# Practice 4 SPSS and RCommander

## Cluster Analysis

It is a class of techniques used to classify *cases* (or *variables*) into **groups** that are relatively **homogeneous** *within* themselves, and **heterogeneous** *between* each other, on the basis of a defined set of variables. These groups are called **clusters**.

Example: clustering of consumers according to their attribute preferences or similar behaviours and characteristics. Then, clusters of similar brands or products can help identifying market opportunities.

### General Steps to conduct a Cluster Analysis

i.   Select a distance measure.
ii.  Select a clustering algorithm.
iii. Determine the number of clusters.
iv.  Validate the analysis.

## Clustering procedures

### Hierarchical procedures

*Agglomerative* (start from *n* clusters, to get to 1 cluster)

*Divisive* (start from 1 cluster, to get to *n* cluster)

### Non hierarchical procedures

K-means clustering

**Agglomerative clustering**

Linkage methods

*Single linkage* (minimum distance)

*Complete linkage* (maximum distance)

*Average linkage*


The distance between two clusters is defined as the difference between the centroids (cluster averages)


# K-means clustering

The number k of cluster is fixed

An initial set of *k seeds* (aggregation centres) is provided

> First *k* elements or

> Other seeds

Given a certain treshold, all units are assigned to the nearest cluster seed

New seeds are computed

Go back to step 3 until no reclassification is necessary

Units can be reassigned in successive steps (*optimising partioning*)

# Hierarchical vs Non hierarchical methods

## Hierarchical clustering

**No** decision about the number of clusters.

Problems when data contain a high level of error.

Can be very slow…

Initial decision are more influential (one-step only).

## Non hierarchical clustering

Faster, more reliable.

Need to specify the number of clusters (arbitrary).

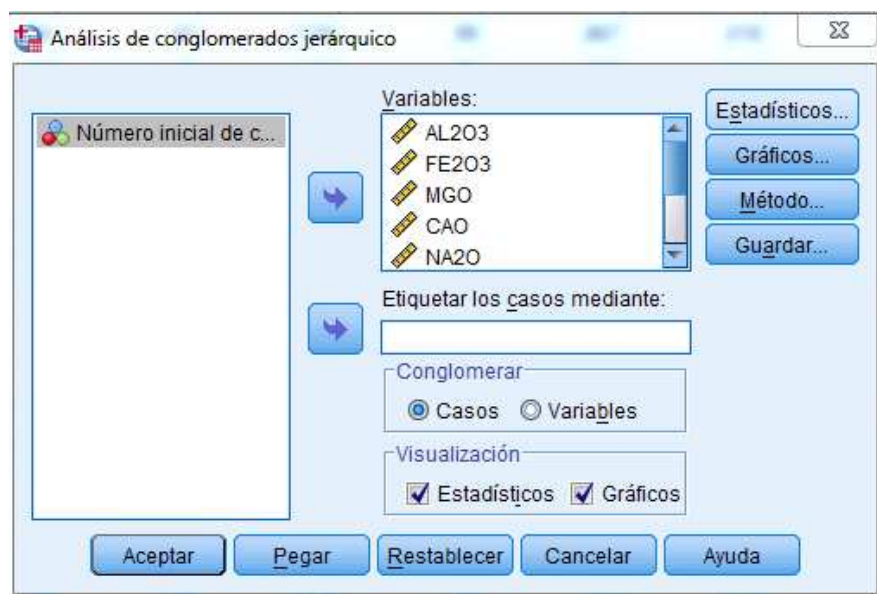Need to set the initial seeds (*arbitrary*).

# Hierarchical cluster analysis

Here, we will be use the file **pottery.txt**

The data give the chemical composition of 48 specimens of *Romano-British* pottery, determined by atomic absorption spectrophotometry, for nine oxides. In addition to the chemical composition of the pots, the kiln site at which the pottery was found is known for these data.

