

# Introduction to **Multivariate Statistical Analysis in Chemometrics**

Kurt Varmuza  
Peter Filzmoser



 **CRC Press**  
Taylor & Francis Group

## **Authors**

Kurt Varmuza  
Peter Filzmoser

Vienna University of Technology,  
Austria

[kurt.varmuza@tuwien.ac.at](mailto:kurt.varmuza@tuwien.ac.at)  
[peter.filzmoser@tuwien.ac.at](mailto:peter.filzmoser@tuwien.ac.at)  
Info: [www.lcm.tuwien.ac.at](http://www.lcm.tuwien.ac.at)

ISBN: 9781420059472

CRC Press (Taylor & Francis),  
Boca Raton, FL, USA, [www.crcpress.com](http://www.crcpress.com)

2009

321 + xiii pages

Ca US\$ 180

(220129)

- Includes important multivariate statistical methods, such as principal component analysis and support vector machines, for analyzing scientific data
- Explains the methods using formulae, graphical illustrations, and schemes
- Demonstrates **R software tools** with fully worked-out, real-world examples
- Emphasizes the use of **robust statistical methods**
- Offers practical advice on applying the methods

Using formal descriptions, graphical illustrations, practical examples, and R software tools,

### **Introduction to Multivariate Statistical Analysis in Chemometrics**

presents simple yet thorough explanations of the most important multivariate statistical methods for analyzing chemical data.

It includes discussions of various statistical methods, such as principal component analysis, regression analysis, classification methods, and clustering.

Written by a chemometrician and a statistician, the book reflects both the practical approach of chemometrics and the more formally oriented one of statistics. To enable a better understanding of the statistical methods, the authors apply them to real data examples from chemistry. They also examine results of the different methods, comparing traditional approaches with their robust counterparts. In addition, the authors use the freely available R package to implement methods, encouraging readers to go through the examples and adapt the procedures to their own problems.

Focusing on the practicality of the methods and the validity of the results, this book offers concise mathematical descriptions of many multivariate methods and employs graphical schemes to visualize key concepts. It effectively imparts a basic understanding of how to apply statistical methods to multivariate scientific data.

# Contents

- 1 Introduction**
  - 1.1 Chemoinformatics - chemometrics - statistics**
  - 1.2 This book**
  - 1.3 Historical remarks about chemometrics**
  - 1.4 Bibliography**
  - 1.5 Starting examples**
    - 1.5.1 Univariate versus bivariate classification
    - 1.5.2 Nitrogen content of cereals computed from NIR data
    - 1.5.3 Elemental composition of aecheaeological glasses
  - 1.6 Univariate statistics - a reminder**
    - 1.6.1 Empirical distributions
    - 1.6.2 Theoretical distributions
    - 1.6.3 Central value
    - 1.6.4 Spread
    - 1.6.5 Statistical tests
- 2 Multivariate data**
  - 2.1 Definitions**
  - 2.2 Basic preprocessing**
    - 2.2.1 Data transformation
    - 2.2.2 Centering and scaling
    - 2.2.3 Normalization
    - 2.2.4 Transformations for compositional data
  - 2.3 Covariance and correlation**
    - 2.3.1 Overview
    - 2.3.2 Estimating covariance and correlation
  - 2.4 Distances and similarities**
  - 2.5 Multivariate outlier identification**
  - 2.6 Linear latent variables**
    - 2.6.1 Overview
    - 2.6.2 Projection and mapping
    - 2.6.3 Example
  - 2.7 Summary**
- 3 Principal component analysis**
  - 3.1 Concepts**
  - 3.2 Number of PCA components**
  - 3.3 Centering and scaling**
  - 3.4 Outliers and data distribution**
  - 3.5 Robust PCA**
  - 3.6 Algorithms for PCA**
    - 3.6.1 Mathematics of PCA
    - 3.6.2 Jacobi rotation
    - 3.6.3 Singular value decomposition
    - 3.6.4 NIPALS
  - 3.7 Evaluation and diagnostics**
    - 3.7.1 Cross validation for determination of the number of principal components
    - 3.7.2 Explained variance for each variable
    - 3.7.3 Diagnostic plots
  - 3.8 Complementary methods for exploratory data analysis**
    - 3.8.1 Factor analysis
    - 3.8.2 Cluster analysis and dendrogram
    - 3.8.3 Kohonen mapping

