**Weierstraß-Institut für**
**Angewandte Analysis und Stochastik**

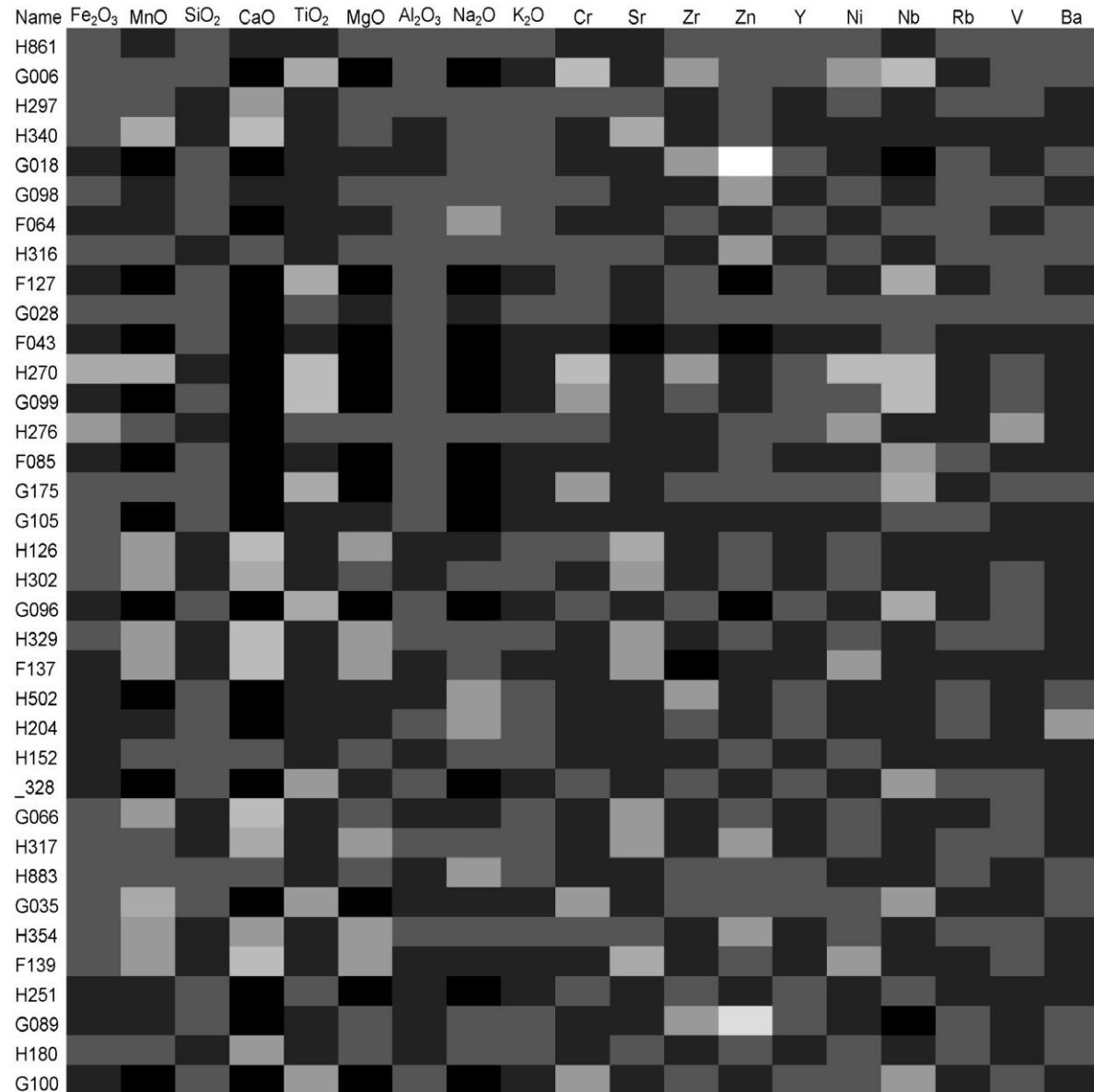# Visualisation and Cluster Analysis

Hans-Joachim Mucha

# Outline

- Introduction

- Hierarchical and partitional clustering

- Multivariate projection

- Dendrograms

- Mapping findspots

- Summary

Usually, the starting point for clustering is a $I \times J$ data matrix $\mathbf{X} = (x_{ij})$ with $I$ observations and $J$ variables.

**Application to archaeometry**
Snapshot of a fingerprint of the data as it comes: measurements of $J = 19$ oxides and elements of $I = 613$ Roman tiles from locations in *Germania Superior*. The darker the gray the higher is the measurement value.

# Introduction

Roman tiles (tegula) with a stamp from different locations in *Germania Superior*.
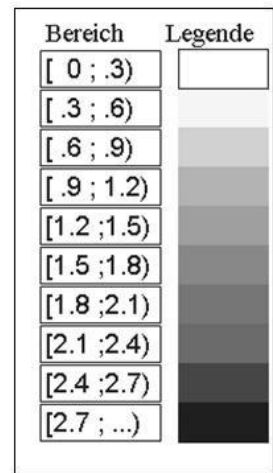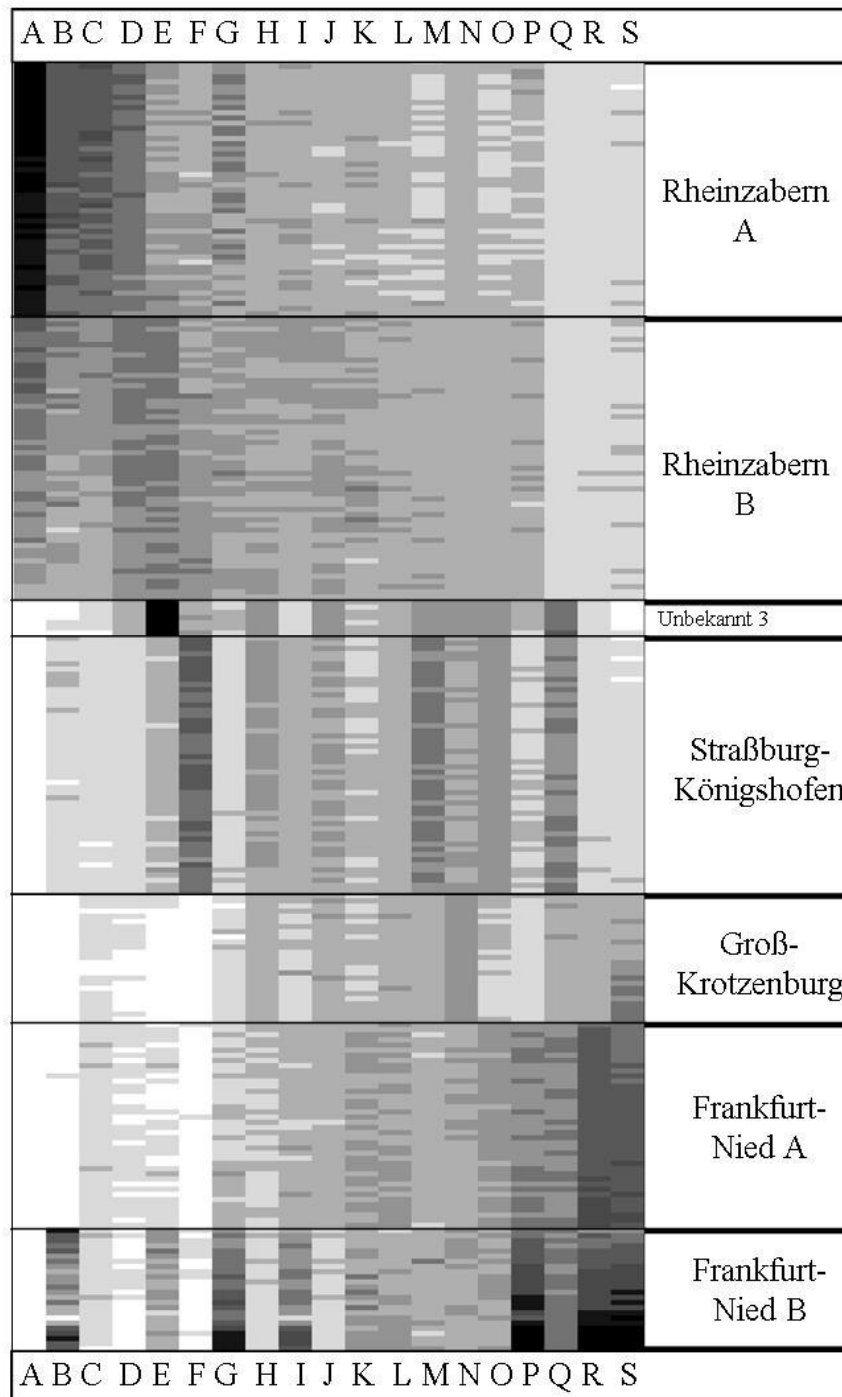
(More about Roman tiles and their multivariate statistical analysis at the website **http://www.ziegelforschung.de** .)