For these data, interest centres on whether, on the basis of their chemical compositions, the pots can be divided into distinct groups, and how these groups relate to the *kiln* site.

Go to Analyze, followed by Classify, and then Hierarchical Cluster. Drag drop all variables except *kiln*.

Analizar  →  Clasificar  →  Conglomerados Jerarquicos

Press in **Estadisticos**



Press in **Graficos**

Press in **Metodo**



Press in **Guardar**



Hierarchic classifications may be represented by a two-dimensional diagram known as a **dendrogram**, which illustrates the fusions made at each stage of the analysis.

We derive the three-cluster solution by cutting the dendrogram at a *height* of 10. Our interest is now a comparison with the *kiln* sites at which the pottery was found.
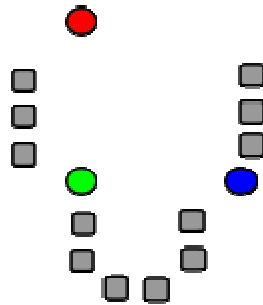
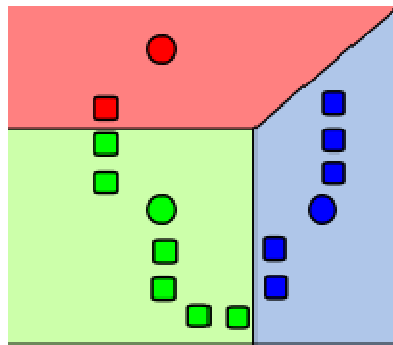Analizar  →  Estadisticos Descriptivos  →  Tablas de Contingencia

The contingency table shows that cluster 1 contains all pots found at *kiln* site number one, cluster 2 contains all pots from *kiln* sites number two, and cluster three collects pots from *kiln* sites four and five. In fact, so the clusters actually correspond to pots from three different regions.
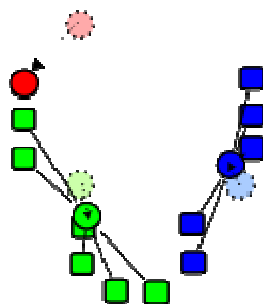
# K-means clustering
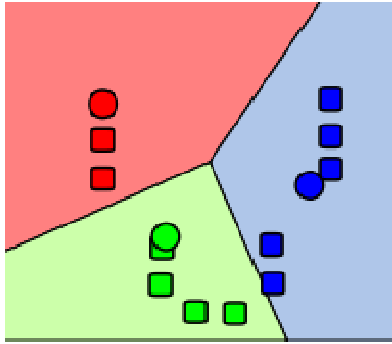
Basic ideas and scheme (from *Wikipedia*):



**1**) *k* initial *means* (in this case *k*=3) are randomly selected from the data set (shown in color).



**2**) *k* clusters are created by associating every observation with the nearest mean. The partitions here represent the *Voronoi diagram* generated by the means.



