pairs of observations are calculated using the same chosen metric. Assuming that the data matrix consists of measurements of  $n$  observations on  $p$  chemical elements, an  $n \times p$  matrix, calculating the distances among all possible pairs will lead to an  $n \times n$  symmetric matrix called dissimilarity matrix, usually denoted as  $D$ . Table 1 lists the most widely used distance metrics as a measure of dissimilarity for continuous variables. At each stage of the algorithm an updated dissimilarity matrix is calculated.

<i>Measure</i>	<i>Definition</i>
Euclidean	$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$
Squared Euclidean	$d_{ij} = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2$
City block or Manhattan	$d_{ij} =  x_{i1} - x_{j1}  +  (x_{i2} - x_{j2})  + \dots +  x_{ip} - x_{jp} $
Minkowski	$d_{ij} = \sqrt[m]{(x_{i1} - x_{j1})^m + (x_{i2} - x_{j2})^m + \dots + (x_{ip} - x_{jp})^m}$
Maximum or Chebyshev	$d_{ij} = \max_k  x_{ik} - x_{jk} $
Pearson	$d_{ij} = (1 - \varphi_{ij})/2$ , where $\varphi_{ij}$ is the Pearson correlation between data vectors $\mathbf{x}_i$ and $\mathbf{x}_j$ .

**Table 1.** Distance measures for continuous variables.

Updating will be necessary as merging for example in agglomerative clustering will occur and the number of clusters will change. Linkage is the method that the algorithm will use in calculations of distances among clusters or clusters and observations. Such calculations will arise at intermediate stages of the algorithm. If we further denote by  $x_{im}$  the  $(i, m)$  element of the data matrix, ( $i = 1, 2, \dots, n$  and  $m = 1, 2, \dots, p$ ), i.e. the measurement of the  $i$ -th observation on the  $m$ -th chemical element, and  $d_{ij}$  the elements of matrix  $D$ , Table 2 lists the definition of the most frequently used in practice linkages. It is worth mentioning that all these definitions are relatively simple, therefore not time consuming, and they can be implemented in most of the statistical packages. At each step of the agglomerative algorithm, units or clusters that correspond to the smallest distance are merged. The adopted choices of both metric and linkage remain the same for all intermediate steps of the algorithm.

Linkage	Definition of inter-group distance of groups A and B
Single Linkage or Nearest Neighbor	$d(A, B) = \min_{\substack{i \in A \\ j \in B}} d(i, j)$
Complete Linkage or Furthest Neighbor	$d(A, B) = \max_{\substack{i \in A \\ j \in B}} d(i, j)$
Average Linkage	$d(A, B) = \bar{d}, \text{ where } \bar{d} \text{ the average distance}$ $\bar{d} = \frac{1}{n_A n_B} \sum_{i \in A} \sum_{j \in B} d(i, j)$ <p><math>n_A, n_B</math> the group sizes of A and B respectively.</p>
Centroid	$d(A, B) = d(\bar{x}(A) - \bar{x}(B))$ <p>where <math>\bar{x}(A)</math> is the point</p> $\bar{x}(A) = \frac{1}{n_A} (x_1(A), x_2(A), \dots, x_p(A))$ <p>and called centroid of group A.</p>
Ward's	<p>The increase in total error sum of squares (ESS) the merge of A and B will cause.</p> <p>If <math>d_e</math> is the Euclidean distance, ESS of a group U with centroid <math>\bar{x}(U)</math> is defined as</p> $ESS(U) = \sum_{i \in U} d_e(x_i - \bar{x}(U))$

**Table 2.** Agglomerative linkages.

Linkages listed in Table 2, have some properties it is useful to know, especially when assessing a clustering result. Single linkage is prone to the chaining phenomenon, a negative property where distant groups may be merged due to the existence of two neighbour measurements. . On the other hand it useful in identifying outliers. Complete, average and ward's methods are used more often in practice and they produce compact clusters. However, they all have problems in practice when ellipsoidal clusters occur in the data -a case quite common in Archaeometry.

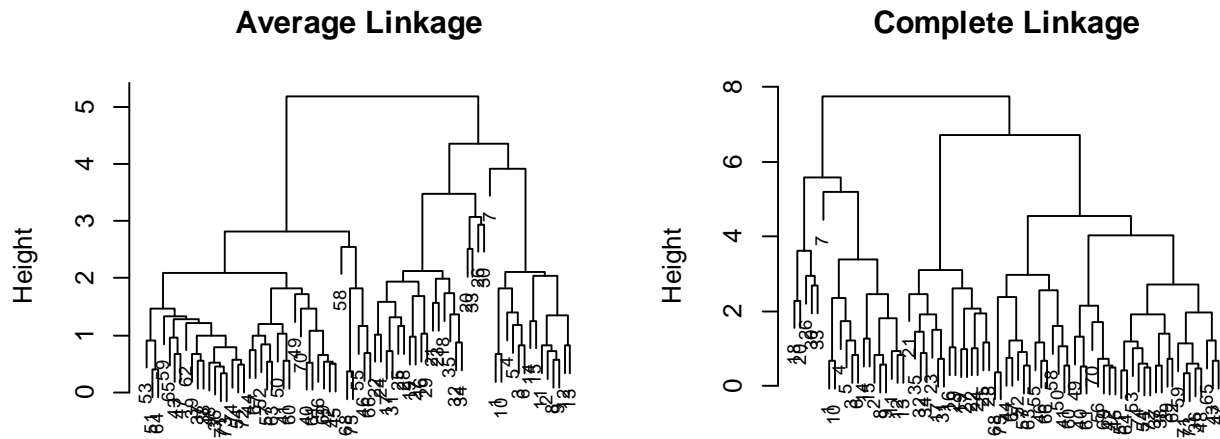
Hierarchical clustering has as a property that once a merging of two observations or groups occurred during an iteration this will remain until the end of the algorithm. From the practical point of view, towards the implementation of hierarchical clustering, choices on: data transformation as discussed in section 2, distance metric and linkage method have to be made. It is apparent that having more than one option for each one of the three factors the number of possible combinations is quite large.

The sequence of hierarchy in merging or partitioning of a hierarchical clustering is represented graphically into a plot called dendrogram. The existing clustering, if any, will result by 'cutting' the dendrogram at a certain height. We illustrate the procedure using some simulated data.

**Example 1.** This first example is consisted from a 95x9 data matrix. The data have been generated by assuming that there are three groups each one distributed by a lognormal distribution with

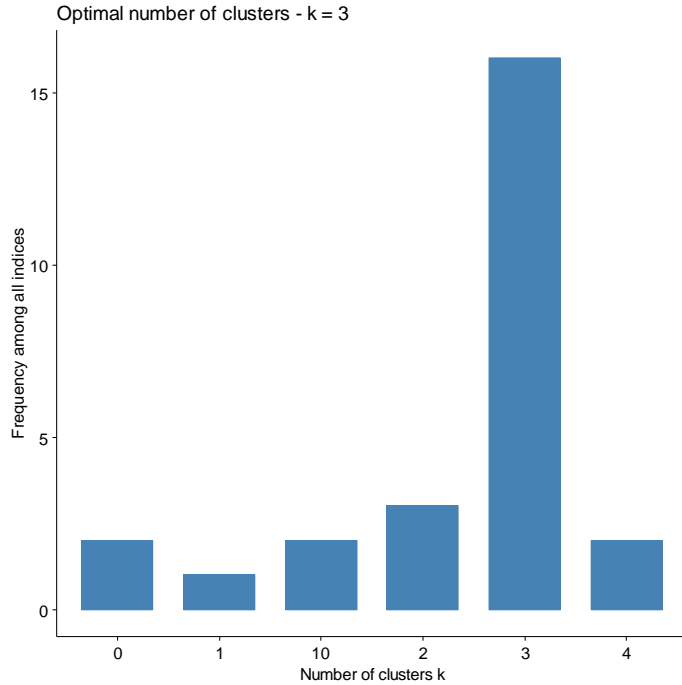
1  
2  
3  
4 different sample sizes and parameters. The sample sizes for the three groups are 15, 20 and 40  
5 respectively.  
6  
7  
8