3.8.4 Sammon's nonlinear mapping

3.8.5 Multi-way PCA

### **3.9 Examples**

3.9.1 Tissue samples from human mummies and fatty acid concentrations

3.9.2 Polycyclic aromatic hydrocarbons in aerosol

### **3.10 Summary**

## **4 Calibration**

### **4.1 Concepts**

### **4.2 Performance of regression models**

4.2.1 Overview

4.2.2 Overfitting and underfitting

4.2.3 Performance criteria

4.2.4 Criteria for models with different numbers of variables

4.2.5 Cross validation

4.2.6 Bootstrap

### **4.3 Ordinary least-squares regression**

4.3.1 Simple OLS

4.3.2 Multiple OLS

4.3.3 Multivariate OLS

### **4.4 Robust regression**

4.4.1 Overview

4.4.2 Regression diagnostics

4.4.3 Practical hints

### **4.5 Variable selection**

4.5.1 Overview

4.5.2 Univariate and bivariate selection methods

4.5.3 Stepwise selection methods

4.5.4 Best-subset regression

4.5.5 Variable selection based on PCA or PLS models

4.5.6 Genetic algorithms

4.5.7 Cluster analysis of variables

4.5.8 Example

### **4.6 Principal component regression**

4.6.1 Overview

4.6.2 Number of PCA components

### **4.7 Partial least-squares regression**

4.7.1 Overview

4.7.2 Mathematical aspects

4.7.3 Kernel algorithm for PLS

4.7.4 NIPALS algorithm for PLS

4.7.5 SIMPLS algorithm for PLS

4.7.6 Other algorithms for PLS

4.7.7 Robust PLS

### **4.8 Related methods**

4.8.1 Canonical correlation analysis

4.8.2 Ridge and Lasso regression

4.8.3 Nonlinear regression

### **4.9 Examples**

4.9.1 GC retention indices of polycyclic aromatic compounds

4.9.2 Cereal data

### **4.10 Summary**

## **5 Classification**

### **5.1 Concepts**

### **5.2 Linear classification methods**

5.2.1 Linear discriminant analysis

5.2.2 Linear regression for discriminant analysis

	5.2.3	Logistic regression
<b>5.3</b>		<b>Kernel and prototype methods</b>
	5.3.1	SIMCA
	5.3.2	Gaussian mixture models
	5.3.3	kNN - classification
<b>5.4</b>		<b>Classification trees</b>
<b>5.5</b>		<b>Artificial neural networks</b>
<b>5.6</b>		<b>Support vector machine</b>
<b>5.7</b>		<b>Evaluation</b>
	5.7.1	Principles and misclassification error
	5.7.2	Predictive ability
	5.7.3	Confidence in classification answers
<b>5.8</b>		<b>Examples</b>
	5.8.1	Origin of glass samples
	5.8.2	Recognition of chemical substructures from mass spectra
<b>5.9</b>		<b>Summary</b>
<b>6</b>		<b>Cluster analysis</b>
<b>6.1</b>		<b>Concepts</b>
<b>6.2</b>		<b>Distance and similarity measures</b>
<b>6.3</b>		<b>Partitioning methods</b>
<b>6.4</b>		<b>Hierarchical clustering methods</b>
<b>6.5</b>		<b>Fuzzy clustering</b>
<b>6.6</b>		<b>Model-based clustering</b>
<b>6.7</b>		<b>Cluster validity and clustering tendency measures</b>
<b>6.8</b>		<b>Examples</b>
	6.8.1	Chemotaxonomy of plants
	6.8.2	Glass samples
<b>6.9</b>		<b>Summary</b>

<b>7</b>		<b>Preprocessing</b>
<b>7.1</b>		<b>Concepts</b>
<b>7.2</b>		<b>Smoothing and differentiation</b>
<b>7.3</b>		<b>Multiplicative signal correction</b>
<b>7.4</b>		<b>Mass spectral features</b>