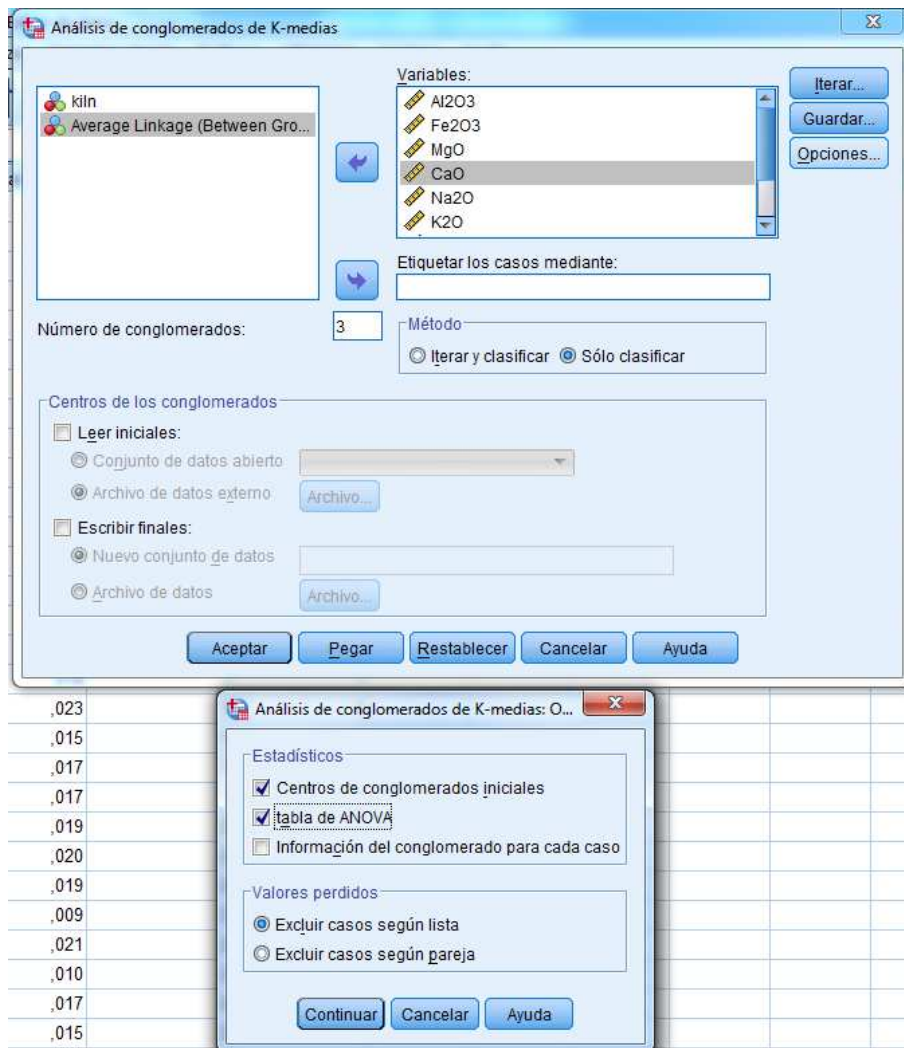**3**) The *centroid* of each of the *k* clusters becomes the new means.

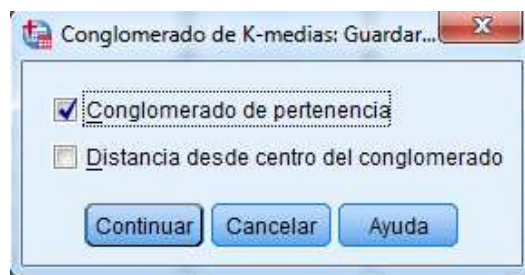**4**) Steps 2 and 3 are repeated until convergence has been reached.

We use the same data **pottery.txt** and set 3 as the number of clusters.

Go to Analyze, followed by Classify, and then K means Cluster. Drag drop all variables except *kiln*.

Analizar → Clasificar → Conglomerados K medias
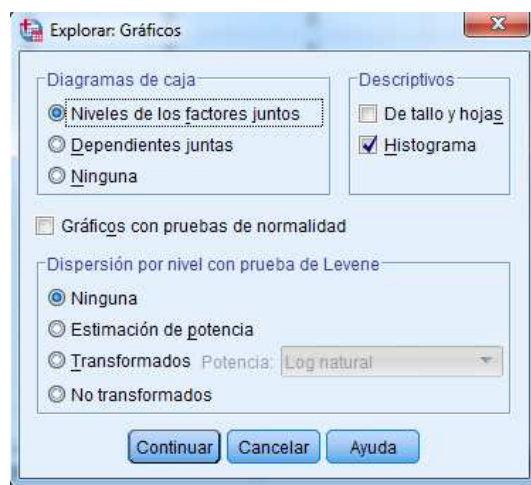
Press in **Guardar**



We derive the three-cluster solution and we make now a comparison with the *kiln* sites at which the pottery was found.