9 Using standardization as a transformation for the data Figure 2 presents the dendrograms  
10 for average and complete linkages. Different software packages may use different ways of  
11 presenting the dendrogram, but the idea is the same. The sequence of the hierarchical clustering  
12 algorithm is captured in one plot and based on that we would like to answer the question ‘how  
13 many clusters are in the data’? Or in other words, at which height we need to ‘cut’ the tree? The  
14 suggested height, at which we cut the tree is where the compact branches of the dendrogram  
15 will remain intact. For example, for the dendrogram according to complete linkage in Figure 2,  
16 cutting the tree at height 6 would be sensible, as the three compact branches of the tree remain  
17 intact. Deciding the number of clusters or the height one cuts the tree is not always  
18 straightforward. There might be two different heights, equally possible, or none. Since the  
19 problem is unsupervised and true clustering is not available to compare with and conclude, the  
20 suggestion is to try all possible scenarios and assess each one of them.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31



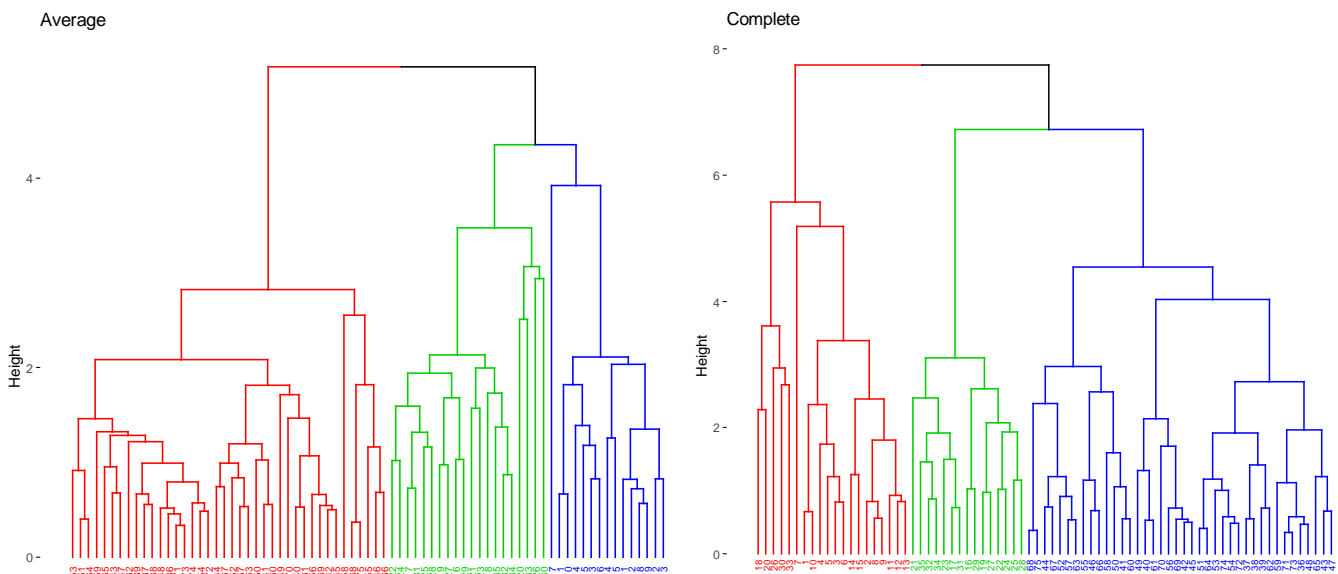
47 **Figure 2.** Dendrograms for the standardized simulated data using average and complete linkage

48  
49  
50 Apart from the dendrogram which is available from a hierarchical clustering, a number of  
51 indices are available in the literature dedicated to suggest the optimal number of groups (Charrad  
52 et al, 2014). Figure 3 plots the bar plot of 30 available indices implemented to the simulated data  
53 under study under the Average Linkage. The number of groups equal to three is the dominant  
54 choice.  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



**Figure 3.** Optimal number of clusters for the simulated data under Average Linkage.

Taking into account both dendrograms and indices, let's assume that we decide to cut both trees at the height where the three branches would be suggested, e.g. height 4 for the average linkage and height 6 for the complete. Figure 4 plots the resulting clusters coloured using different colour for each cluster.



**Figure 4.** Average and complete linkage dendrograms with the proposed clusters given in different colors.

### K-means

1  
2  
3  
4 Another category of clustering techniques is partitioning or optimization method (see Everitt et  
5 al, 2011). According to this approach the groups are formed applying a partition to the data into  
6 a certain number of groups, using an optimization criterion. There is no hierarchy associated to  
7 the group membership, as in hierarchical clustering, and number of groups in the data has to be  
8 known in advance. The most known method in this category is k-means which uses as  
9 optimization criterion the minimum sum squared error. K-means as a method is very popular in  
10 Archaeology.  
11

12  
13  
14  
15 Baxter (2015) states that this popularity is mainly because it is readily understood, is  
16 perceived as being geared to archaeological needs, and was rendered accessible at a time when  
17 computational resources were limited compared to what is now available. Baxter in the same  
18 paper proposes variations of k-means method. This is because it is well established in the  
19 literature that k-means tends to produce spherical shape clusters and k-means is more  
20 appropriate for equally sized and spherical shape clusters (see Baxter 2015, Banfield and Raftery  
21 1993, Papageorgiou et al 2001 among others). This is not the case in Archaeometric data,  
22 especially for geochemical data collected for provenance studies, where ellipsoidal shape and  
23 unequal size clusters exist. Another disadvantage of k-means is the effect of the presence of  
24 outliers.  
25  
26  
27  
28

29  
30 Therefore, although popular, k-means may be used as one method, but needs to be cross  
31 examined with other techniques and check the validity of the proposed clustering. Within the  
32 framework of partitioning methods, k-medoids and Trimmed k-means are possible alternatives  
33 with improved behavior for both disadvantages of k-means method (Steinley, 2006).  
34  
35  
36

### 37 **Model-Based methods**

38 Model-based clustering assumes a model, i.e. statistical distribution; to describe each distinct  
39 group in the data. The data matrix is assumed to be generated from a mixture of distributions  
40 and each component of the mixture captures one group. Fitting the mixture model to the data  
41 and estimating the parameters for every component will allow the scientist to identify the  
42 clusters and assign the observations to one of the resulting clusters based on the maximum  
43 posterior probability (see Banfield and Raftery, 1993). The estimation of parameters is  
44 implemented by Expectation Maximization (EM) algorithm, a computational statistical technique  
45 for deriving the maximum likelihood estimator. As a consequence, model-based clustering  
46 method is more demanding with respect to computational effort and it can be implemented via  
47 a statistical package. The computational demand and technical problems will be more intense as  
48 the dimension of the data observation is increase. The use of specialized statistical package to  
49 implement model-based cluster may be the reason for a limited use of model-based clustering in  
50 Archaeological applications.  
51  
52  
53  
54  
55  
56

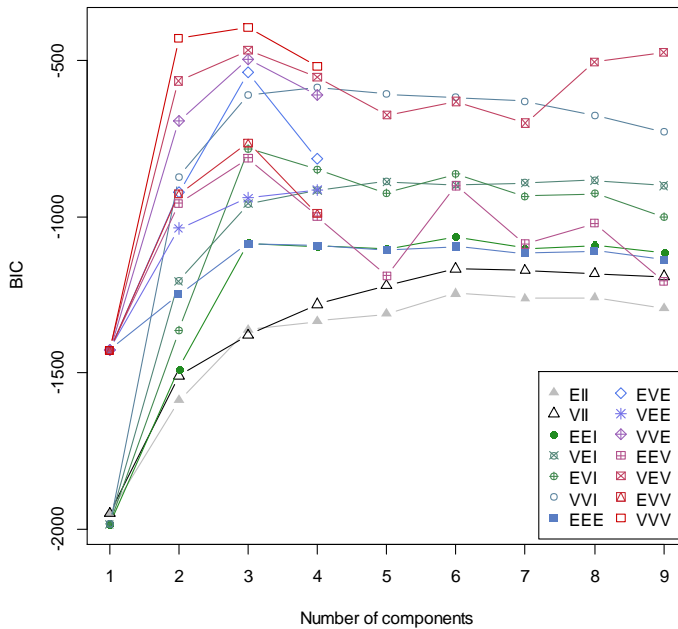
57  
58 The advantages of model-based clustering are many. The probabilistic model that is assumed  
59 provides the framework to assume ellipsoidal, different size or different orientation clusters. It is  
60  
61  
62

therefore more general and more appropriate for compositional data. Moreover, estimating the model parameters results to a fully defined probabilistic model and statistical inference is now feasible. Statistical tests for comparison the model fitting are available and can be used in order to decide the number of components, i.e. the number of clusters in the dataset.