# Introduction

Aim: Clustering is used to group objects such that similar ones are collected in the same cluster.
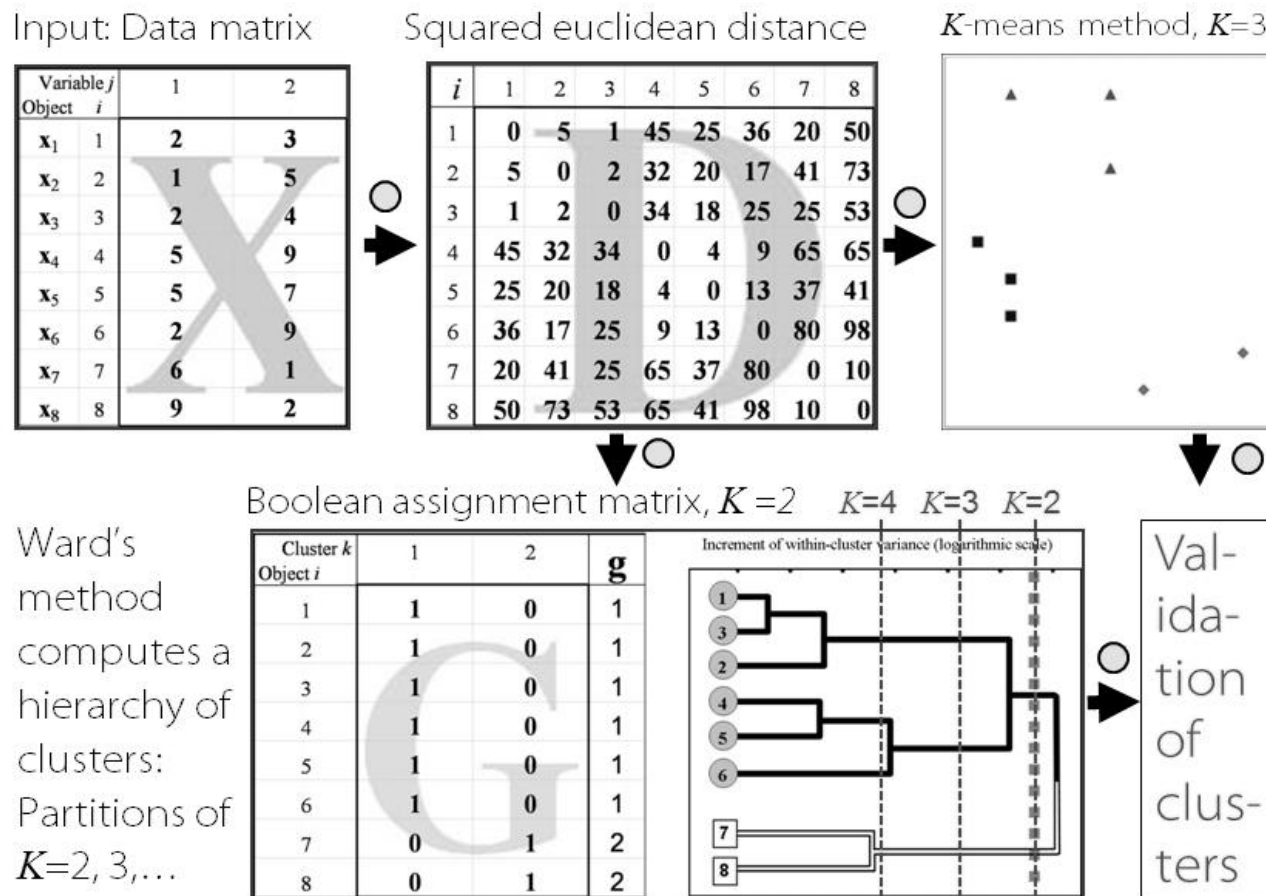
**Tiles** (same data as before): Here, the final aim of clustering is to find proveniences. Then the data can be ordered by both the clusters and some projection scores such as the first principal component to make an ordered fingerprint.
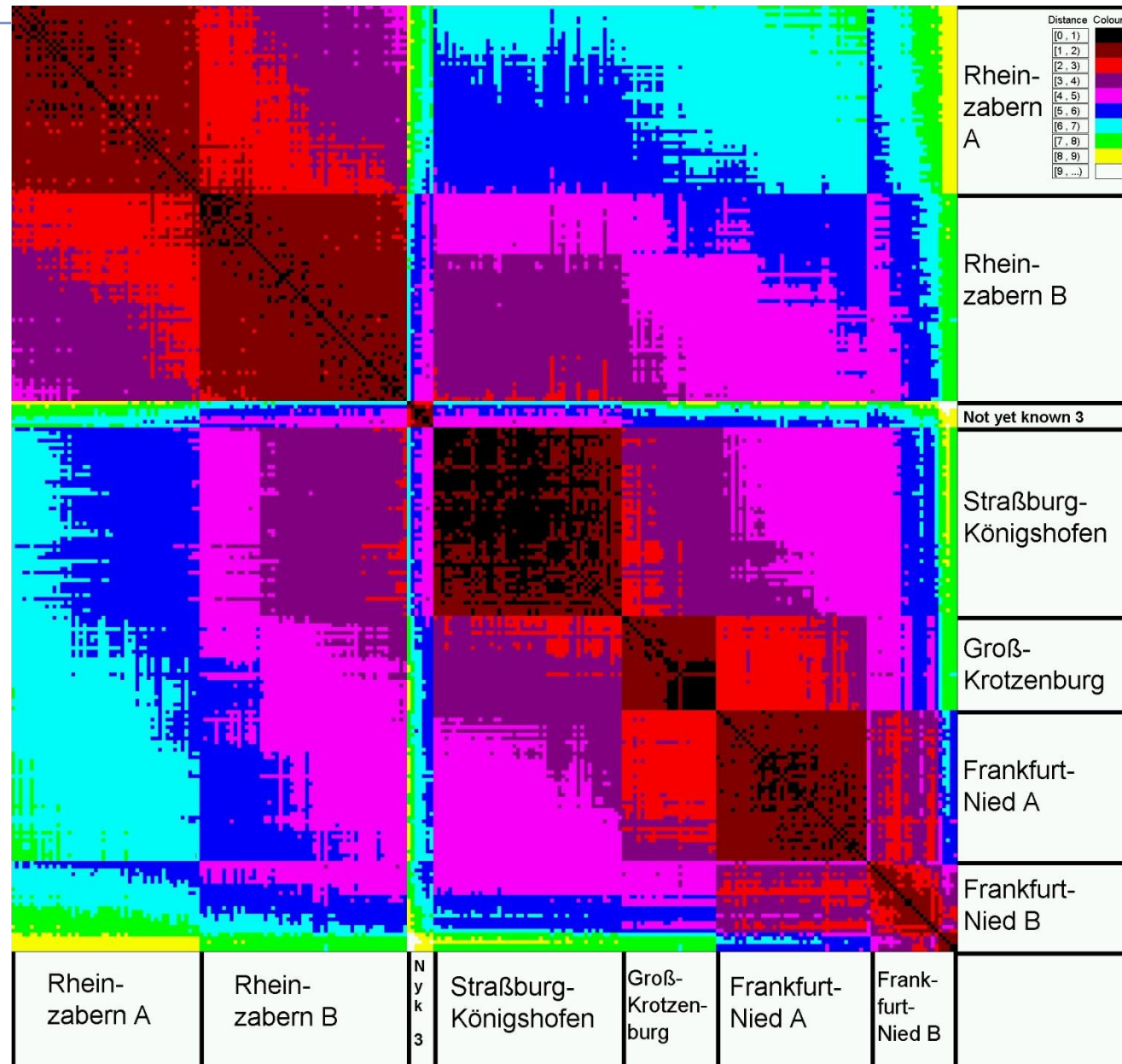
Sometimes, both hierarchical and partitional clustering can be formulated as pairwise data clustering, where instead of $\mathbf{X}=(x_{ij})$ a distance (proximity) matrix $\mathbf{D}=(d_{il})$ is used.

# Introduction

A heatmap of pairwise distances allows two-dimensional visualisations of arbitrary high-dimensional data.

**Tiles** (same data as before): Ordered heatmap of the Euclidean distance matrix $\mathbf{D}$ between the observations.

# Introduction

A heatmap of the Euclidean distances of random generated 20 dimensional Gaussian data.

A built-in analysis function of Excel 2003 was used here to generate the Gaussian data. Obviously, these generated data are not random.

# Hierarchical and partitional clustering

To handle both **hierarchical and partitional** methods at the same time in an unique fashion, we would like to focus on Gaussian model-based cluster analysis in its simplest setting. Concretely, the **sum of squares** (SS) has to be minimised. (By the way, a similar formulation can be derived for the logarithmic SS criterion.)

Starting point is a distance matrix $\mathbf{D_Q} = (d_{il})$ with pairwise **(weigthed) squared euclidean distances** as elements:

$$d_{\mathbf{Q}}(\mathbf{x}_i, \mathbf{x}_l) = (\mathbf{x}_i - \mathbf{x}_l)^T \mathbf{Q}(\mathbf{x}_i - \mathbf{x}_l).$$