Analizar → Estadisticos Descriptivos → Tablas de Contingencia

We obtain similar results as in the hierarchical clustering case.

It may be very informative to compute the descriptive statistics on the original variables for each cluster in order to interpret them.

Analizar → Estadisticos Descriptivos → Explorar

# Cluster Analysis with RCommander

## Hierarchical cluster analysis

Here, we will be use the file **pottery.txt**

The data give the chemical composition of 48 specimens of *Romano-British* pottery, determined by atomic absorption spectrophotometry, for nine oxides. In addition to the chemical composition of the pots, the kiln site at which the pottery was found is known for these data.

For these data, interest centres on whether, on the basis of their chemical compositions, the pots can be divided into distinct groups, and how these groups relate to the *kiln* site.
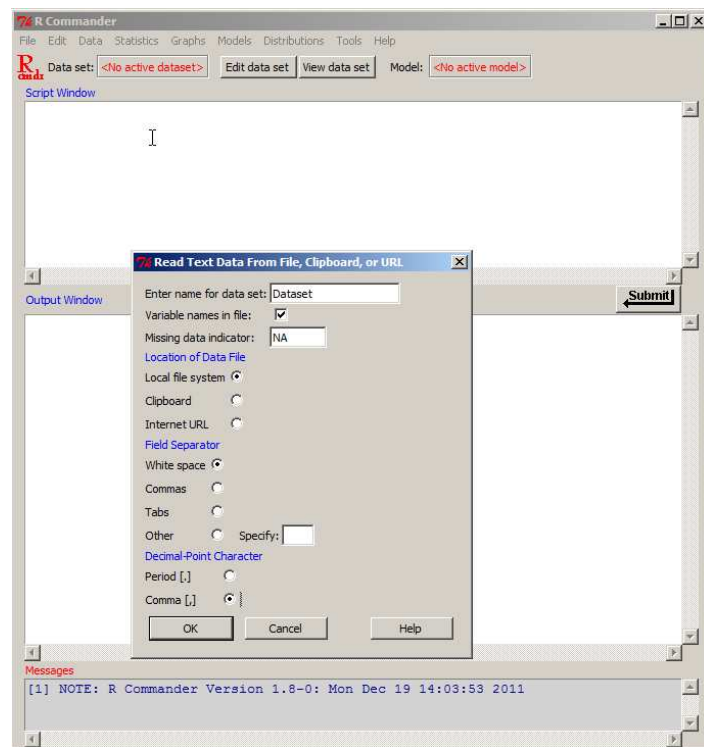
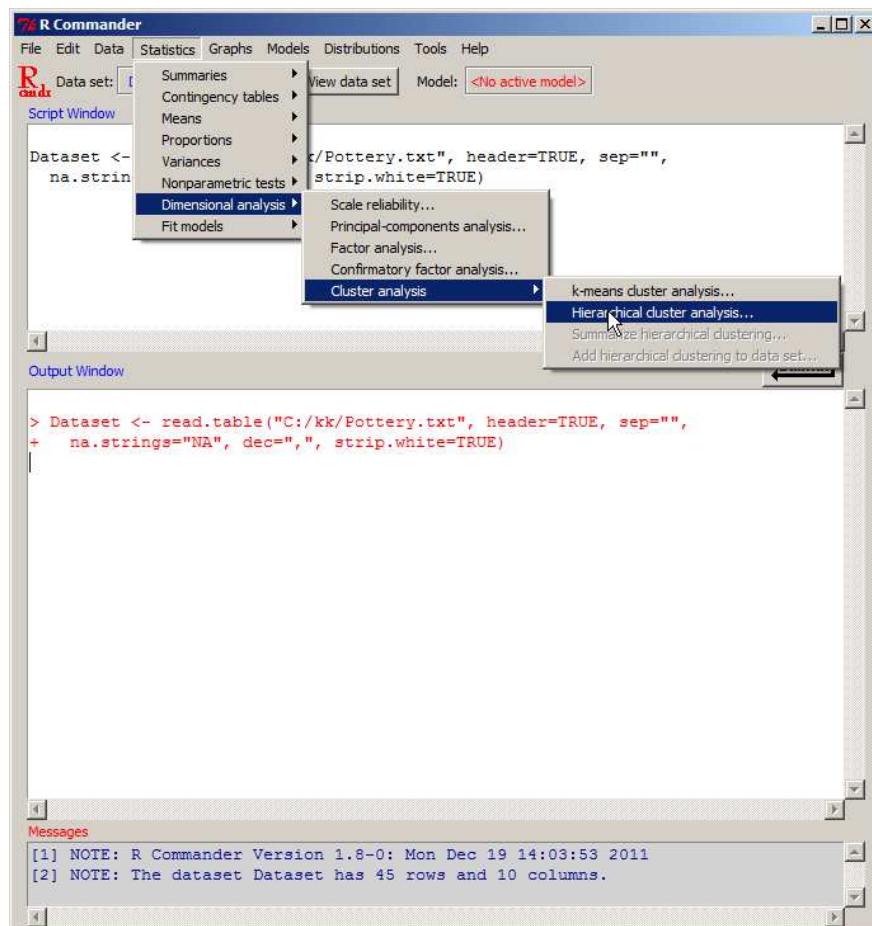In the working panel of R, type

