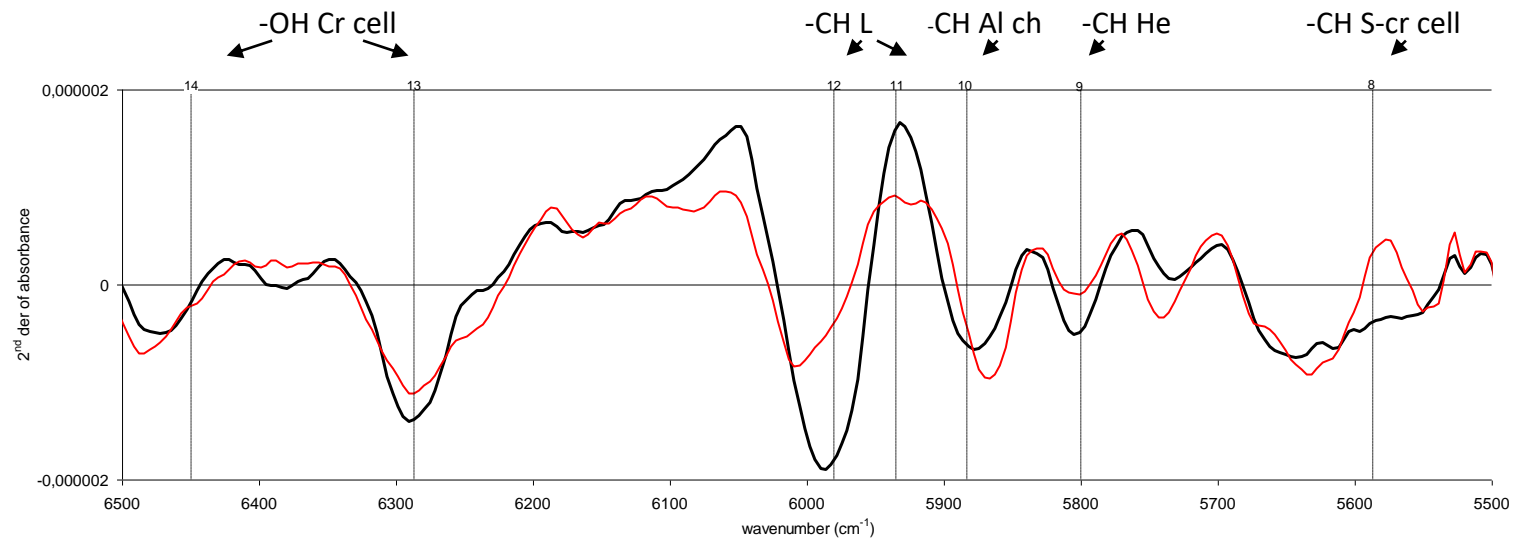
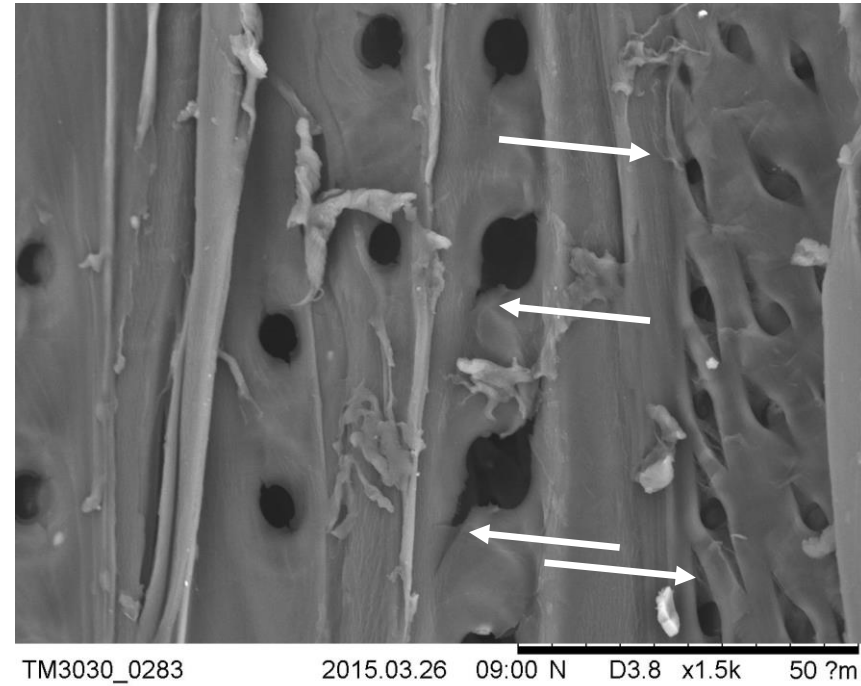