Here the matrix $\mathbf{Q}$ is diagonal with weights $q_{jj} > 0$. The SS criterion $V_k$ based on $\mathbf{D_Q}$ has to be minimised with respect to a fixed number of clusters $K$:
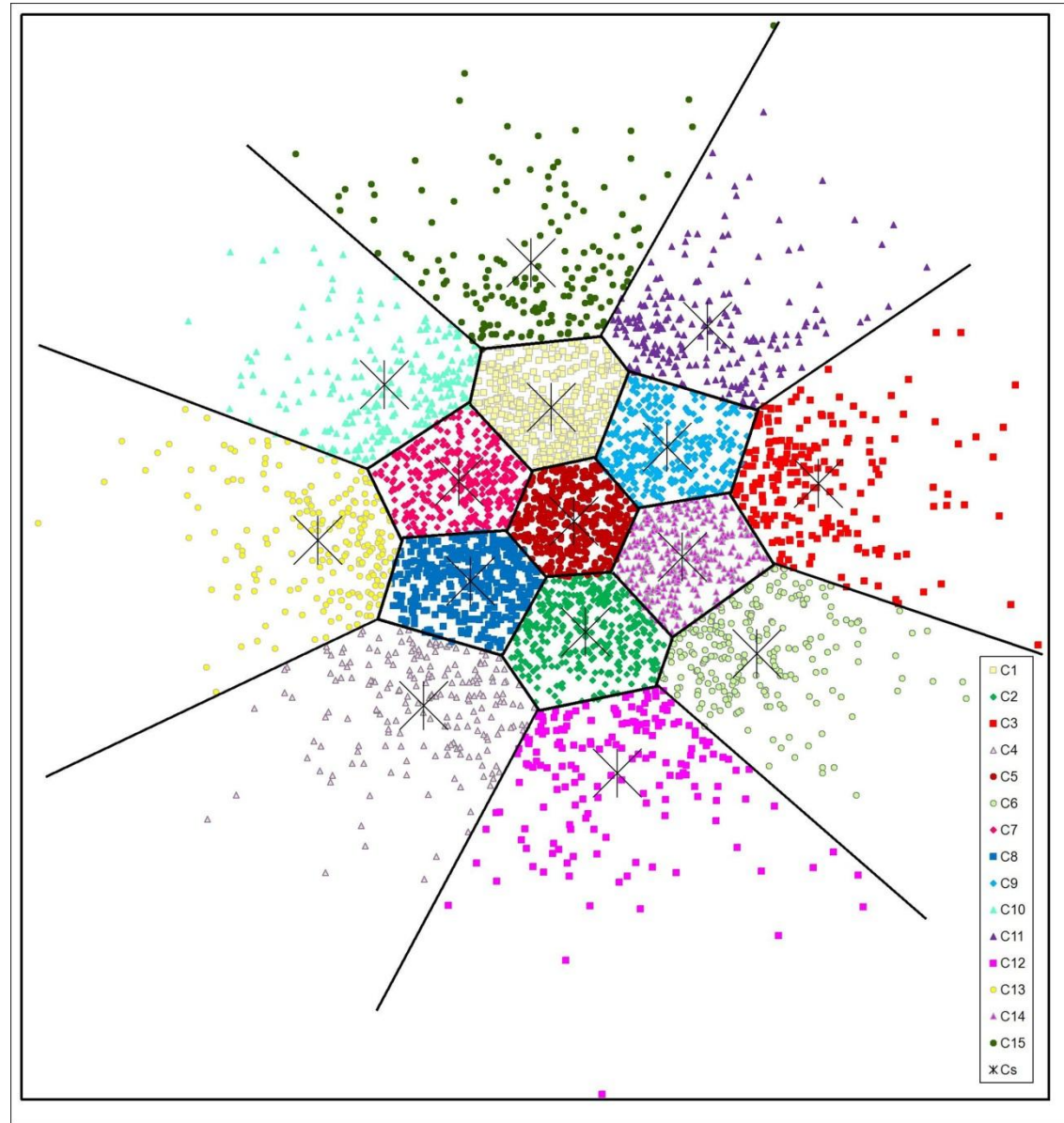
($m_i$ and $M_k$: weight of object $i$ and cluster $C_k$, respectively) .

$$V_K = \sum_{k=1}^{K} \frac{1}{M_k} \sum_{i \in C_k} m_i \sum_{\substack{l \in C_k \\ l > i}} m_l d_{\mathbf{Q}}(\mathbf{x}_i, \mathbf{x}_l)$$

Scatterplot of the result of the partitional $K$ means clustering of bivariate Gaussian data without a cluster structure into $K$ = 15 clusters. The result is a so-called Voronoi (or Dirichlet) tessellation. In the plot, additionally, the cluster centres are marked by stars.

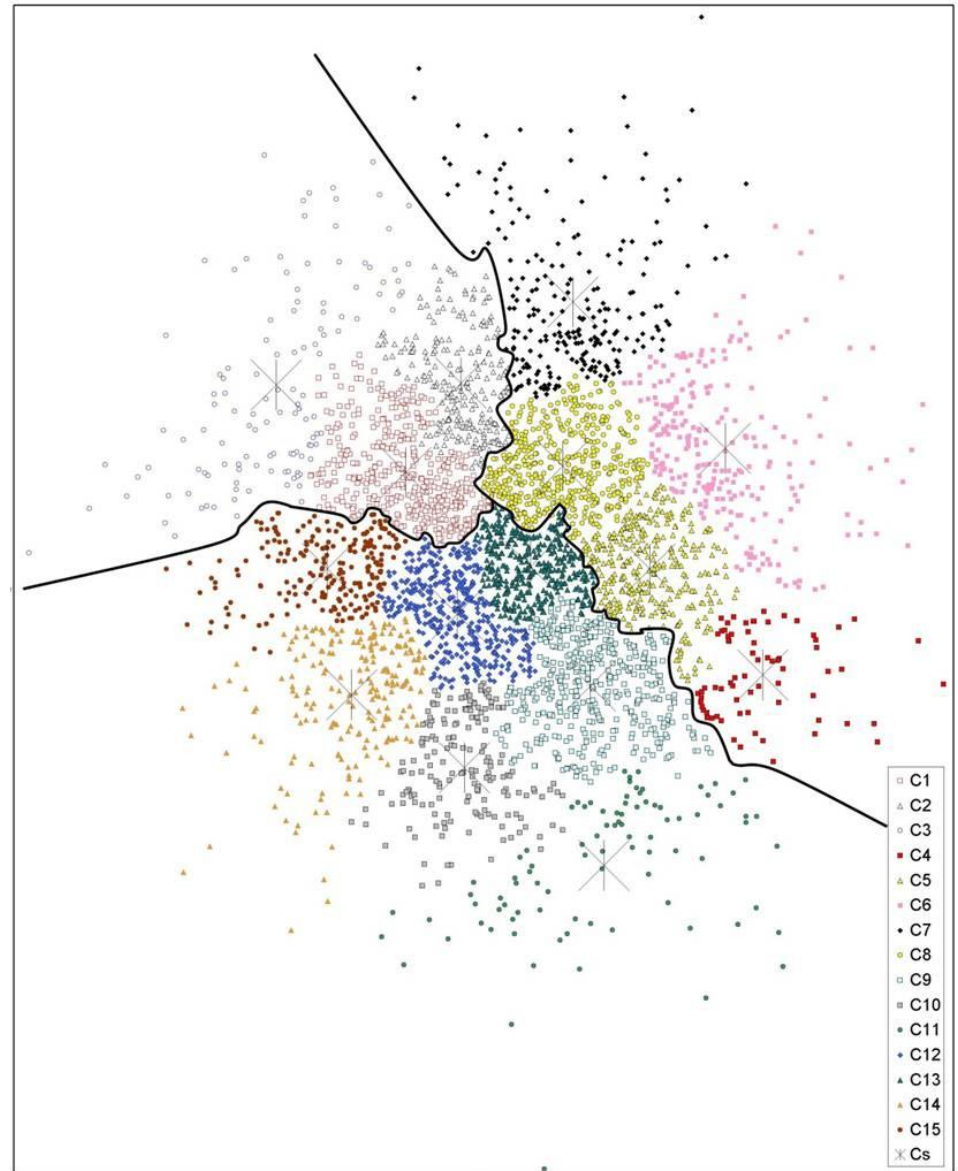(Criterion used for clustering: minimum sum of squares.)

# Hierarchical and partitional clustering

Scatterplot of two results of the hierarchical ***Ward*** method of bivariate Gaussian data without a cluster structure (same data as before). The solution for $K = 15$ clusters is marked by colour, the three cluster solution by freehand lines. Obviously, quite "different kinds" of clusters occur compared to the $K$ means clustering.

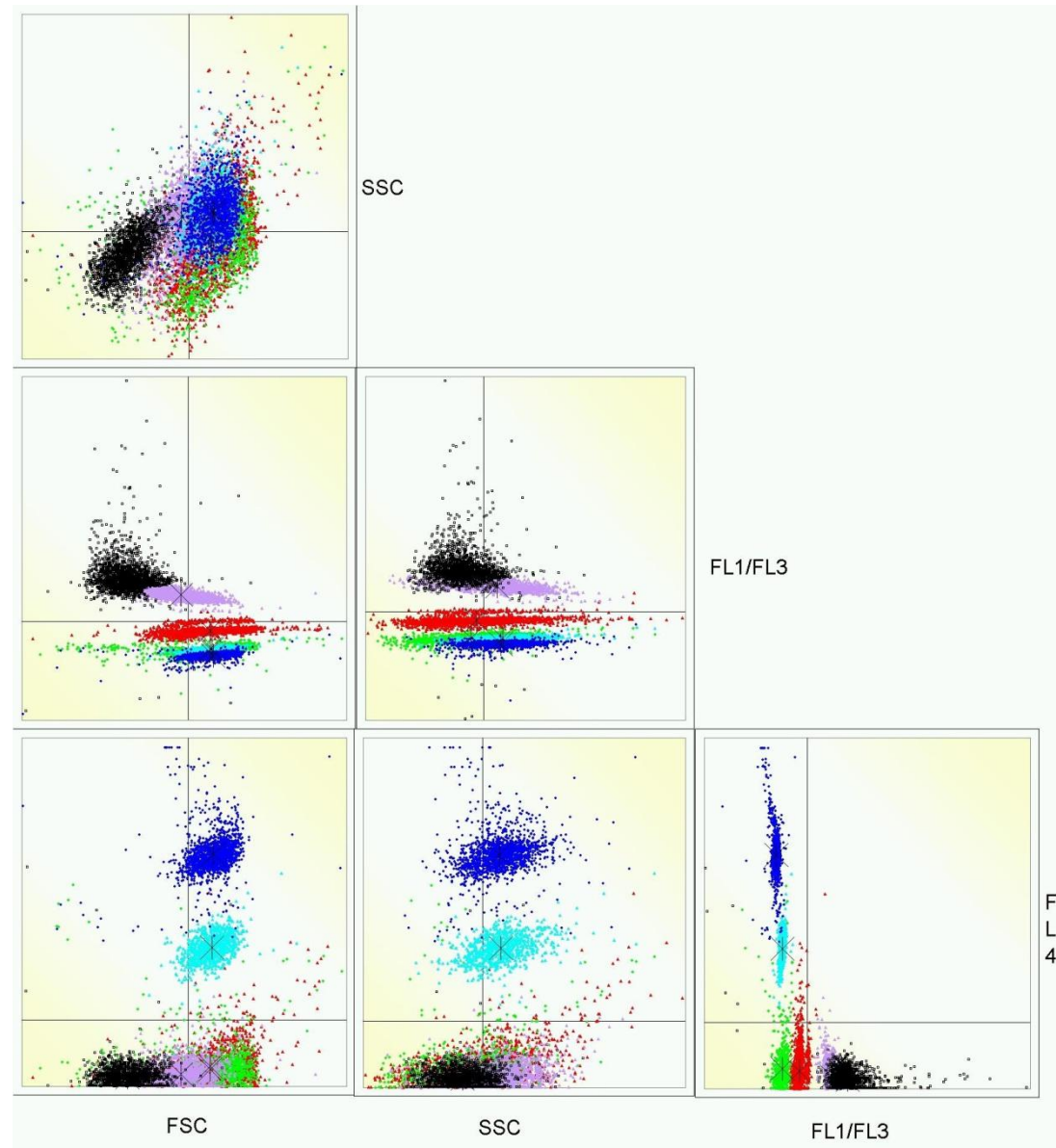(Same criterion as before: minimum sum of squares.)