## **Appendix 1      Symbols and abbreviations**

## **Appendix 2      Matrix algebra**

A.2.1	Definitions
A.2.2	Addition and subtraction of matrices
A.2.3	Multiplication of vectors
A.2.4	Multiplication of matrices
A.2.5	Matrix inversion
A.2.6	Eigenvectors
A.2.7	Singular value decomposition

## **Appendix 3      Introduction to $\mathbb{R}$**

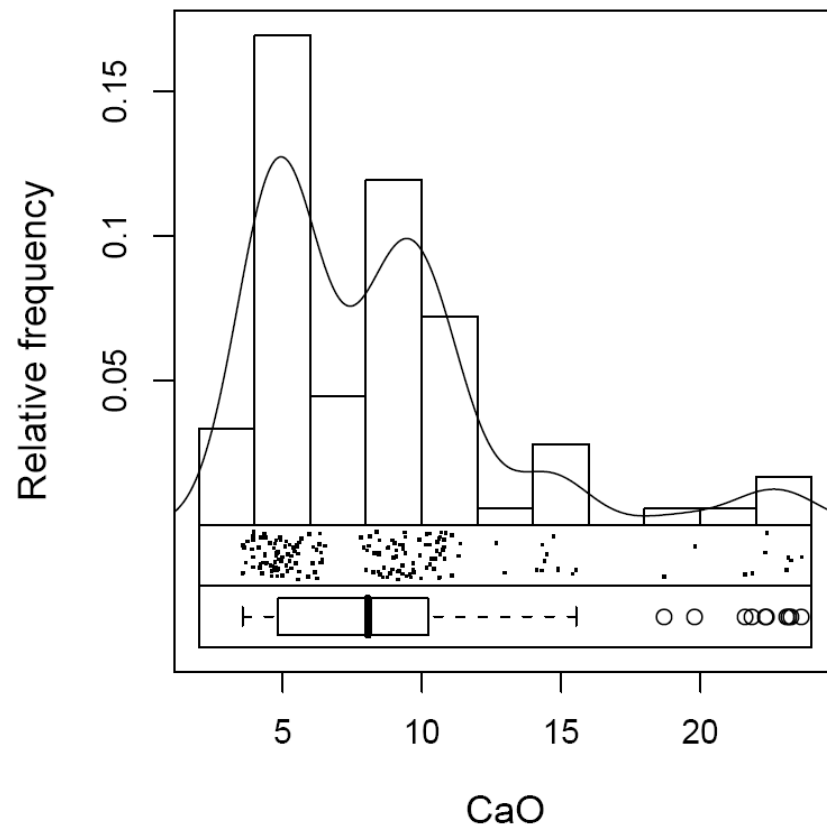
A.3.1	General information on $\mathbb{R}$
A.3.2	Installing $\mathbb{R}$
A.3.3	Starting $\mathbb{R}$
A.3.4	Working directory
A.3.5	Loading and saving data
A.3.6	Important $\mathbb{R}$ functions
A.3.7	Operators and basic functions
A.3.8	Data types
A.3.9	Data structures
A.3.10	Selection and extraction from data objects
A.3.11	Generating and saving graphics

## **Index**

## Examples from the book

### 1. Introduction

#### 1.6. Univariate Statistics - A Reminder



The BOXPLOT is an informative graphics to display a data distribution, based on median and quartiles. The boxplot can be defined as follows [Frank and Todeschini 1994]. The height of the box is given by the first and third quartile, and the mid line shows the median; the width of the box has usually no meaning. One whisker extends from the first quartile to the smallest data value in the interval  $Q_1$  to  $Q_1 - 1.5 \cdot IQR$  and is called the lower whisker. The other whisker extends from the third quartile to the largest data value in the interval  $Q_3$  to  $Q_3 + 1.5 \cdot IQR$  and is called the upper whisker. Outliers - not within the range  $[Q_1 - 1.5 \cdot IQR, Q_3 + 1.5 \cdot IQR]$  - are plotted as individual points.