**library(Rcmdr)**

Go to

Data  →  Import Data  →  From text File…

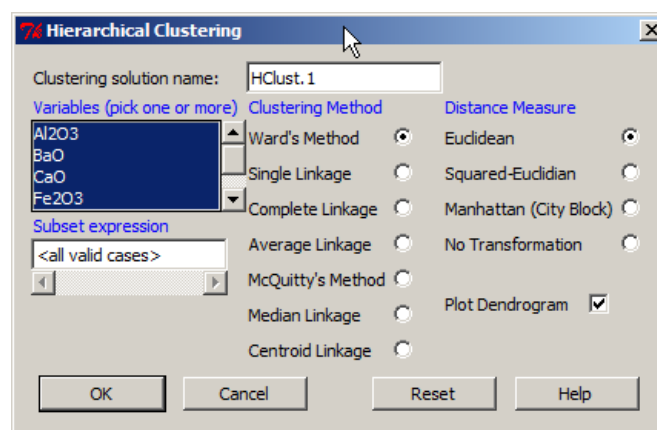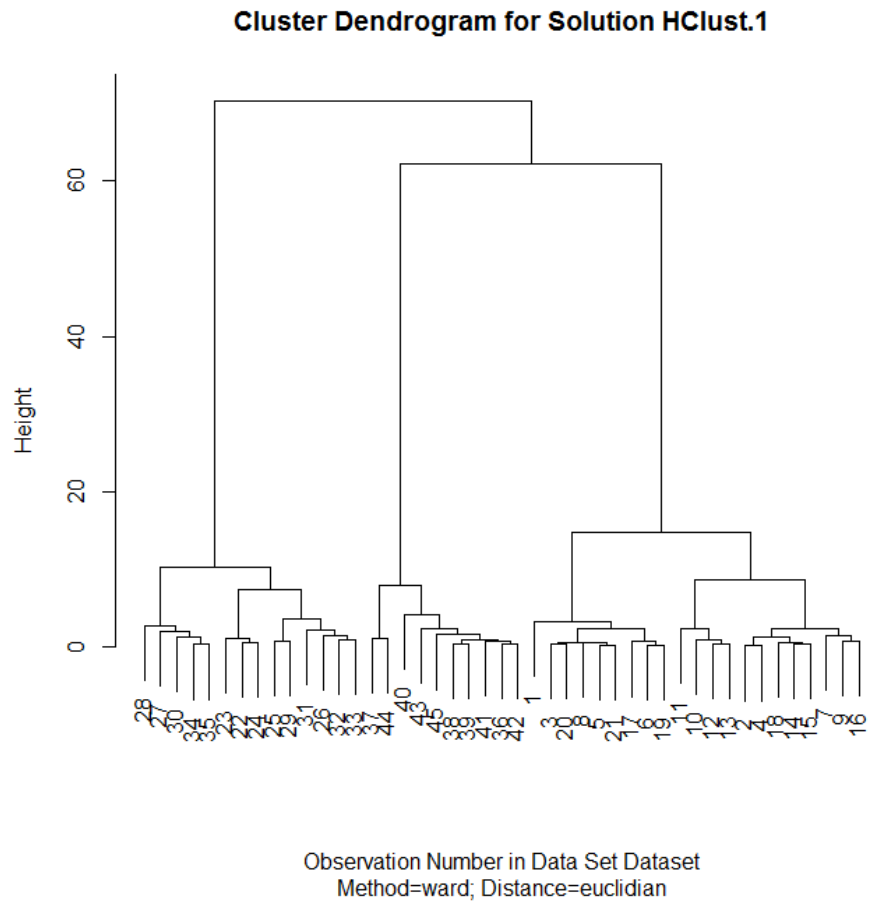Take the file **pottery.txt**