Other clustering methods, for example density-based methods, fuzzy clustering, kernel clustering are also available in the literature. Moreover, recent methods especially designed to cope with the presence of outliers, high dimensionality and large-scale data have also been proposed (Xu and Wunsch II, 2008).

**Example 1 (continued). K-means:** Let us assume that as suggested from dendrograms in Figure 2 and the majority of indices in Figure 3 the number of clusters is three. Implementing k-means with  $k=3$  and using the same transformation for the data results into three quite different clusters in comparison with those in hierarchical and the true known origin since the data are simulated. More analytically, the first two clusters of 15 and 20 observations are merged to one cluster, and the third cluster of 40 observations is divided into two smaller.

**Model-based:** Implementing model-based clustering without any prior information about the number of clusters in the data the optimal model based to BIC criterion was the model with three components and variable volume, shape and orientation. Figure 5 plots the BIC value for a range of number of components (one to nine) and a range of possible models. All seventy five observations are correctly assigned, i.e in complete agreement with the true origin.



**Figure 5.** BIC criterion for model selection in model-based clustering.

#### 4. Compare clustering results and assess their validity

The clustering methods and possible modifications of each method, together with data transformation choices, lead to a great number of possible clustering. Clustering in general, is an unsupervised technique, meaning that the true classification is not known in advance and there is not available a labelling for the observations indicating the cluster they belong to. Lacking of this information, assessment of the plethora of possible clustering results that may be available for the same case study may be a very difficult task. The difficulty is to decide which clustering is more valid compared to others or which clustering is more valid in general and therefore more probable to represent the true classification. In this paragraph we focus on how we compare clustering results obtained from different methods or choices in implementation and secondly how we assess a specific clustering with respect to its validity.

##### Correlation and Similarity Indices for clustering comparison.

One way to compare two or more clustering with respect to their similarity is to calculate the correlation of the corresponding dendrograms. Cophenetic and Baker correlation coefficients are appropriate measures (Saraçlı, S. et al. 2013). Adjusted Rand index (Hubert and Arabie, 1985) is a popular measure of clustering agreement and other indices used for the same purpose are: Rand, Fowlkes and Mallows, Wallace and Jaccard index (Gordon, 1998 and Shotwell, 2013). Adjusted Rand index is taking values from 0 to 1 and the closer to one is the better agreement between the two clustering results. Such measures are useful, especially when a large number of clustering techniques is considered because if the majority of some clusterings agree, then it is more probable this particular clustering to hold in reality.

**Example 1 (continued).** Let us consider the two clustering results as presented in Figure 4, clustering from k-means with  $k=3$  and clustering as suggested from the optimal model on model-based clustering. Comparison between the two hierarchical methods produces 0.90 and 0.77 Cophenetic and Baker correlation coefficients respectively. These results suggest that the two dendrograms are similar, but not identical. This comparison can be extended to as many dendrograms are considered. Table 3 lists the adjusted Rand and Fowlkes and Mallows indices for the clusterings under consideration. Based on these coefficients we can conclude that hierarchical average and model-based clustering give almost identical result, complete is close, but not identical to the first two, whereas k-means clustering has a relatively low degree of similarity with all other clustering results. In a real analysis situation, the list of clustering results under consideration is expected to be much longer than four. If one particular clustering is systematically apart from the majority, it would suggest that this clustering is weak. If on the contrary, a clustering is confirmed with its comparison with many others, in the sense of agreement, this clustering is a strong candidate.

	Adjusted Rand				Fowlkes and Mallows			
	Average	Complete	k-means	Model-based	Average	Complete	k-means	Model-based
Average	1				1			
Complete	0.89	1			0.93	1		
k-means	0.46	0.46	1		0.66	0.66	1	
model-based	0.97	0.87	0.47	1	0.98	0.92	0.66	1

**Table 3.** Adjusted Rand and Fowlkes-Mallows indices calculated for the four clustering results on simulated data.

After calculating as many as possible correlation/similarity indices for the derived clusterings, we proceed with examining the source of disagreement. For the numerical example, if we calculate the classification table between the four examined clusterings using labels 1, 2 and 3 to denote the three clusters, Table 4 lists the results. The degree of similarity between Average and Complete hierarchical clusterings is explained because they only differ at the cluster membership of five observations. Average linkage classifies those in cluster 2, while complete linkage in cluster 1. The data here are simulated and we know that three groups of sizes 15, 20 and 40 have been generated. It is easy to confirm here that average linkage agrees with the grouping as simulated. Model-based which scores high in similarity with hierarchical clusterings, differs with the Average linkage in classification of one out of seventy five observations. For k-means, the two smaller groups 1 and 2 are merged to one, and the third group of forty observations is split into two others. This is the more distant to the true clustering among the four proposals.

		Complete			k-means			Model-based		
		1	2	3	1	2	3	1	2	3
Average	1	15	0	0	0	15	0	14	1	0
	2	5	15	0	0	20	0	0	20	0
	3	0	0	40	21	0	19	0	0	40

**Table 4.** Cluster membership comparison for average and complete linkage dendrograms.

### Internal and External cluster validation.

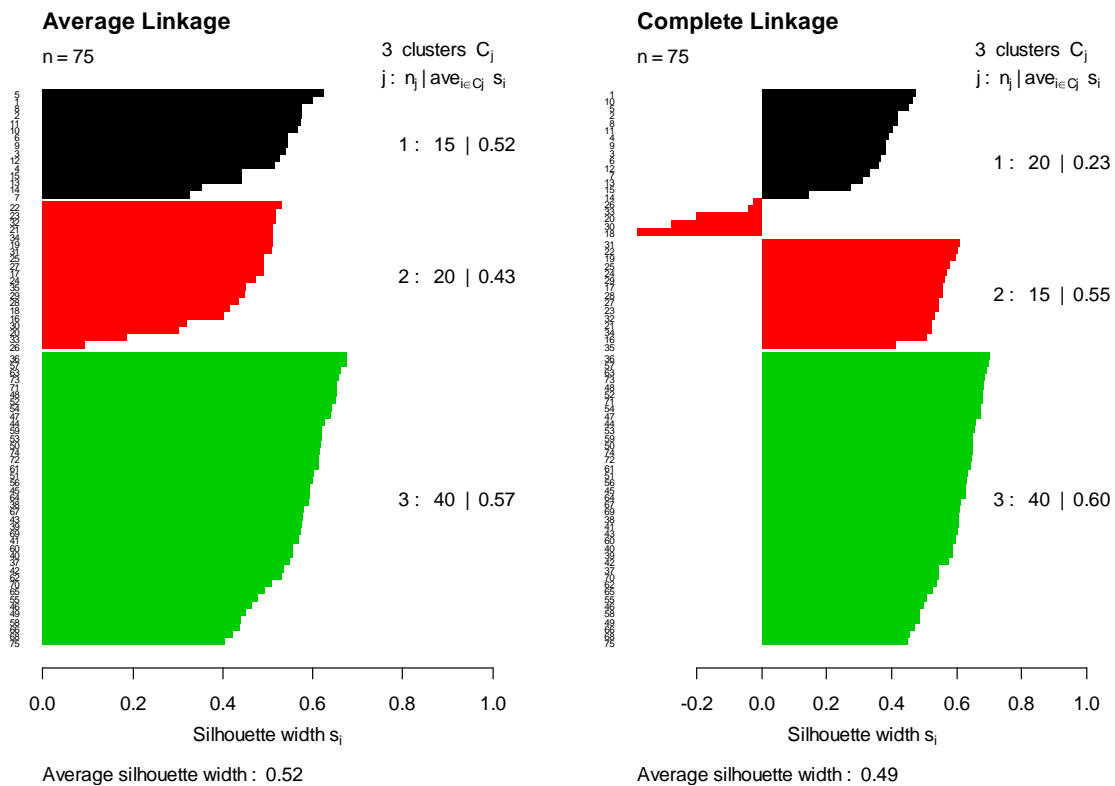
Apart from comparing the clustering results to conclude on which classification is more dominant, a validation of a specific clustering would be valuable to assist in deciding which analysis is the analysis that produce trustworthy results, closer to reality. Clustering validation or assessment is not easy nor has a unique solution. Two approaches are proposed here as dominant for the problem. The first approach is based on specialized indices or measures that