Figure 1.7. The EDAPLOT (Exploratory Data Analysis plot) of data combines one-dimensional scatter plot, histogram, probability density trace, and boxplot. Data used are CaO concentrations (%) of 180 archaeological glass vessels [Janssen *et al.* 1998].

## Examples from the book

### 2. Multivariate Data

#### 2.6. Linear Latent Variables

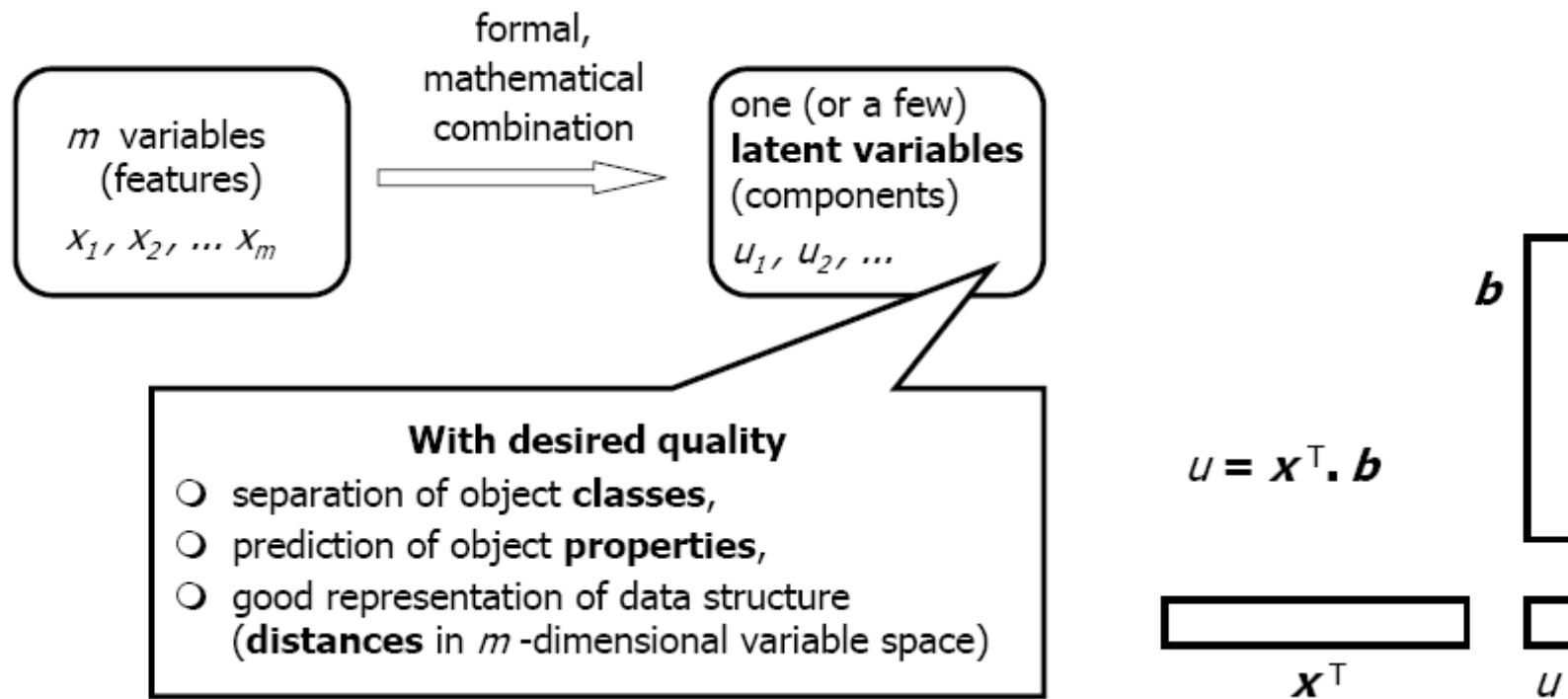


Figure 2.14. General concept of a latent variable (left), and calculation of the score  $u$  of a linear latent variable from the transposed variable vector  $\mathbf{x}^T$  and the loading vector  $\mathbf{b}$  as a scalar product (right).