# Multivariate projection

A scatterplot matrix can be recommended when the number of variables ranges from 3 to low numbers such as 20.
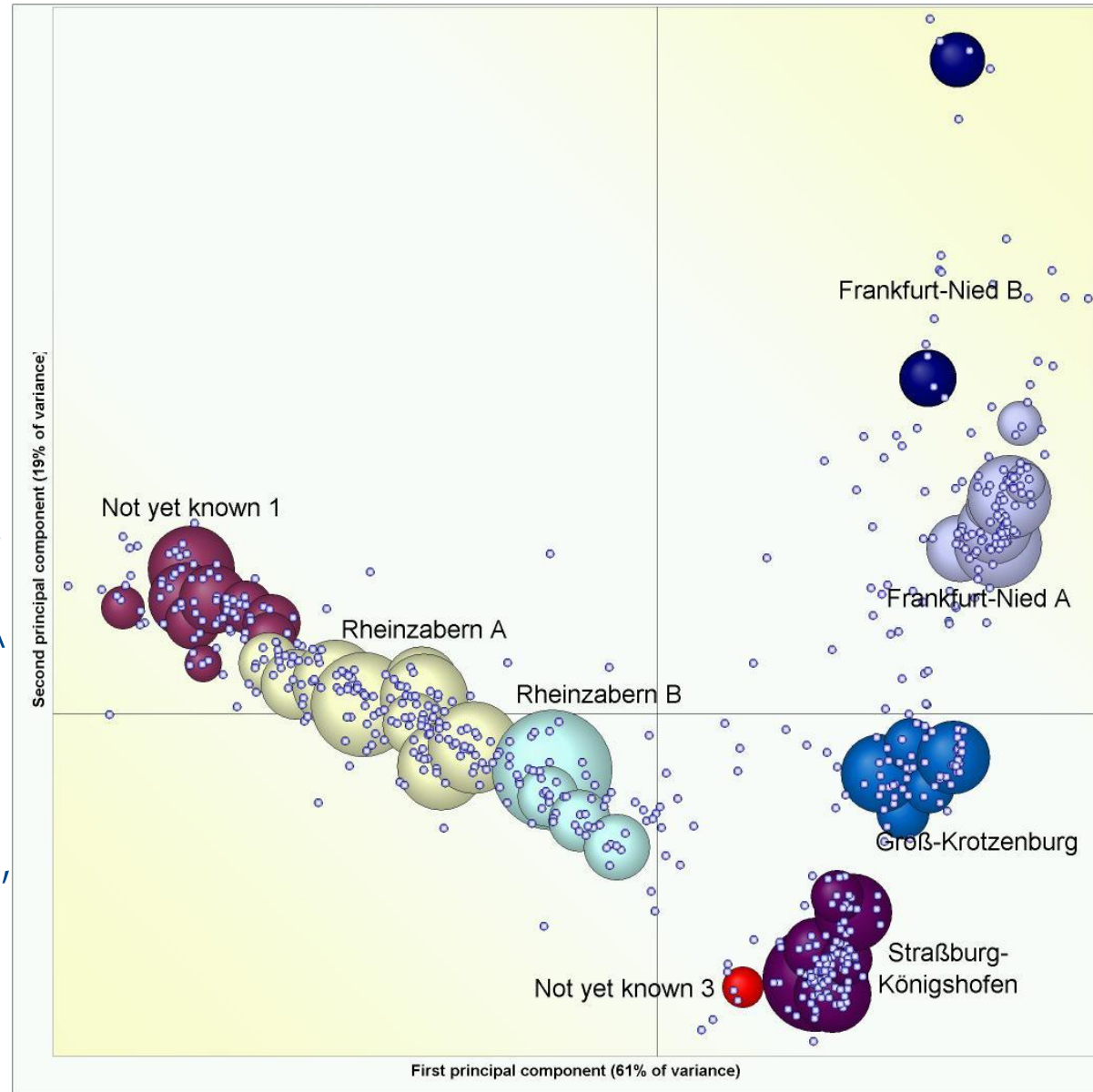
**Application to ecology**: Visualisation of particles (in water) based on measurements of their fluorescence properties. Their cluster membership is marked by colour (Mucha et al. 2002: WIAS Technical Report **5**).