measure the quality of clustering. Secondly, a graphical representation of a clustering result under consideration against the data may be used.

Silhouette coefficient (Rousseeuw, 1987) is the most popular measure to assess the internal validity of a clustering. Silhouette coefficient is a measure which shows how well the objects lie within the clusters they are assigned to. The silhouette coefficient takes values from minus one to one and a value close to one corresponds to a well clustered observation whereas a negative value indicates the worst situation with respect to clustering.

**Example 1 (continued).** Figure 6, plots the Silhouette coefficients for all the observations of the simulated dataset at Example 1, according to the clustering results presented in Figure 3. The average score of Silhouette coefficient is higher for the Average Linkage (0.52) in comparison with Complete (0.49) indicating a relatively better clustering. Moreover the Silhouette coefficient for five observations according to complete linkage is negative, i.e. these observations are not well clustered. It is easy to verify that these are the same five observations that the two clustering results assign in different groups as listed in Table 4. The Silhouette coefficient for k-means method is 0.31, a lower value than 0.52 for the Average linkage, a fact that advocate for the inferior performance of k-means in this certain example.



**Figure 6.** Silhouette coefficients for a three cluster result according to average and complete linkage.

Lambda-Wilk's test is another statistic that in this context can provide a measure of compactness within each proposed cluster and how well clusters are separated from each other.

1  
2  
3  
4 The smallest value for Lambda-Wilk's test the better separation among the groups. For the data  
5 at Example 1, Lambda-Wilk's test for all considered methods are very small, within the range of  
6 0.02 to 0.05.  
7  
8  
9

10 The second approach towards the assessment of a clustering result is a graphical method. A  
11 clustering proposal is visualized by plotting the data using the labeling provided from the  
12 clustering. The plot can either give grounds to accept a clustering when this gives a sensible result  
13 when plotted, or on the contrary, it can be used to explain why a clustering result performs  
14 poorly. In order to plot a clustering result on data the need of a lower dimension coordinates  
15 system is required. For this reason, a technique of data dimension reduction is necessary and  
16 Principal Component Analysis is the most frequently used in Archaeometry.  
17  
18  
19  
20  
21

### 22 **Principal Component Analysis**

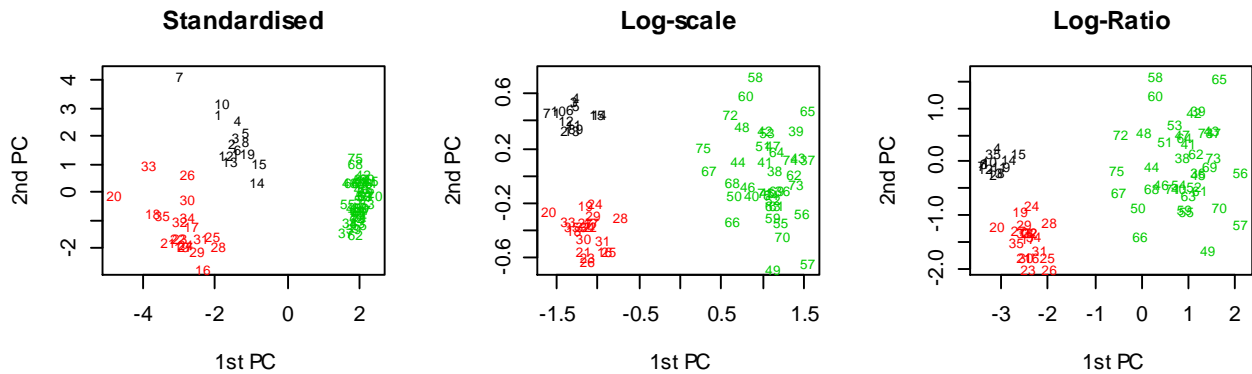
23 Principal Component Analysis (PCA) (see for example Jolliffe, 2002 as a general reference and  
24 Rogers et al. 2016 for a practical guide) is a data transformation technique which exploits the  
25 correlation among the data variables, in order to construct a new set of variables, the Principal  
26 Components (PC) with some good properties. The first property is that the two data sets are  
27 equivalent with respect to the amount of information included as this measured from the total  
28 variability of the data points.  
29  
30  
31

32 The advantage of PCA is that the variability reconstruction is such that there are no  
33 correlations among the principal components, i.e. the new set of variables is uncorrelated which  
34 means that a univariate study of all variables will reveal all included information. Moreover, there  
35 is a hierarchy in Principal Components importance. The leading ones explain the majority of the  
36 total variability allowing us to retain only a subset of PCs and explain a significant percentage of  
37 the total variability based only on those PCs. Especially in applications, such as Archaeometry,  
38 that data variables are highly or moderately correlated, two or three of the leading PCs can  
39 explain a sufficient amount of the total data variability. Another characteristic of PCs is that they  
40 are linear combinations of the original variables. This very simple form of transformation enables  
41 the scientist to interpret the PCs with respect to the original variables. This is based on the  
42 loadings, as they called the coefficients of the linear combinations of PCs.  
43  
44  
45  
46  
47  
48

49 Before applying PCA in archaeometric data, transformations as discussed in Section 2 will be  
50 necessary. If raw data are used, and not logged or log-ratio transformed the standardization of  
51 the data in zero mean and variance one will be needed. This is because if the variances among  
52 the data variable vary a lot, the leading PCs will be dominated from those with large variance. In  
53 this case the leading PCs will not explain the majority of the data set variability across all variables,  
54 but will explain the variability of the data with respect to the variables with large variance only.  
55 Especially in compositional data where the measurements for some oxides account around 50%  
56 of the total composition and for some trace elements the measurements are particles in millions  
57  
58  
59  
60  
61  
62  
63  
64  
65