## Examples from the book

### 3. Principal Component Analysis

#### 3.1. Concepts

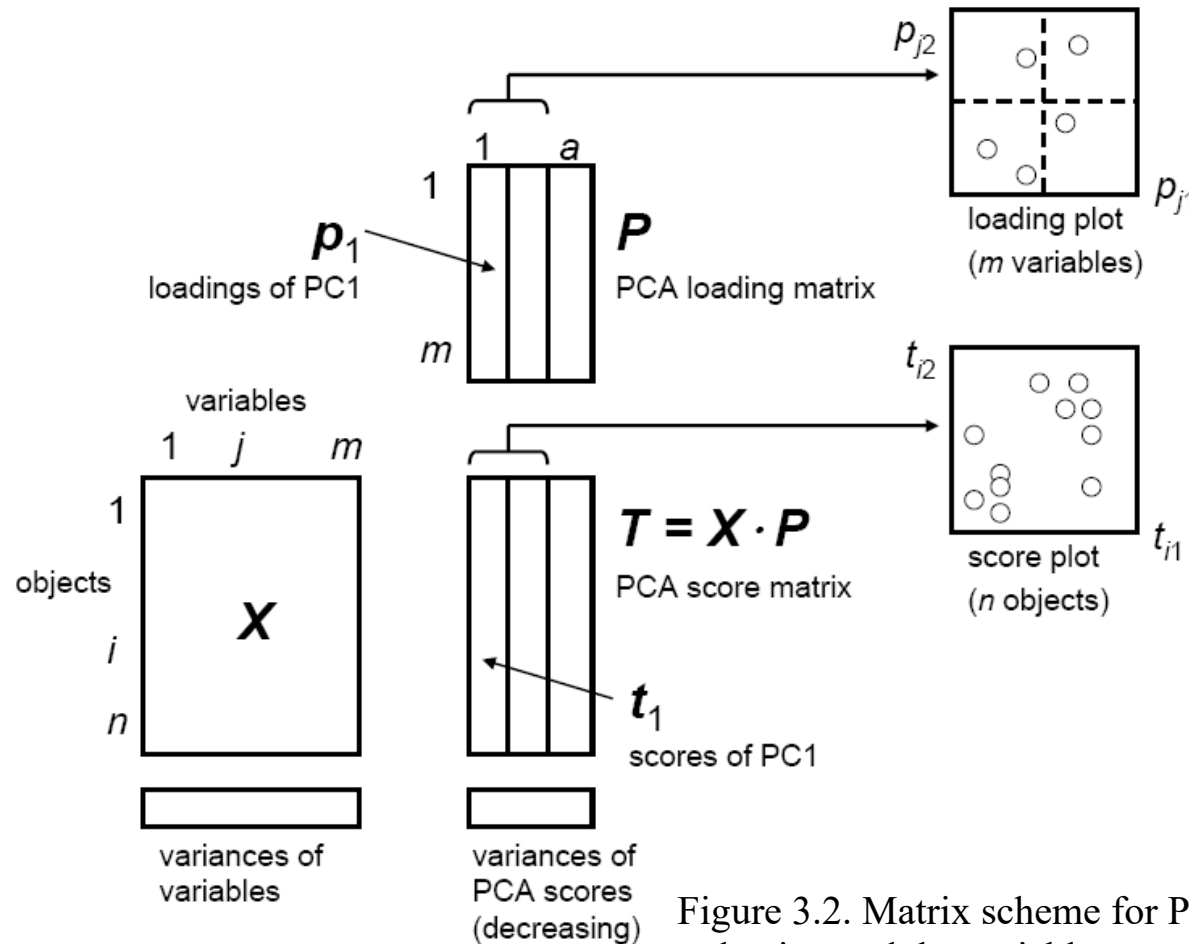


Figure 3.2. Matrix scheme for PCA. Since the aim of PCA is dimension reduction and the variables are often highly correlated  $a \leq \min(n, m)$  principal components are used.

Principal component analysis (PCA) can be considered as "the mother of all methods in multivariate data analysis". The aim of PCA is dimension reduction and PCA is the most frequently applied method for computing linear latent variables (components). PCA can be seen as a method to compute a new coordinate system formed by the latent variables, which is orthogonal, and where only the most informative dimensions are used. Latent variables from PCA optimally represent the distances between the objects in the high-dimensional variable space - remember, the distance of objects is considered as an inverse similarity of the objects.