Go to

Statistics → Dimensional analysis → Cluster analysis → Hierarchical cluster analysis
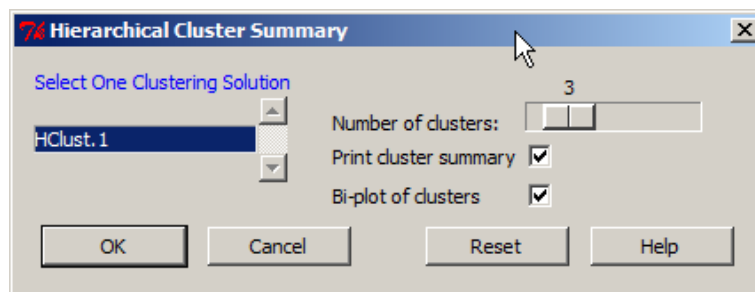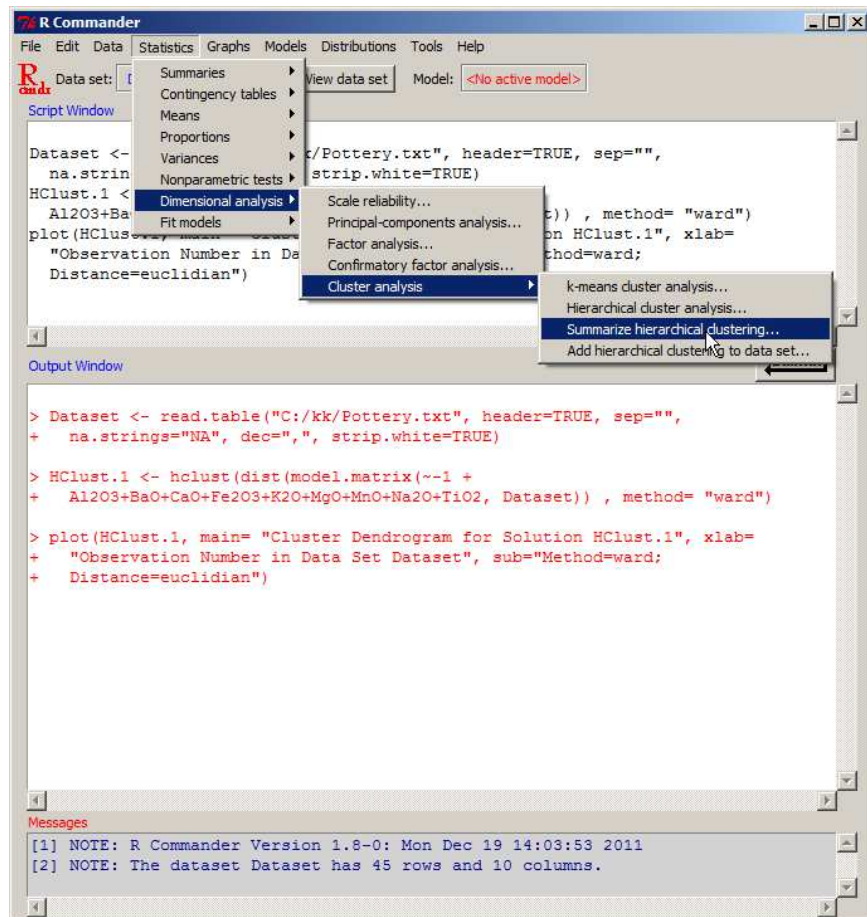
Include all variables **except** *kiln*