# Multivariate projection

Principal components analysis (PCA) and linear discriminant analysis are well-known multivariate projection techniques. The latter requires in-formation about classes.
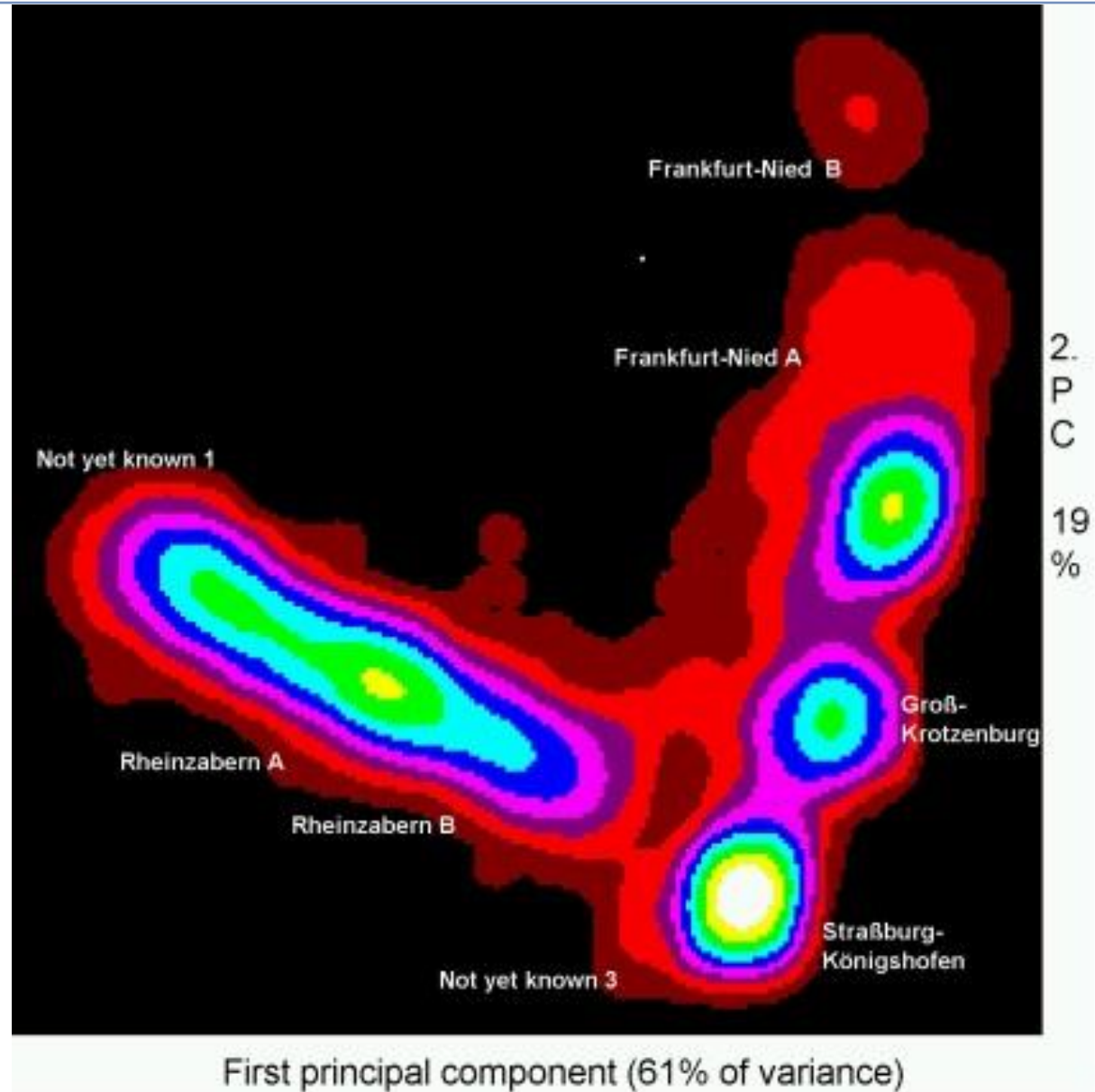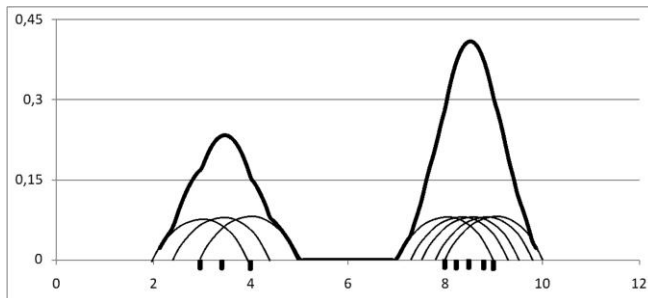
**Tiles:** Simultaneous PCA plot of observations (marked by points), dense regions (bubbles), and the final clusters (colour).
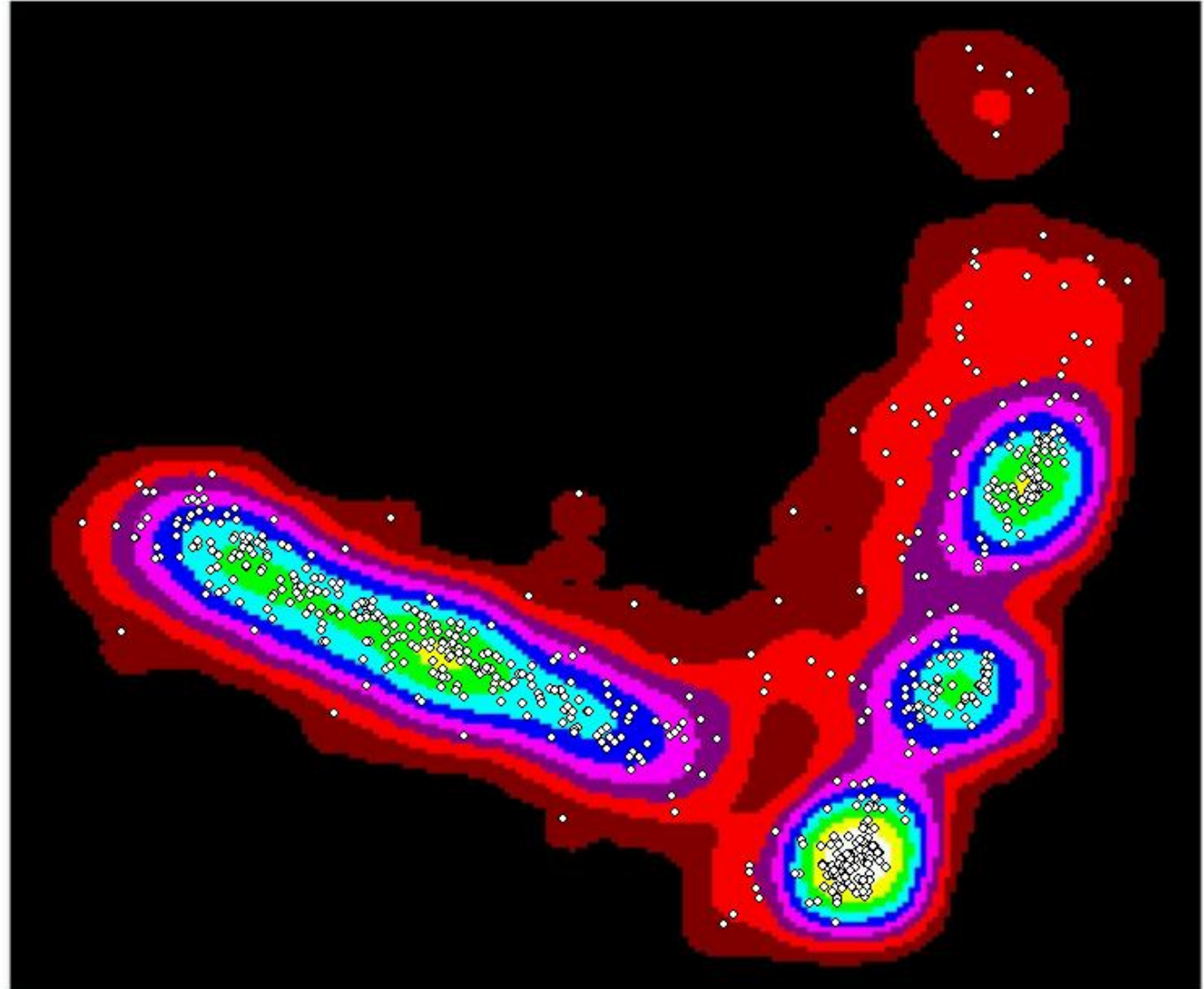
# Multivariate projection

**Tiles:** Several cuts of a non-parametric density estimate in the plane of the first two principal components (PC): this is another (smooth) view at the data.

The idea of density estimation (smoothing) is the superposition of elementary kernels:
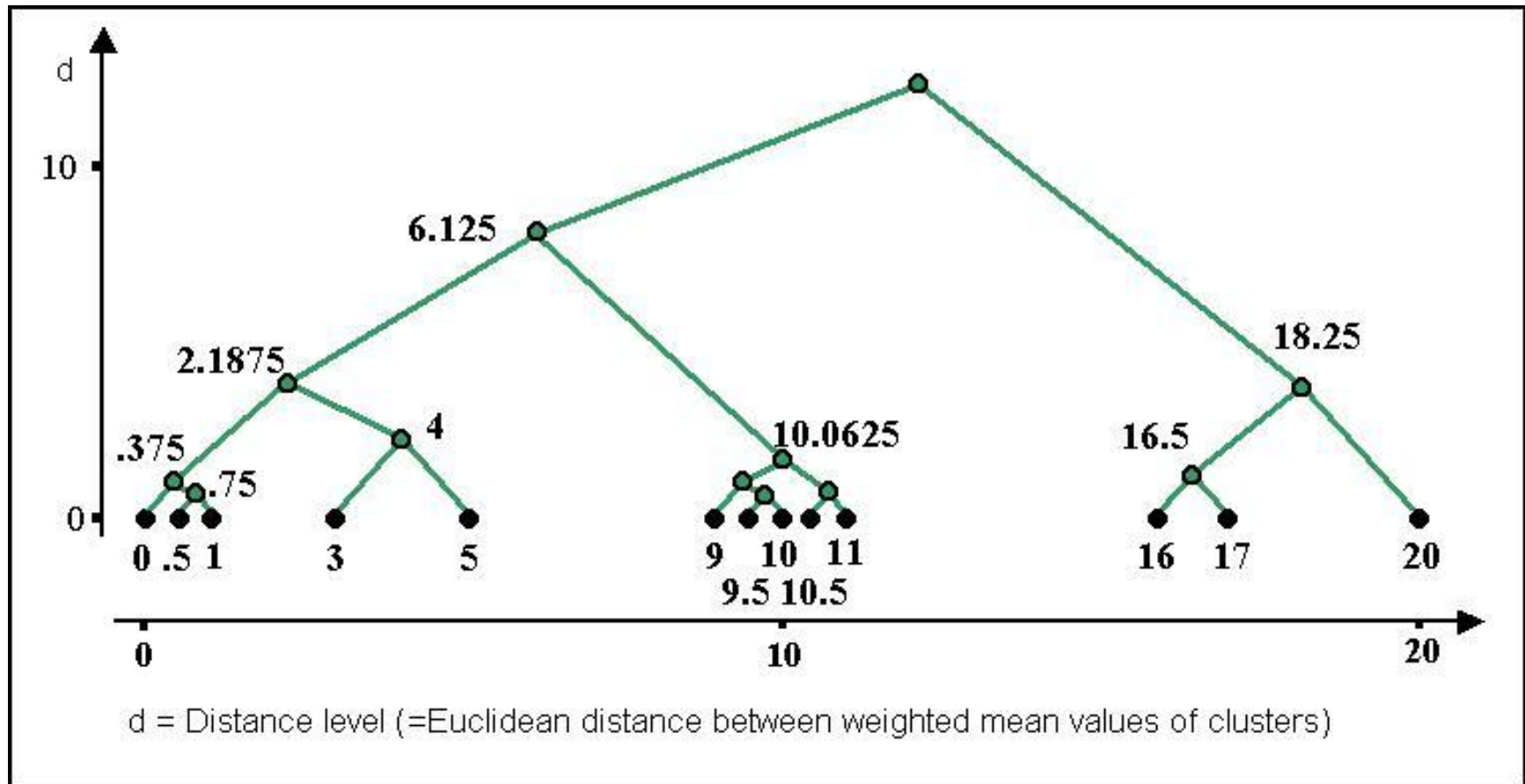




First principal component (61% of variance)