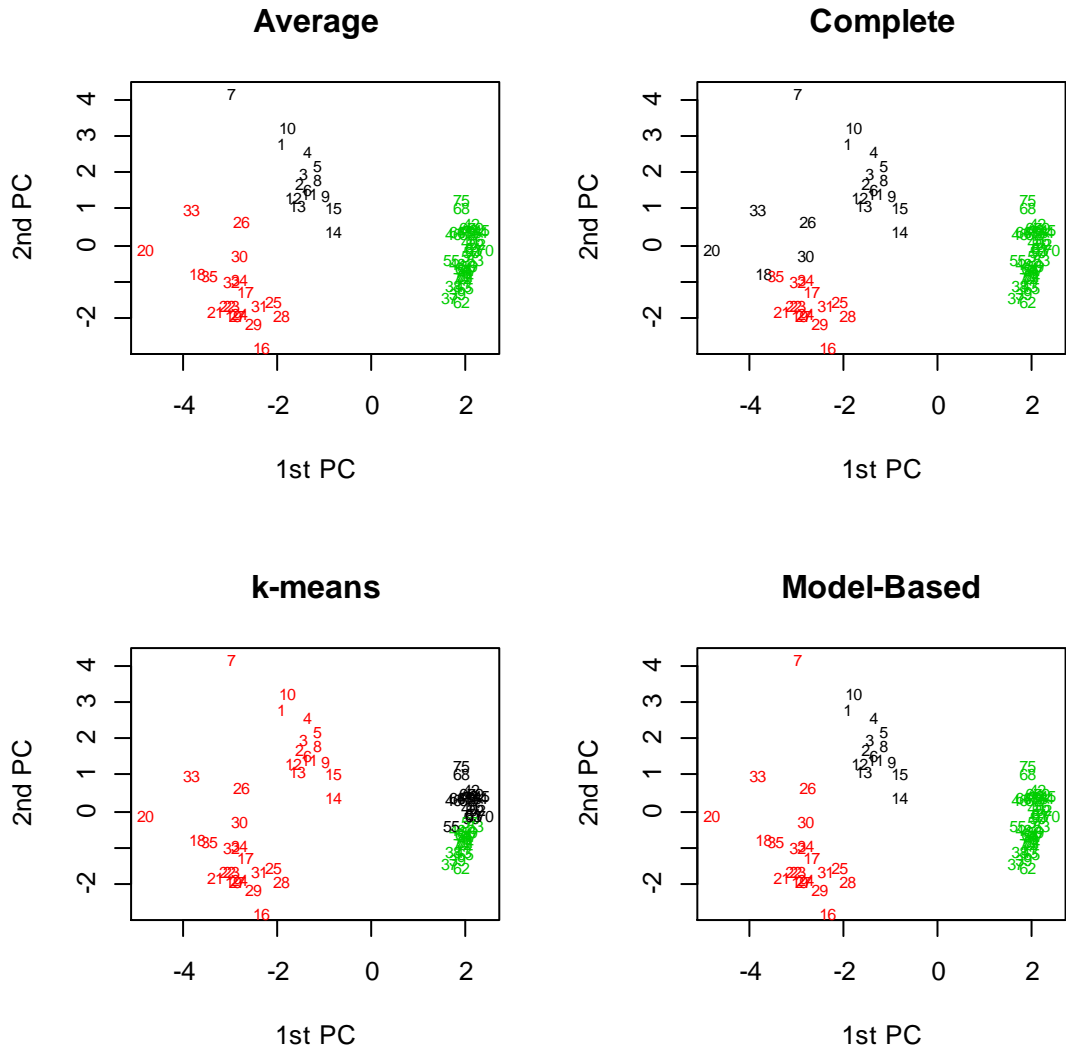
their corresponding variance will inevitably vary and standardization will be essential. If log transformation or log-ratio transformation results to a data set that the variables have comparable variances, a further standardization is not necessary, at least with respect to PCA performance. Another practical issue that may arise in PCA implementation is the presence of outliers. We discuss this issue at the dedicated to outliers paragraph as they affect not only PCA, but also cluster analysis performance.

**Example 1 (continued).** Using the simulated data of Example 1, we apply PCA on the standardized data and the first PC explains 57.1% of the total variation, while the second 19.8%. A graph on the two first PC will combine 77% of the total variability of the original data. This is a vast amount of the variability of the original data and it could be extremely useful in visualization of a clustering result, but yet a remaining 23% may be informative as well. This is why a plot of the first against the third PC may be also needed to be checked. Figure 6 plots the simulated data of dimension nine, at a two dimensional space using the first two principal components. The PC analysis on standardized, log-transformed with basis 10 and log-ratio transformed data is plotted and observations are labeled with respect to the true origin, since the data are simulated.



**Figure 6.** Data plotted in the first two PCs for standardized data (left) and log-ratio transformed data (right).

To use PCA as a graphical way of clustering assessment, we plot all under consideration clustering results into the first two PCs. Figure 7 plots clusterings listed in Table 4. In practice this list may be quite long. It is easy in our working example to verify that Average linkage and model-based method propose a compact result that captures the heterogeneity. According to Complete linkage, five observations seem to have been misclassified and k-means produces the less sensible result because the very compact group at the right of the plot is separated into two smaller groups.



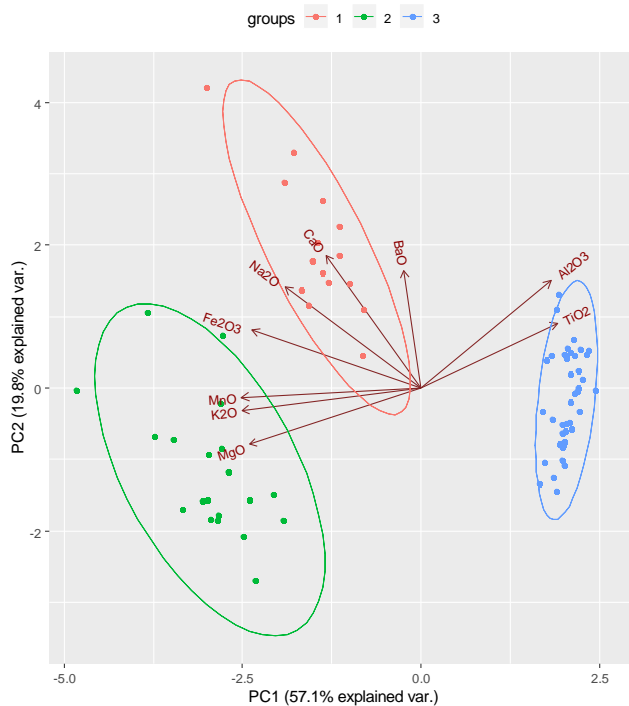
**Figure 7.** Clustering results from previous analysis on simulated data of Example1 plotted on the first two PCs.

PCA can also be used for characterizing the derived provenance groups. This can be done by using the loadings of PCA with respect to the original data. For the example, the loadings for the first two PCs are given in Table 5. Graphically, the loadings together with a scatter plot of the observations into the two first principal components are plotted in a biplot. The biplot of the simulated data is plotted in Figure 8. By inspection of Table 5 and Figure 8, we conclude that, for example, measurements in group 1 are characterized by higher values in CaO, Na<sub>2</sub>O and BaO. Group 2 is characterized by higher values in MgO, MnO and K<sub>2</sub>O and at the same time lower values in Al<sub>2</sub>O<sub>3</sub> and TiO<sub>2</sub>. The special characteristics of group 3, or in other words the factors which discriminate this group from the others is the high measurements in Al<sub>2</sub>O<sub>3</sub> and TiO<sub>2</sub>. Figure 8

presents the box-plots of three selected chemical elements for the simulated dataset, using the labels of the group membership as derived from Average Linkage which coincides with the true clustering. It is apparent that these selected elements, among others, can play the role of discriminating factors of the three groups.

	Al <sub>2</sub> O <sub>3</sub>	Fe <sub>2</sub> O <sub>3</sub>	MgO	CaO	Na <sub>2</sub> O	K <sub>2</sub> O	TiO <sub>2</sub>	MnO	BaO
PC1	0.30	-0.39	-0.40	-0.22	-0.32	-0.42	0.32	-0.42	-0.04
PC2	0.42	0.23	-0.22	0.52	0.40	-0.09	0.26	-0.04	0.46

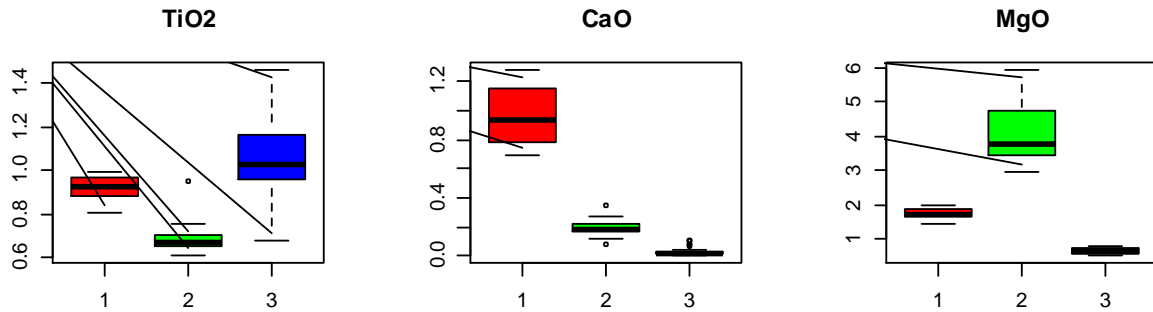
**Table 5.** Loadings of the first two PCs.