## Examples from the book

### 4. Calibration

#### 4.9. Examples

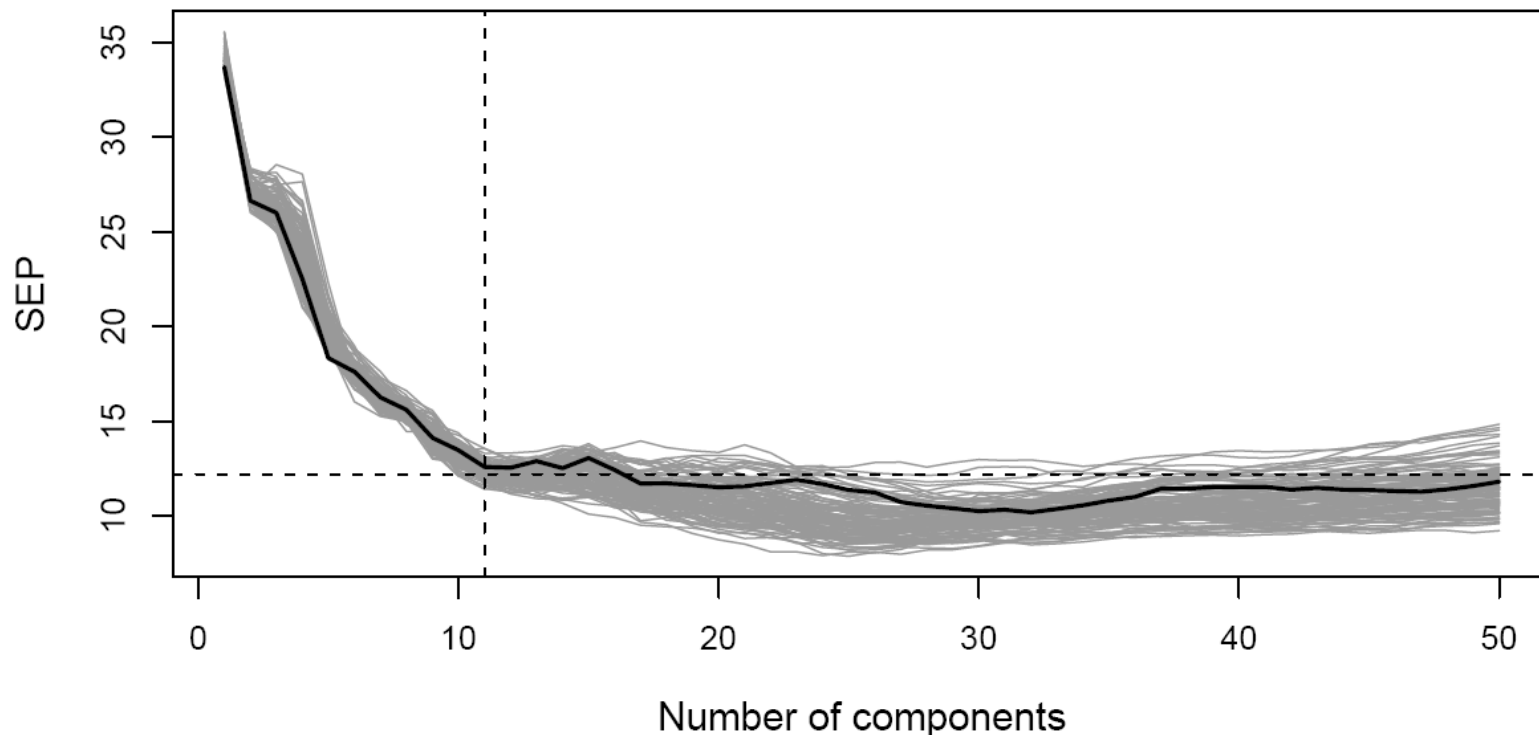


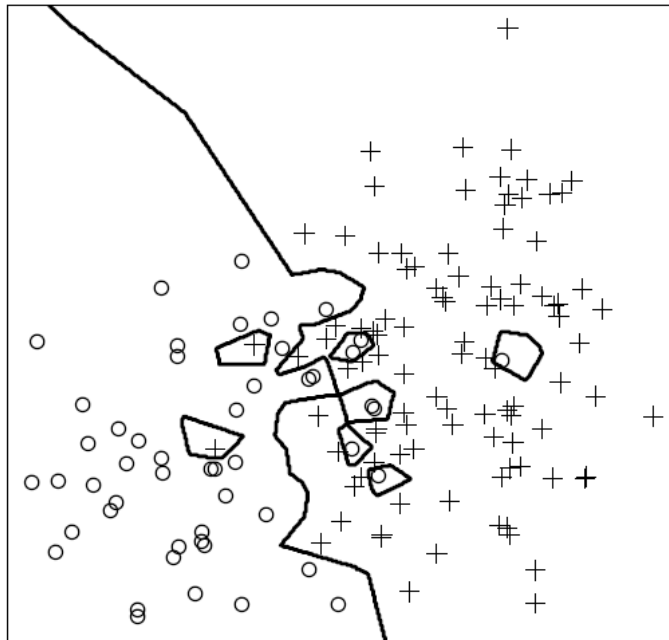
Figure 4.34. Results of PLS for the PAC data set. The black line results from a single 10-fold cross validation, the gray lines from repeating the 10-fold cross validation 100 times. The choice of the optimal number of principal components is based on repeated double cross validation.