# Bivariate density estimation

**Tiles:** from individual observations (points) to a smooth bivariate non-parametric density. The latter was cut at different levels.
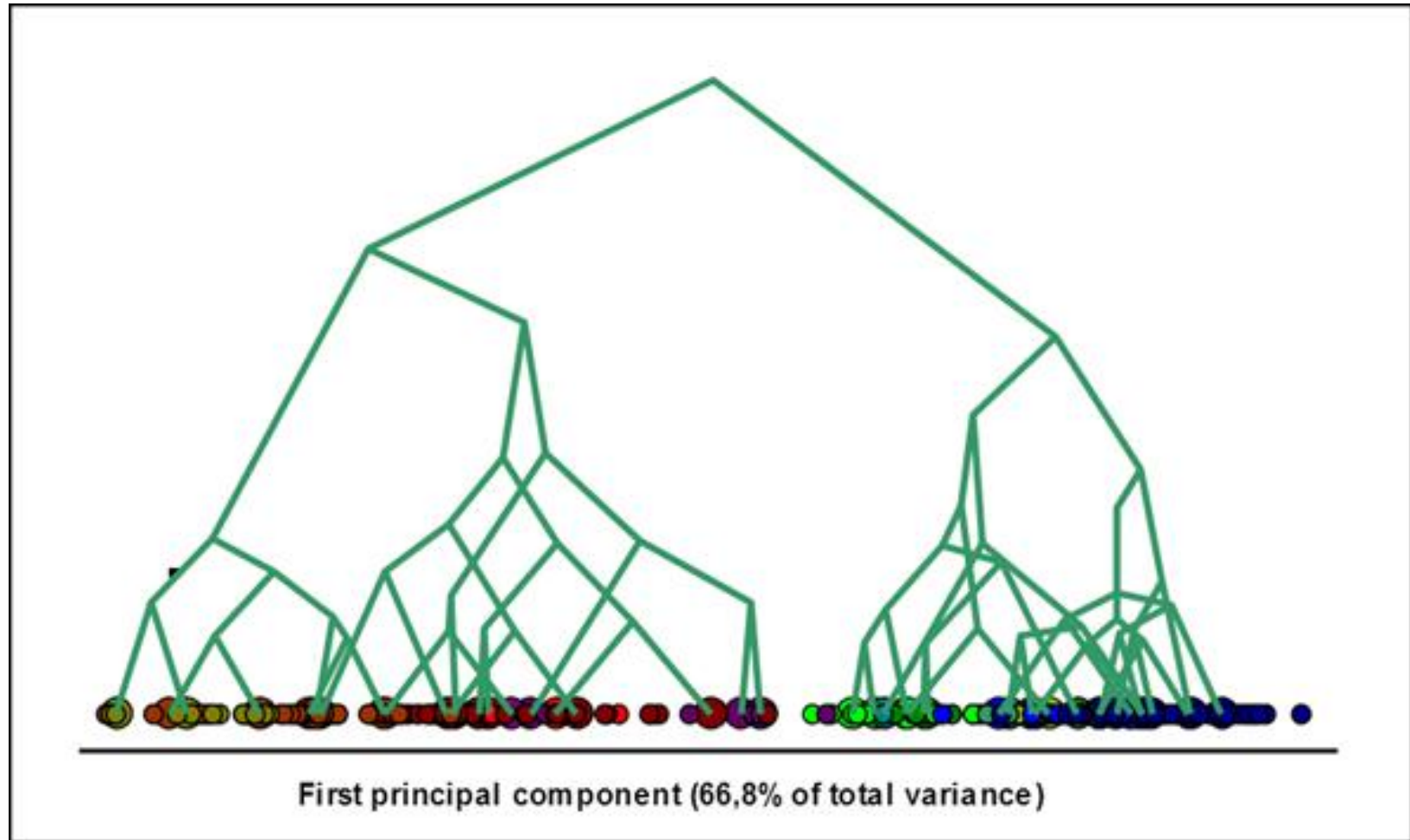
# Dendrograms

Dendrograms report about the process of hierarchical clustering. Example: Non-equidistant dendrogram of Centroid clustering of 20 data point {0, 0.5,1, 3, 5, …, 20} on the real line (abscissa).



d = Distance level (=Euclidean distance between weighted mean values of clusters)
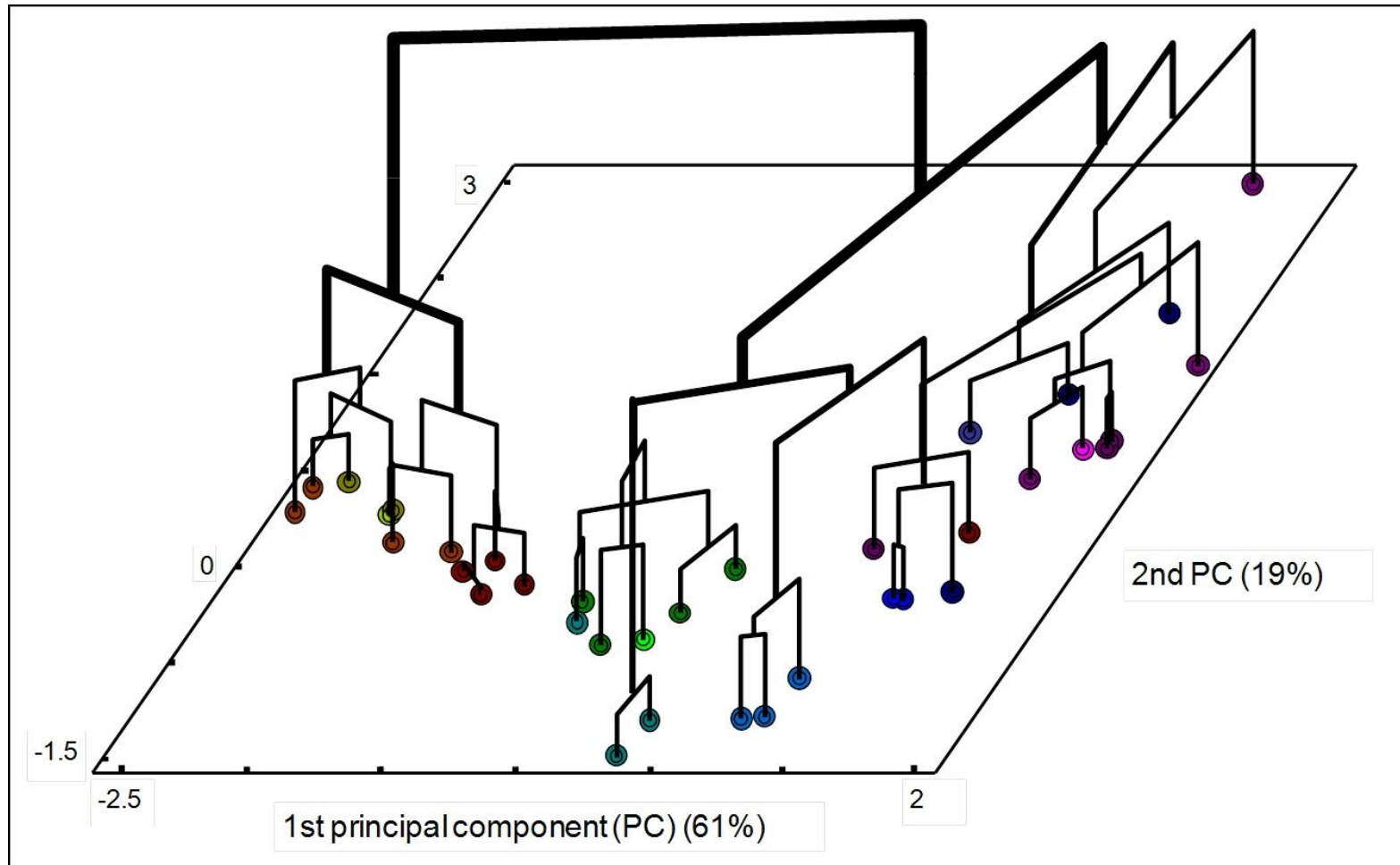
# Dendrograms

Example: Non-equidistant dendrogram of the hierarchical **Ward** clustering of 200 Swiss bank notes based on 6 measurements. The genuine bank notes on the right hand side look more homogeneous than the forged ones.



First principal component (66,8% of total variance)