**Figure 8.** Biplot of the first two PCs obtained from the simulated data of Example 1.

Summarizing, PCA can be used as a first graphical test to investigate and possibly detect chemical compositional structure in the data. If this is present a cluster analysis would be a sensible next step of the analysis. Labeling cannot be possible at this stage, but still a certain structure, if present, can be detected. This suspected structure can further pursued by a cluster analysis. Moreover, if a structure is present and coherent all possible transformations would be able to capture this. If a structure is suggested only with some data transformations usually means something.

1  
2  
3  
4 PCA at a later stage of the analysis can be used for cluster validation and lastly, PCA can assist to  
5 derive conclusions on the special characteristics that each group possesses.  
6  
7  
8  
9  
10  
11



12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25 **Figure 8.** Box-Plots of selected original measurement variables.  
26

27  
28 Apart from PCA technique used as a data projection in to a lower dimension system of  
29 coordinates, other statistical learning or machine learning techniques can also be used. For  
30 example, factor analysis, multidimensional scaling.  
31  
32

### 33 **Categorical Data and Mixed-mode data**

34  
35  
36  
37 Apart from the geochemical continuous data, other source of information leading to quantitative  
38 discrete data may also be available. Such information may result from petrographic  
39 examinations, e.g. optical microscopy examination of thin-sections.  
40

41 To accommodate this type of information in the analysis, if quantitative, one can proceed with  
42 two ways. The information provided from discrete data can be used complimentary. A cluster  
43 analysis based on continuous data can be conducted independently and a cross examination from  
44 groups suggested from the discrete data can be used as confirmation or explanation of the cluster  
45 findings obtained from the chemical information only. Alternatively, but less often used in  
46 practice, one can merge both types of variables to a single, integrated analysis, that takes into  
47 account both types of data continuous and discrete. This type of analysis is introduced in  
48 archaeological application from Baxter et al. (2008) with the name mixed-mode analysis. The  
49 merits of such analysis are that it is more informative and avoids the problem of contradictive  
50 results suggested from two separate cluster analysis.  
51  
52

53 Baxter et al. (2008) methodology for the mixed-mode approach of the problem is based on  
54 a generalization of Gower coefficient of similarity and weighting the contribution of continuous  
55 and discrete data. This weighting aims to avoid the domination of binary data over the continuous  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 towards clustering. Papageorgiou and Moustaki (2005) propose latent class models as a model-  
5 based method appropriate to answer the cluster problem while the input variables can be  
6 continuous, binary, categorical ordinal or nominal. They successfully implement this  
7 methodology to an Archaeological dataset. Recently, Androulakis et al 2018 proposed a model-  
8 free approach for the mixed-mode problem based on a modification of Gower coefficient of  
9 similarity.

10  
11  
12  
13 However, although methodologies have proposed in the literature for dealing with the  
14 problem of different type of information adopting the integrated approach, in practice the use  
15 of this approach is limited. This is because special software, sometimes not freely available, is  
16 necessary or they are computational methods and therefore are demanding in technical and/or  
17 computational effort.  
18  
19  
20  
21  
22

### 23 **Outliers**

24  
25 Another characteristic of Archaeometric data is the presence of outliers, i.e. observations that  
26 deviate from the remaining sample. This may be due a deviate measurement of this observation  
27 in one or two variables, or a multivariate outlier where the difference of this measurement is  
28 more general and it may be due to an unexpected relation among the measured variables for this  
29 particular observation. For example, two or more variables may be strongly positively related,  
30 i.e. strong positive correlation, and the measurements on these variables indicate negative  
31 relation for a particular data point. If a plot of the data at these dimensions was possible, the bulk  
32 of data points would be expected to form an ellipsoidal (not spherical shape) due to the strong  
33 correlation, and the outliers would lie outside this ellipsoidal. Detection of univariate or  
34 multivariate outlier is essential for a statistical analysis because most of the methodologies  
35 presented in previous paragraphs are affected from the presence of outliers and their  
36 performance degenerate.  
37  
38  
39  
40  
41  
42

43 Detection of univariate variables is relatively straightforward problem and univariate  
44 techniques, such as boxplots, histograms etc. can identify those outliers, or better extreme values  
45 on certain variables. Detection of multivariate outliers is however a more challenging problem  
46 and a number of approaches have been proposed in the literature (see Filzmoser et al. 2005 and  
47 Rousseeuw et al 1990 amongst others). The most commonly used technique to identify  
48 multivariate outliers in the literature is the Mahalanobis distance. Assuming a set of  $p$ -  
49 dimensional observations  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ ,  $i = 1, 2, \dots, n$  Mahalanobis distance for  
50 observation  $i$  is defined as  
51  
52  
53  
54  
55

$$56 \quad d_M(\mathbf{x}_i) = \sqrt{(\mathbf{x}_i - \hat{\boldsymbol{\mu}})' \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x}_i - \hat{\boldsymbol{\mu}})}$$

57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 where  $\mu, \Sigma$  are the mean and variance covariance matrix of the theoretical population that has  
5 generated the sample and  $\hat{\mu}, \hat{\Sigma}$  are some estimates of those parameters based on the data. Usual  
6 estimates for  $\mu, \Sigma$  are the sample mean and sample variance covariance matrix  $\bar{x}, S$  respectively.  
7 If the population distribution is multivariate normal the distribution of  $d_M^2$  can be proved to be  
8 chi-square with  $p$  degrees of freedom,  $X_p^2$ . Under this assumption any observation that has a  $d_M$   
9 value at the outer, say 95% probability area, for the  $X_p^2$  distribution is characterized as extreme  
10 and therefore outlier. Note that  $d_M$  expression takes into account all measured variables,  
11 although the distribution is univariate, and therefore conclusions on multivariate outliers can be  
12 obtained. For example, if  $p = 14$  the cut-off point for a 95% confidence ellipsoidal is  $X_{14,0.95}^2 =$   
13 23.68. According this measure, any data point  $x_i$  for which the square of  $d_M(x_i)$  is greater than  
14 23.68 is considered as outlier.

15 However, estimates  $\bar{x}, S$  are sensitive to the presence of outliers themselves, or in statistical  
16 terminology, they are not robust as estimates when one or a small set of observations deviate  
17 from the rest. This has as a result Mahalanobis distance in turn to be affected from the presence  
18 of outliers. Since, its practicality as a measure, is essential when outliers are indeed present in  
19 the data, various modifications of Mahalanobis distance as a tool for outlier detection are  
20 proposed in the literature. Modifications are mainly consisted of using other, robust estimates,  
21 of  $\mu, \Sigma$  instead of  $\bar{x}$  and  $S$  (Rousseeuw and Van Zomeren, 1990).

22 Once outliers are identified we proceed with the cluster analysis by keeping those special  
23 observations aside. A further examination of the set of outliers with respect to their compositions  
24 will possibly reveal the reason why they diverse from the remaining data set. Outliers may be  
25 quite important observations and meaningful in archaeological terms, but as far as the clustering  
26 problem is concerned they do not cluster with any of the existing groups and they will only  
27 confound any clustering technique.