## Examples from the book

### 5. Classification

#### 5.3.3. k-NN Classification

kNN classification for  $k=1$



kNN classification for  $k=15$

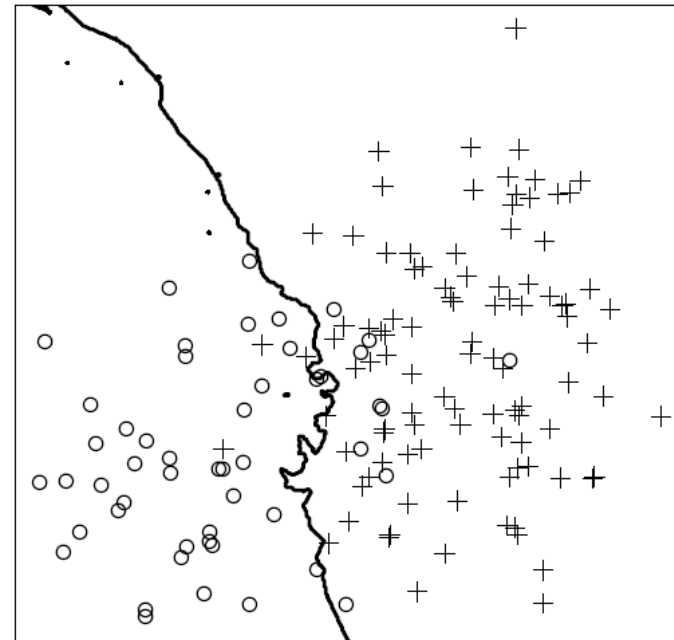


Figure 5.13. kNN classification for two groups of 2-dimensional data. The training data are shown with the symbol corresponding to the group membership. Any new data point would be classified according to the presented decision boundaries, where  $k = 1$  in the left plot, and  $k = 15$  in the right plot has been used.