# Dendrograms

Plot-dendrogram of **Ward** 's clustering of Roman tiles. The (reduced) dendrogram is projected onto the plane of the first two PC.
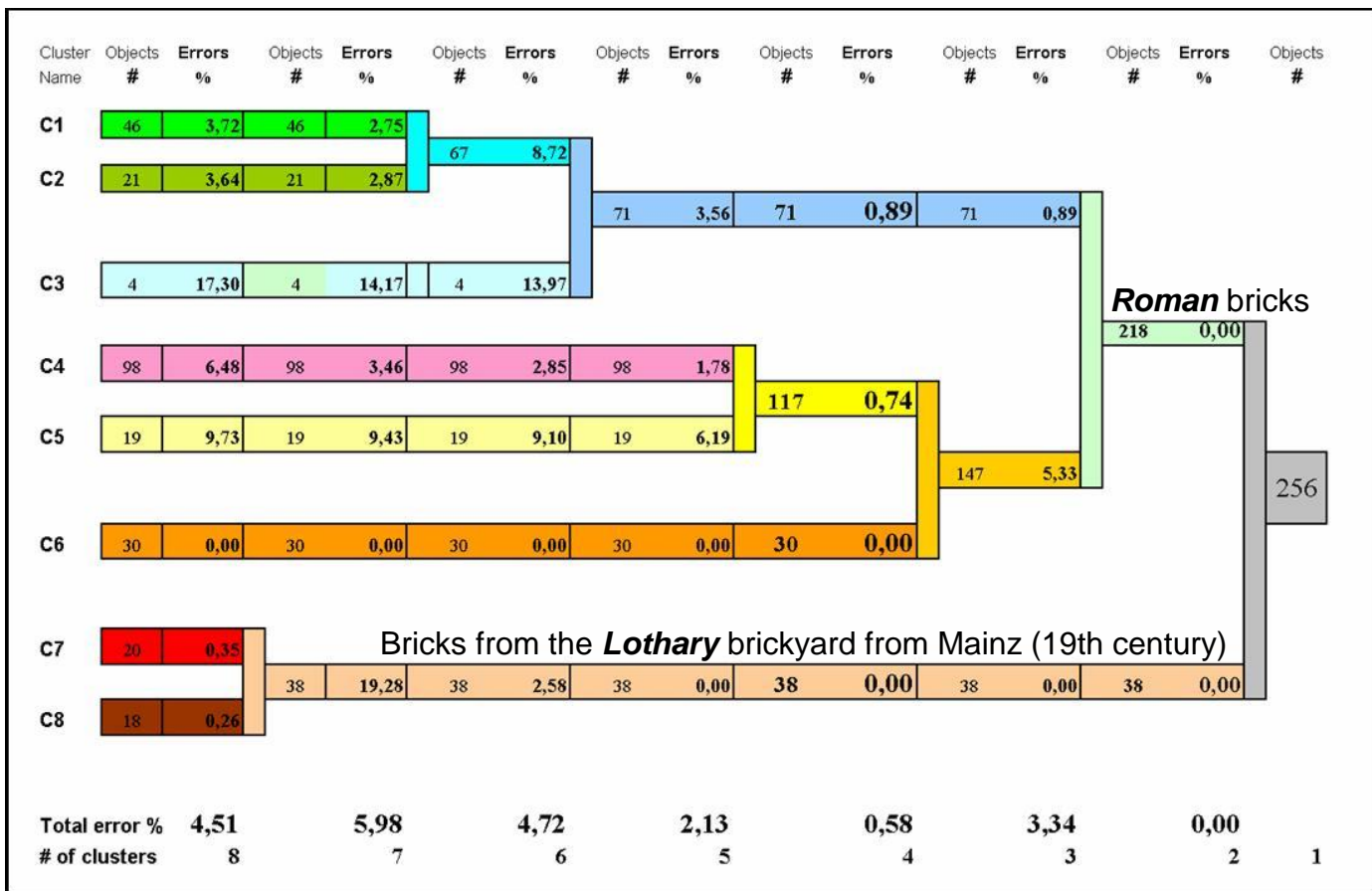
Informative dendrograms are ordered binary trees that show additional information such as *Jaccard*'s stability values. To define stability with respect to the individual clusters, the following measures of similarity between two clusters (sets) $E$ and $F$ are used: $\gamma(E,F) = \dfrac{|E \cap F|}{|E \cup F|}$ (Jaccard), $\tau(E,F) = \dfrac{2|E \cap F|}{|E| + |F|}$ (Dice), and $\eta(E,F) = \dfrac{|E \cap F|}{|E|}$ .

**Roman bricks versus modern ones.** Ward's hierarchical clustering clearly separates Roman bricks from modern ones. But, what about the stability of the clusters? In the dendrogram, the results of validation are presented using the "rate of recovery" $\eta$. ($\eta$ is related to the "errors" by $(1-\eta)*100\%$.) Here, 250 bootstrap clusterings were compared with the original clustering.

# Mapping findspots

Spatial and statistical analysis of findspots of Roman military brickstamps in Mainz. The background:

- The findings come from the first four centuries A.D. The Roman bricks and tiles have been classified based on their stamp. In addition, new types of stamps have been defined.
- The database itself is under development yet: many hundreds of findings are not yet identified by their coordinates of location.
- Mainz has been a great urban area for thousands of years. A planned excavation is not possible. Therefore most of the places of discovery are found by "accidents" such as the excavation around the main railway station.
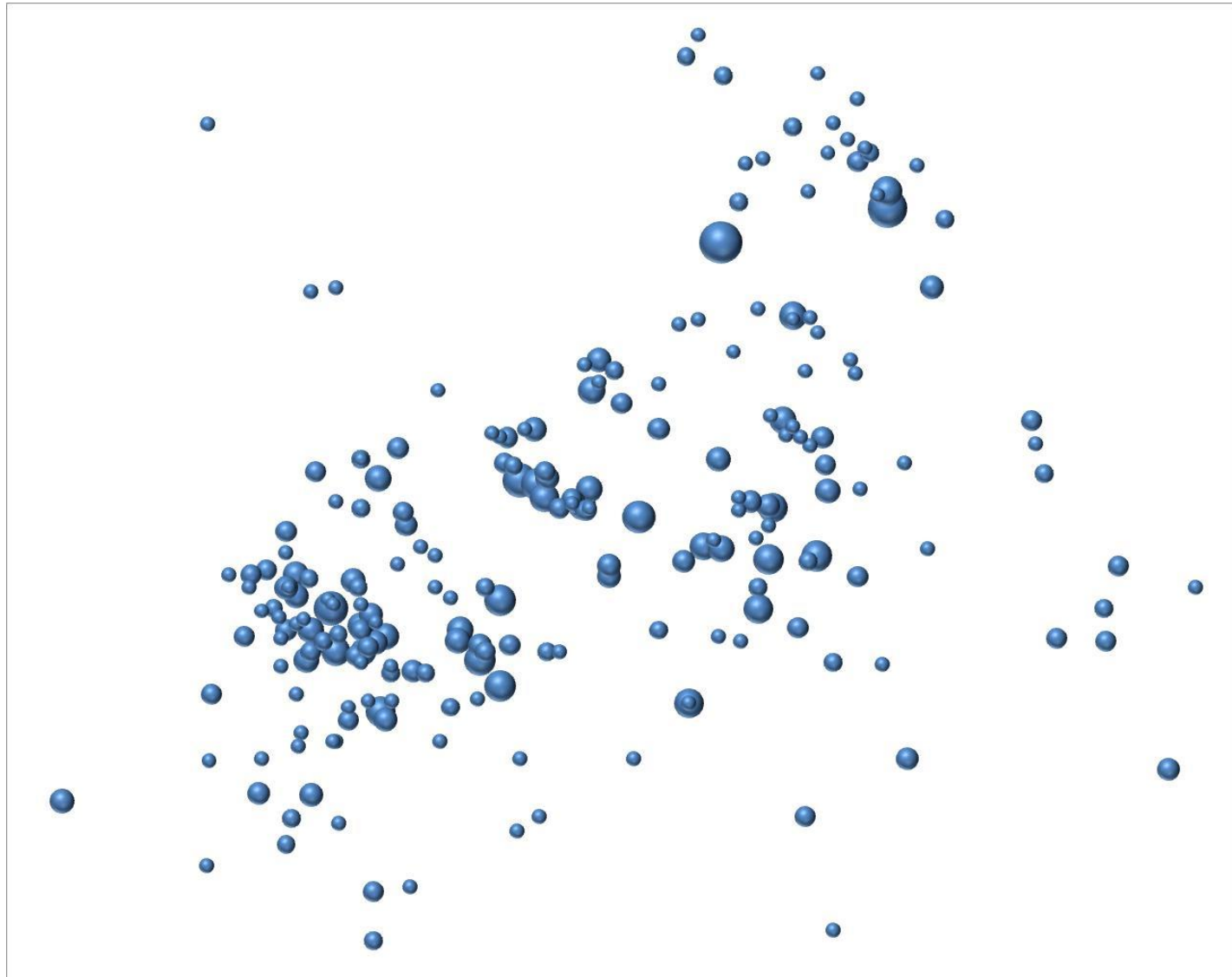
# Mapping findspots

Example of a mapping of a zoomed area of Mainz (city) with locations of findspots.
The river Rhine is located in the right upper corner. All in all there are 246 sites until now.
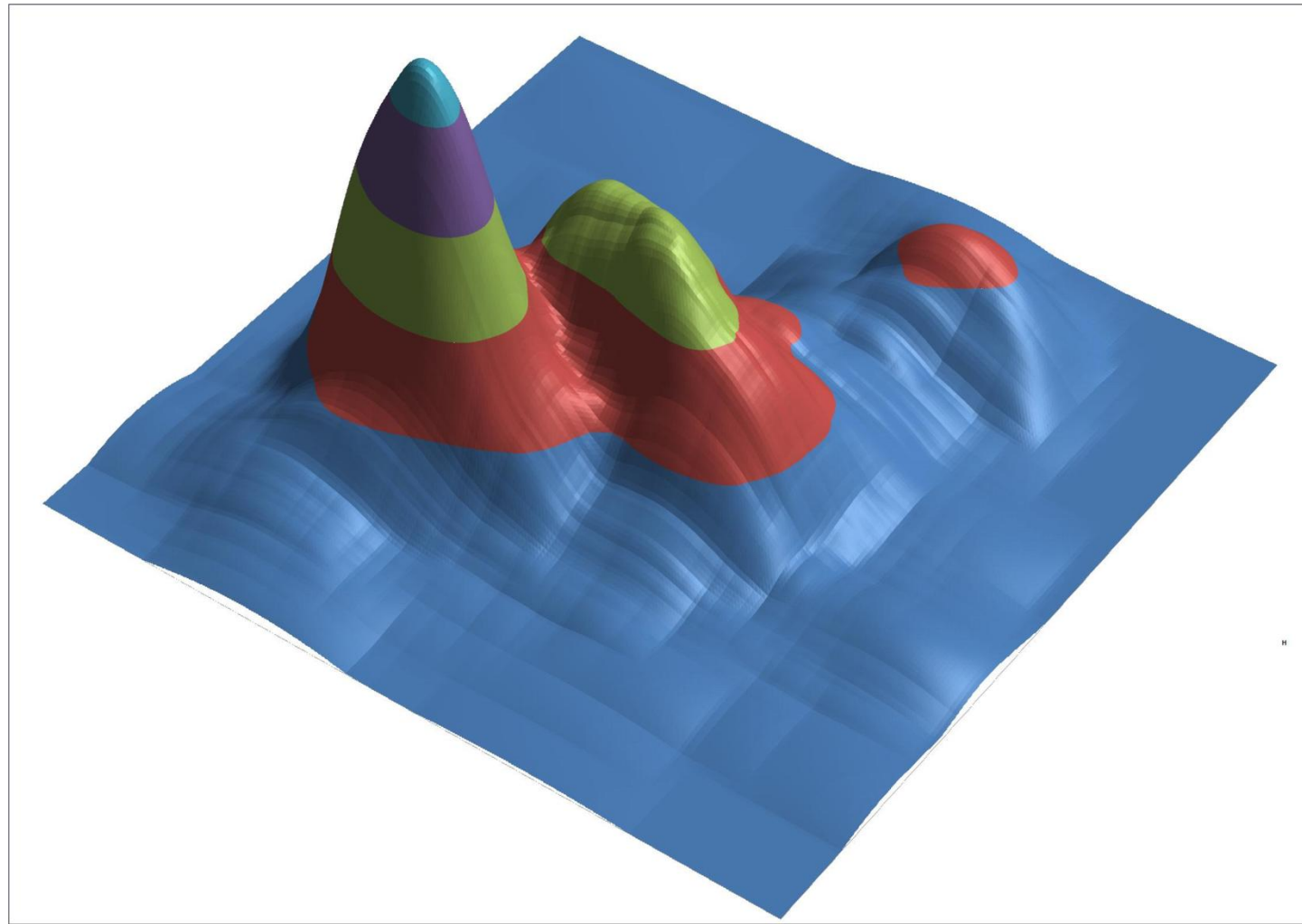
# Mapping findspots