## 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 **5. Classification**

46  
47 Classification is a supervised multivariate technique and refers to the problem of classifying a  
48 new observation to one of the known identified groups in the data. It is an interesting problem  
49 in Archaeometry when a clustering at previous stage has been conducted and a number of groups  
50 corresponding to different origin have been established. If new samples, artifacts in this case, are  
51 available from another study classification will be the appropriate method to assign the new  
52 artifact or artifacts to one of the existing groups and consequently make conclusions about its  
53 origin.

54  
55 As in clustering, classification is a multivariate statistical problem and a number of  
56 approaches and techniques can be applied in this context. One approach is to consider the



1  
2  
3  
4 problem at a probabilistic manner where either a prediction model is used to predict the  
5 classification category for the new observation or the probability for the unclassified  
6 observation to belong to each one of the existing groups is calculated and assignment is based  
7 on the maximum probability. Within this framework belong the methodologies of linear  
8 regression, logistic regression, multiple logistic regression methods, Bayes classifier and Naive  
9 classifier.

10  
11  
12  
13 Discriminant Analysis (DA) is another widely used technique for classification problem  
14 where the task is to construct a classification rule which discriminates the groups. The rule can  
15 be linear or quadratic leading to linear discriminant and quadratic discriminant analysis  
16 respectively. The resulting rule may also be useful to visualize the data in a lower dimension  
17 space than the original dimension of the data. This visualization can be considered as an  
18 alternative to PCA projection. PCA and DA utilize a different criterion according to which the  
19 components are constructed. In particular, DA takes into account the information of existed  
20 groups and DA components result as the projection of the original data which maximizes the  
21 between groups variation. Other classification techniques are the K nearest neighbor and  
22 classification trees. James et al. (2013) is a reference book with emphasis the applications of  
23 the classification.

24  
25  
26  
27  
28  
29 We present here only an illustration of classification trees technique, because this is the  
30 less mathematical, no assumptions about the population distributions are required, it can  
31 accommodate both continuous and discrete data and most importantly the results are easy to  
32 interpret. A software is necessary, but the implementation is straight forward.

33  
34  
35  
36  
37 **Example 2.** For this experiment we assume that the dataset consists of a 73x24 data matrix  
38 with ten variables to represent continuous variables of chemical measurements and 14  
39 variables as categorical, binary or ordinal. Three groups have been identified and labels 'gr1',  
40 'gr2' and 'gr3' have assigned to those groups.  
41 Applying a classification tree method on the data using the group label results to the tree  
42 plotted in Figure 9. From the 24 variables, only three variables corresponding to equal number  
43 of nodes are sufficient in order to construct a discriminant rule. The rule is apparent and  
44 written on each branch of the tree. For example is the new subject has measurement feldspar  
45 variable 2,4 or 6, the right branch is suggested where another test will follow with respect the  
46 value of the new subject in Na<sub>2</sub>O and more specifically if it less than 0.925 or not.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

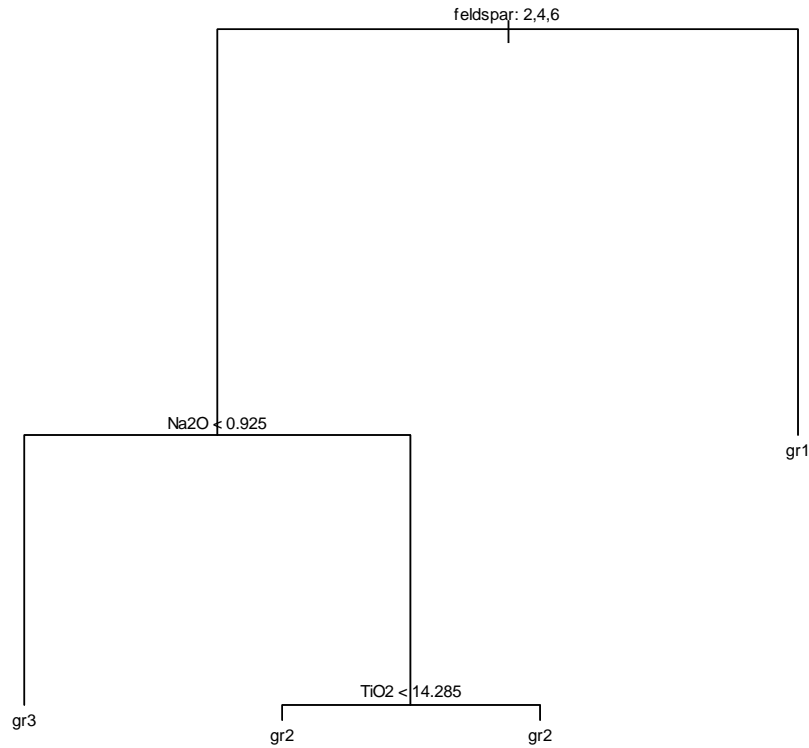


Figure 9. Classification tree.

## 6. Conclusions

Summarizing, we list the structure of the statistical analysis in steps.

**Step 1.** Run a preliminary analysis of the data, without any transformation, including univariate and bivariate study. Identify departure from normality, long tail for variables, outliers (univariate or multivariate using Mahalanobis distance) and possible structure in the data.

**Step 2** Data transformation.

**Step 3.** Implement PCA to transformed data in order to (i) identify or confirm outliers located at step 1 and (ii) check if data (using all measured variables) suggest the existence of different groups.

**Step 4** If PCA gives grounds to a non-homogeneous dataset, implement cluster analysis by using various methods and ways of data transformations. Verify outlier indicated by PCA

1  
2  
3  
4 or Mahalanobis distance and confirm with data inspection. Remove those and repeat  
5 cluster analysis. Identify clear groups, inference about these groups and justify the result.  
6 Set aside the distinct groups and repeat cluster analysis.  
7  
8

9  
10 **Step 5.** In the process of cluster analysis, comparison and assessment of the clustering  
11 results as proposed from various methods is performed.  
12

13  
14 **Step 6.** For the compact groups of the analysis proceed with their characterization and  
15 determine the discriminating factors among groups.  
16

17  
18 **Step 7.** If discrete data from microscopy study are also available and have a quantitative  
19 nature, a cluster analysis for this type of data can also be implemented and results with  
20 respect to clusters suggested can be seen in comparison with the clustering suggested from  
21 the continuous data.  
22

23  
24 **Step 8.** Mixed-mode cluster analysis is suggested if quantitative discrete data are also  
25 available.  
26

27  
28 **Step 9.** A classification analysis is appropriate when the aim is to classify a new subject,  
29 an artifact, into one of the existing identified groups.  
30  
31

32  
33 Closing this article, I would like to borrow a paragraph from Whallon (1984) that Baxter (2015)  
34 also uses as an introduction to chapter 11. The paragraph is “Archaeologists are ill-trained to, the  
35 rigorous and logical thought necessary form an informed use of quantitative methods, while the  
36 rare statisticians who have tried their hands at archaeology, typically have understood the nature  
37 of archaeological data, questions, and models only partially, vaguely, or incorrectly, so that their  
38 efforts are usually no better than the archaeologist’s own”. As mentioned in the Introduction  
39 section, Archaeologists now-days are more educated in using quantitative tools. However,  
40 interaction and collaboration between the archaeologist and the statistician are key components  
41 that cannot be substituted from technology and training. To my opinion any result is safe only  
42 when it can be confirmed from both sciences.  
43  
44  
45  
46  
47  
48  
49  
50  
51

## 52 **Bibliography**

53 Aitchison, J (1986) *The Statistical Analysis of Compositional Data*. Chapman & Hall, London.  
54

55 Binford, L. R. (1964) A Consideration of Archaeological Research Design, *American Antiquity*,  
56 **29**, 4, 425-441.  
57  
58  
59  
60  
61

1  
2  
3  
4 Banfield, J.D. and Raftery, A.E. (1993) Model-Based Gaussian and Non-Gaussian Clustering.  
5 *Biometrics*, Vol. 49, 803-821.  
6