## Examples from the book

### 6. Cluster Analysis

#### 6.1. Concepts

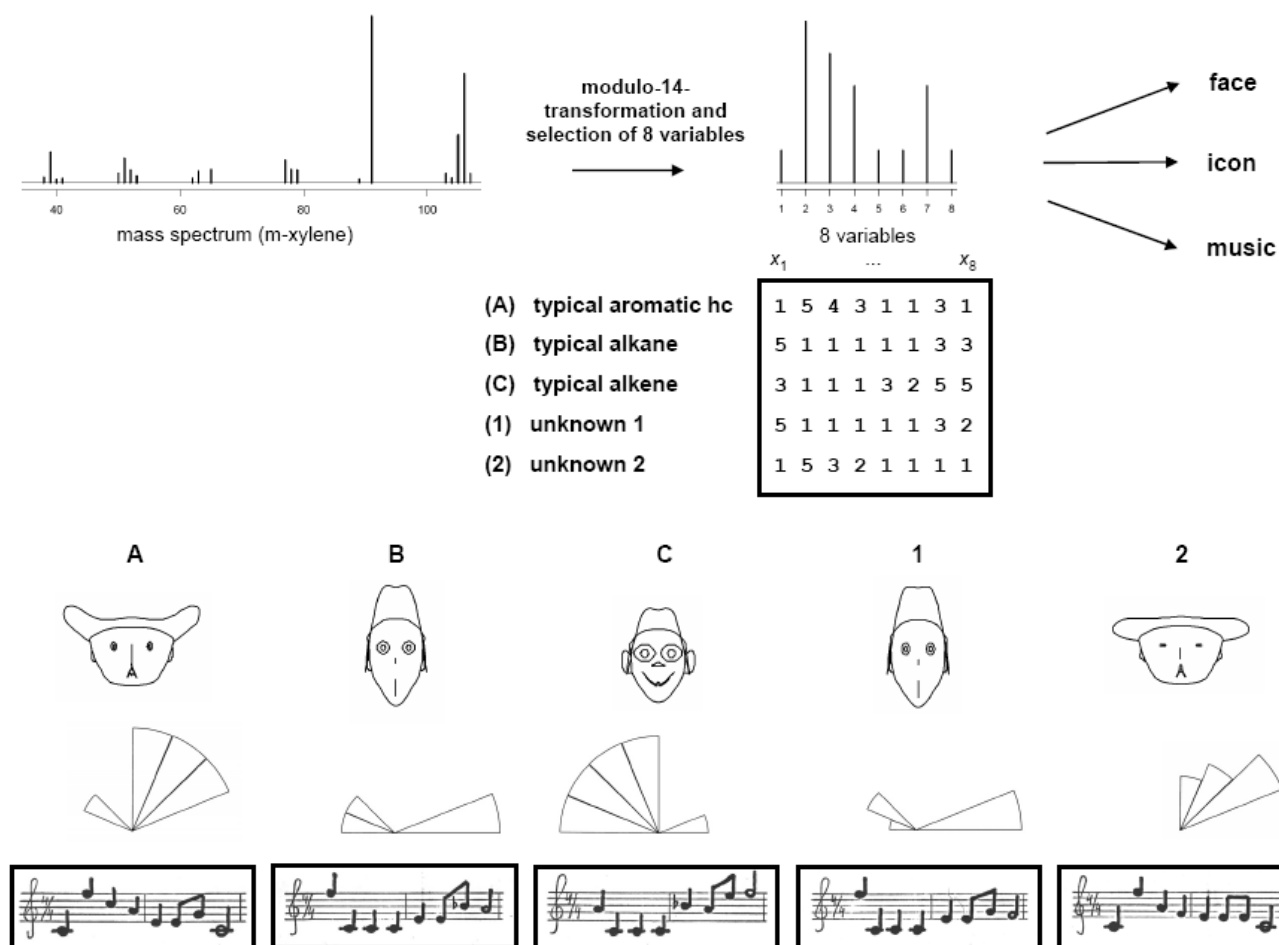


Figure 6.2. Representation of multivariate data by icons, faces and music for human cluster analysis and classification in a demo example with mass spectra. Mass spectra have first been transformed by modulo-14 summation (see Section 7.4) and from the resulting 14 variables 8 variables with maximum variance have been selected and scaled to integer values between 1 and 5. **A**, typical pattern for aromatic hydrocarbons; **B**, typical pattern for alkanes; **C**, typical pattern for alkenes; **1** and **2**, "unknowns" (2-methyl-heptane and meta-xylene). The  $5 \times 8$  data matrix has been used to draw faces (by function "faces" in the R-library TeachingDemos), segment icons (by R-finction "stars"), and to create small melodies [Varmuza 1986]. Both unknowns can be easily assigned to the correct class by all three representations.