**Cluster Dendrogram for Solution HClust.1**

Height

60

40

20

0

Observation Number in Data Set Dataset
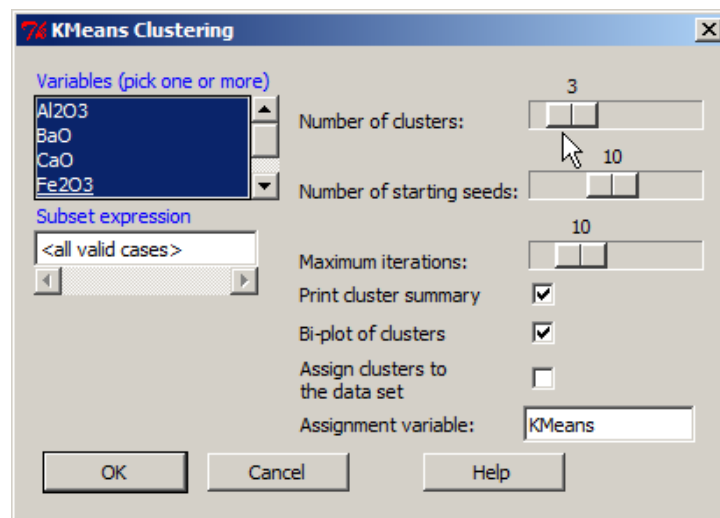Method=ward; Distance=euclidian

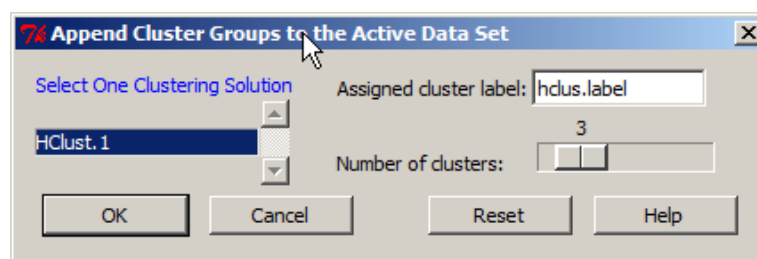Summarize results:

# K-means clustering

Statistics  →  Dimensional analysis  →  Cluster analysis  →  K-means cluster analysis
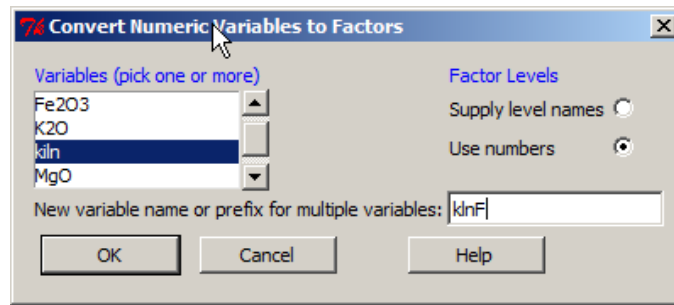
Include all variables **except** *kiln*



It may be very informative to compute the descriptive statistics on the original variables for each cluster in order to interpret them.



Convert variable *kiln* to type **factor**

Compute a table to check groups