7  
8 Baxter, M.J. (2015) Exploratory Multivariate Analysis in Archaeology (Foundations of  
9 Archaeology). 2<sup>nd</sup> Edition. Eliot Werner Publications/Percheron Press.  
10

11  
12 Baxter, M.J. (2015) Spatial k-means clustering in archaeology – variations on a theme. Working  
13 paper – November 2015 (accessed in Academia.edu).  
14

15  
16 Baxter, M.J., Beardah, C.C., Papageorgiou, I., Cau, P.M., Day, P.M. And Kilikoglou, V. (2008). On  
17 Statistical Approaches To The Study Of Ceramic Artefacts Using Geochemical And Petrographic  
18 Data. *Archaeometry*, **50**, 142–157. <https://doi.org/10.1111/J.1475-4754.2007.00359.X>  
19  
20

21  
22  
23 Baxter, M. J. (2006) A review of supervised and unsupervised pattern recognition  
24 in archaeometry. *Archaeometry*, **48**, 4, 671–694.  
25

26  
27 Baxter, M.J. (2001) Statistical modelling of artefact compositional data. *Archaeometry* **43**, 1, 131-147.  
28

29  
30 Baxter, M.J. (2008) Mathematics, Statistics and Archaeometry – The last 50 years or so. *Archaeometry*  
31 **50**, 6, 968–982.  
32

33  
34 Baxter, M.J. (1995) Standardization and Transformation In Principal Component Analysis, with  
35 Applications to Archaeometry. *Applied Statist.* (1995) **44**, No.4, pp. 513-527.  
36

37  
38 Beardah, C. C., Baxter, M. J., Cool, H. E. M., and Jackson, C. M. (2003) Compositional data analysis  
39 of archaeo-logical glass: problems and possible solutions, in CoDaWork'03: Compositional Data  
40 Analysis Workshop, Girona, Spain; available at  
41 [http://ima.udg.es/Activitats/CoDaWork03/paper\\_baxter\\_Beardah2.pdf](http://ima.udg.es/Activitats/CoDaWork03/paper_baxter_Beardah2.pdf)  
42  
43

44  
45 Bieber, A.M. JR., Brooks, D.W., Harbottle, G. and Sayre, E.V. (1976) Application of Multivariate  
46 Techniques to Analytical Data on Aegean Ceramics. *Archaeometry* 18, 59-74.  
47

48  
49 Buxeda i Garrigós, J. (1999) Alteration and Contamination of Archaeological Ceramics: The  
50 Perturbation Problem. *Journal of Archaeological Science*, 26, 295-313.  
51

52  
53 Charrad, M., Ghazzali, N., Boiteau, V. and Niknafs, A. (2014). NbClust: An R Package for  
54 Determining the Relevant Number of Clusters in a Data Set. *Journal of Statistical Software*, 61(6),  
55 1-36. URL <http://www.jstatsoft.org/v61/i06/>.  
56

57  
58 Everitt, B.S., Landau, S. and Leese, M. (2011) *Cluster Analysis*, 5<sup>th</sup> Edition. John Wiley  
59 & Sons.  
60

1  
2  
3  
4 Everitt, B.S. and Dunn, G. (2010) *Applied Multivariate Data Analysis*. 2<sup>nd</sup> Edition. John Wiley &  
5 Sons.  
6

7  
8 Filzmoser, P., Garrett, R.G. and Reimann, R. (2005) Multivariate outlier detection in exploration  
9 geochemistry. *Computers & Geosciences*, **31**, 579-587.  
10

11  
12 Gordon A.D. (1998) Cluster Validation. In: Hayashi C., Yajima K., Bock HH., Ohsumi N., Tanaka  
13 Y., Baba Y. (eds) *Data Science, Classification, and Related Methods*. Studies in Classification,  
14 Data Analysis, and Knowledge Organization. Springer, Tokyo.  
15

16  
17 Hubert, L. and Arabie, P. (1985) Comparing Partitions, *Journal of the Classification*, **2**, 193-218.  
18

19  
20 James, G., Witten, D., Hastie, T. and Tibshirani, R. (2013). *An Introduction to Statistical Learning*  
21 *with Applications in R*. Springer, New York.  
22

23  
24 Jolliffe, I. T. (2002). *Principal Component Analysis*. Springer.  
25

26  
27 Kassambara, A. (2017) *Practical Guide To Cluster Analysis in R. Unsupervised Machine Learning:*  
28 *Volume 1 (Multivariate Analysis)*. STHDA Publishing.  
29

30  
31 Martín-Fernández, J.A., Buxeda i Garrigós, J. and Pawlowsky-Glahn, V. (2015) Logratio Analysis in  
32 Archeometry: Principles and Methods. In Barcelo, J.A. and Bogdanovic, I. (Eds) *Mathematics and*  
33 *Archaeology*. CPC press, Boca Raton FL, 178-189.  
34

35  
36 Palarea-Albaladejo, J. and Martin-Fernandez, J.A. (2015) zCompositions — R package for  
37 multivariate imputation of left-censored data under a compositional approach. *Chemometrics*  
38 *and Intelligent Laboratory Systems*. **143**, 85-96.  
39

40  
41 Papageorgiou, Ioulia. (2018) Cluster Analysis. *The SAS Encyclopedia of Archaeological Sciences*,  
42 John Wiley & Sons. <https://doi.org/10.1002/9781119188230.saseas0099>  
43

44  
45 Papageorgiou, I., Baxter, M.J. and Cau, M.A (2001) Model-Based Cluster Analysis of Artefact  
46 Compositional Data. *Archaeometry*, **43**, 571-588.  
47

48  
49 Papageorgiou, I. and Moustaki, I. (2005) Latent class models for mixed variables with  
50 applications in Archeometry. *Computational Statistics & Data Analysis*, **48**, 659 – 675.  
51

52  
53  
54  
55 Pawlowsky-Glahn, V. & Buccianti, A. (2002) *International Journal of Earth Sciences*, **91**, 357–  
56 368. <https://doi.org/10.1007/s005310100222>  
57

1  
2  
3  
4 Pawlowsky-Glahn, V. & Buccianti, A. (2011) *Compositional Data Analysis: Theory and*  
5 *Applications*. Wiley.

6  
7  
8 Rogers, S. and Girolami, M. (2016) *A First Course in machine Learning*. 2<sup>nd</sup> Edition. Chapman and  
9 Hall.

10  
11  
12 Rousseeuw, P. and Van Zomeren, B. (1990). Unmasking multivariate outliers and leverage points.  
13 *Journal of the American Statistical Association*, **85**, 633-639.

14  
15  
16 Rousseeuw, P.J. (1987) Silhouettes: A graphical aid to the interpretation and validation of cluster  
17 analysis. *Journal of Computational and Applied Mathematics*, 20, 53-65.

18  
19  
20 Saraçlı, S., Doğan, N. & Doğan, İ. (2013) Comparison of hierarchical cluster analysis methods by  
21 cophenetic correlation. *J Inequal Appl* 2013, 203. doi:10.1186/1029-242X-2013-203

22  
23  
24 Shotwell, M.S. (2013). profdpm: An R Package for MAP Estimation in a Class of Conjugate Product  
25 Partition Models. *Journal of Statistical Software*, 53, 1-18. URL  
26 <http://www.istatsoft.org/v53/i08/>.

27  
28  
29 Steinley, D. (2006) K-means clustering: A half-century synthesis. *British Journal of Mathematical*  
30 *and Statistical Psychology*, 59, 1–34.

31  
32  
33 Whallon, R. (1984) Unconstrained Clustering for the Analysis of Spatial Distributions in  
34 Archaeology. In *Intrasite Spatial Analysis In Archaeology*, edited by H. J. Hietala, pp. 242-277.  
35 Cambridge University Press, New York.

36  
37  
38 Xu, R. and Wunsch-II, D.C. (2008). *Clustering*. Wiley, John & Sons, Inc.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65