The locations of findspots are marked by bubbles. The size of a bubble is proportional to the number of findings at the location. So the plot becomes more informative.
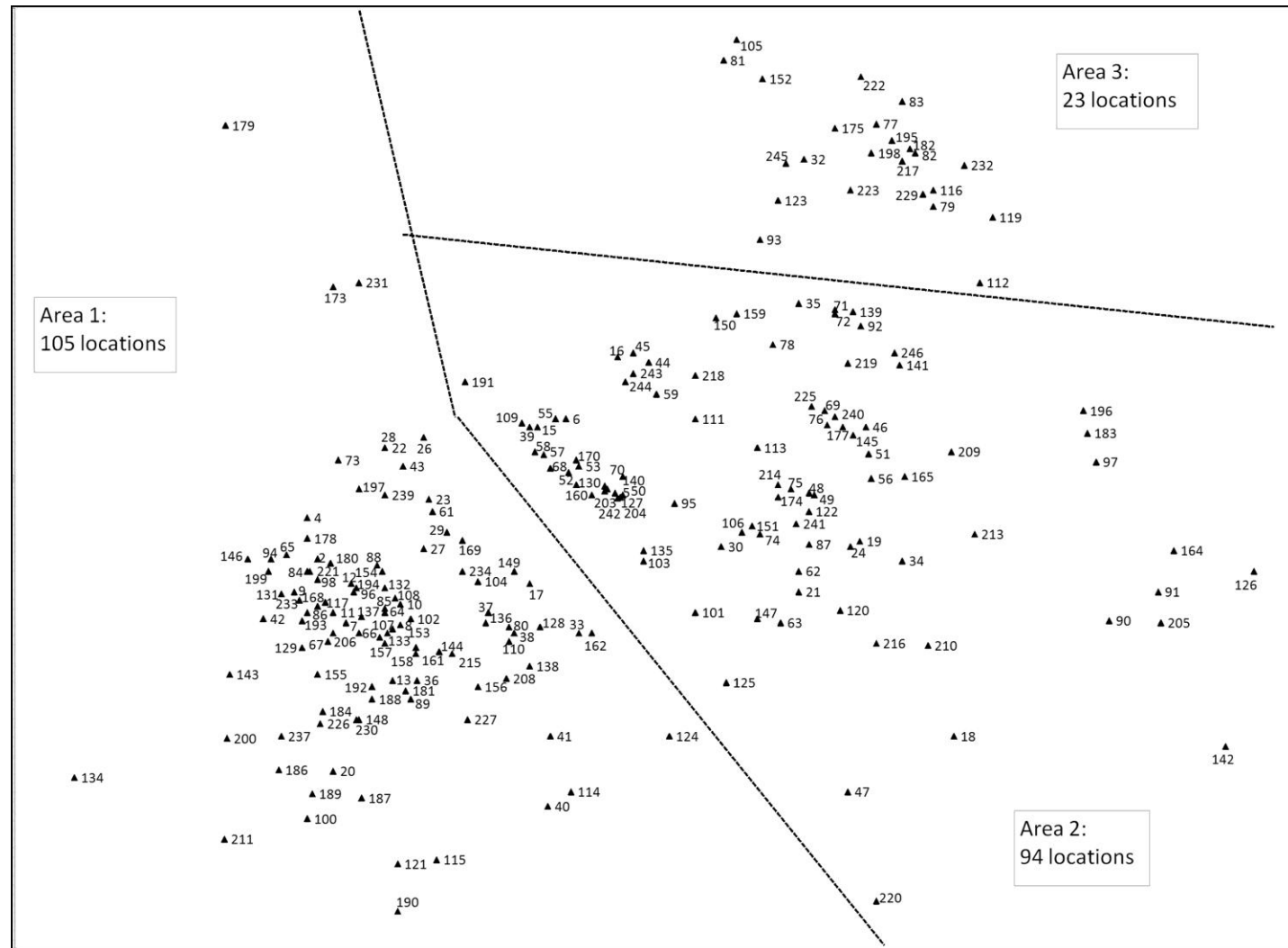
Bivariate density estimation of all locations. Herein they are weighted by the logarithm of the number of findings plus 1.

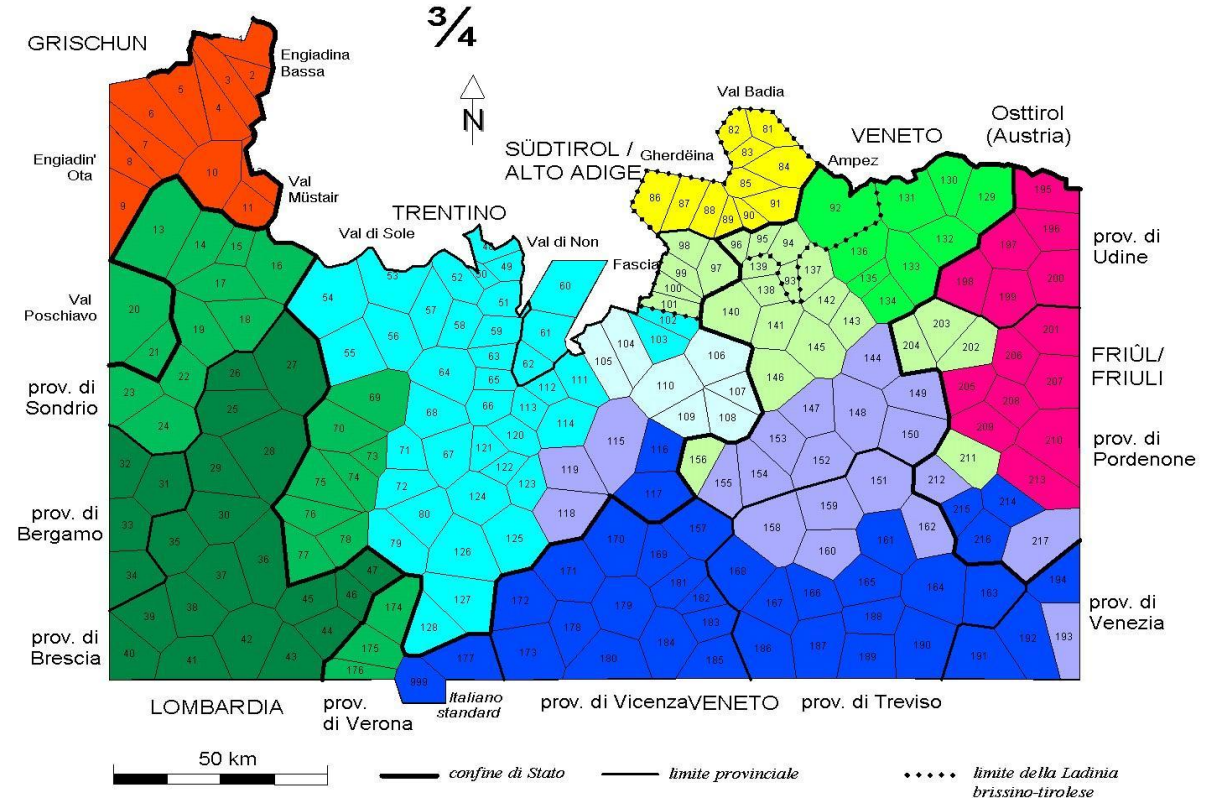Archaeological dissection of the ancient Mainz into the three main spatial regions *Area 1, Area 2,* and *Area 3.* It turns out that these regions coincide with three time periods significantly.

The visualisation of the reliability of the cluster membership of each individual observation is of special interest in the field of geographical mapping. The reliability can be assessed based on the investigation of stability of individual clusters. The figure shows an application to quantitative linguistics where significant regions (95% level) in North Italy are in dark colour (Mucha and Haimerl 2005).

# Summary

Visualisation is essential for a better understanding of data analysis and cluster analysis.

*Thank you for your attention!*

Bibliography

Mucha, H.-J., Bartel, H.-G. and Dolata, J. (2005): Techniques of rearrangements in binary trees (dendrograms) and applications. MATCH Commun. Math. Comput. Chem., 54, pp. 561--582.

Mucha, H.-J. and Ritter, G. (2009): Classification and clustering: Models, software and applications. WIAS Report No. 26, WIAS, Berlin.