Modelling practice: data mining

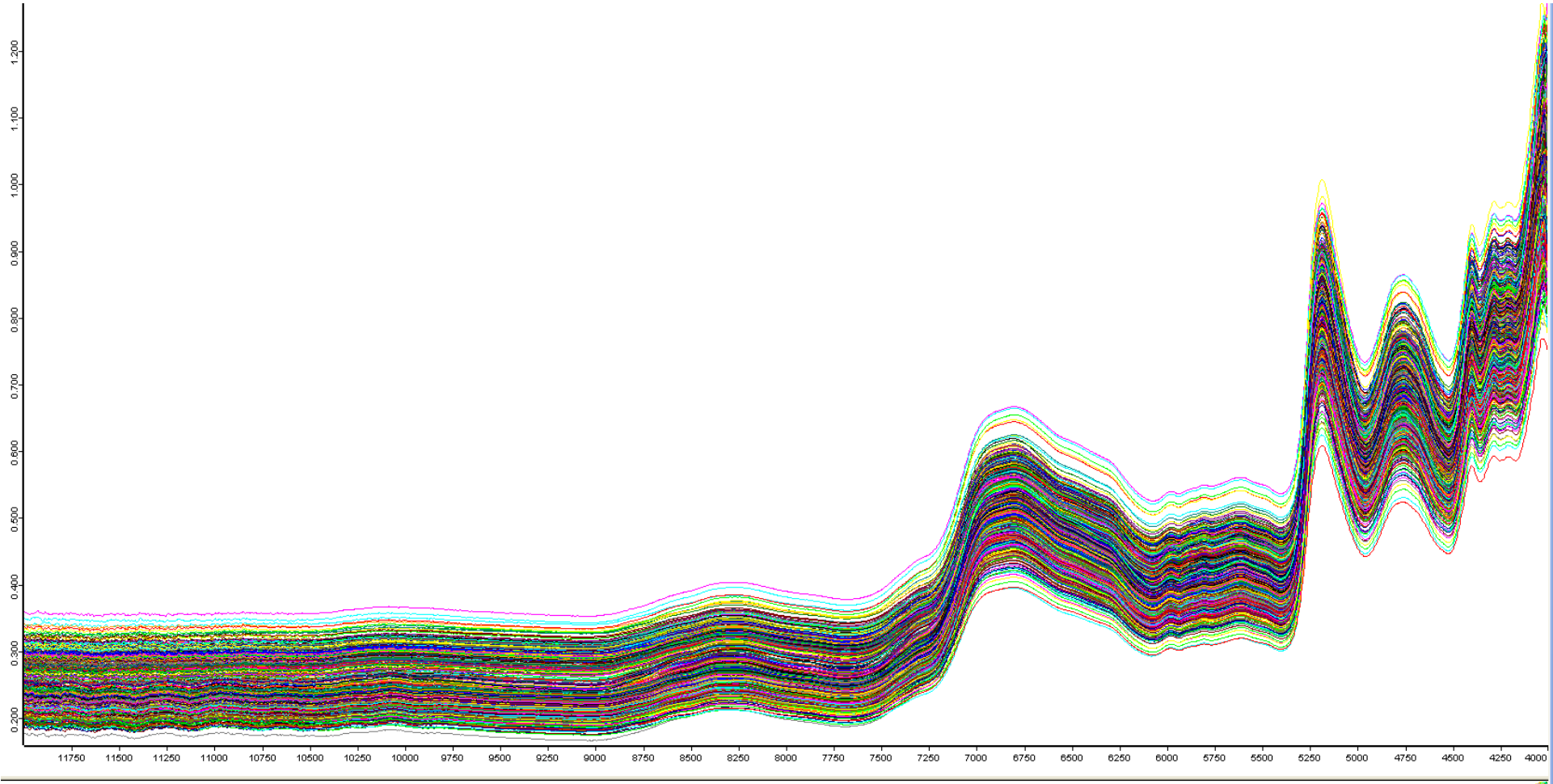
Jakub & Anna Sandak

CNR-IVALSA, San Michele all Adige, Italy

University of Primorska, Koper, Slovenia



When we have more spectra to be analyzed...





but...

Data is not the same as information

Too much data - too little information

(Harald Martens)

Chemometrics

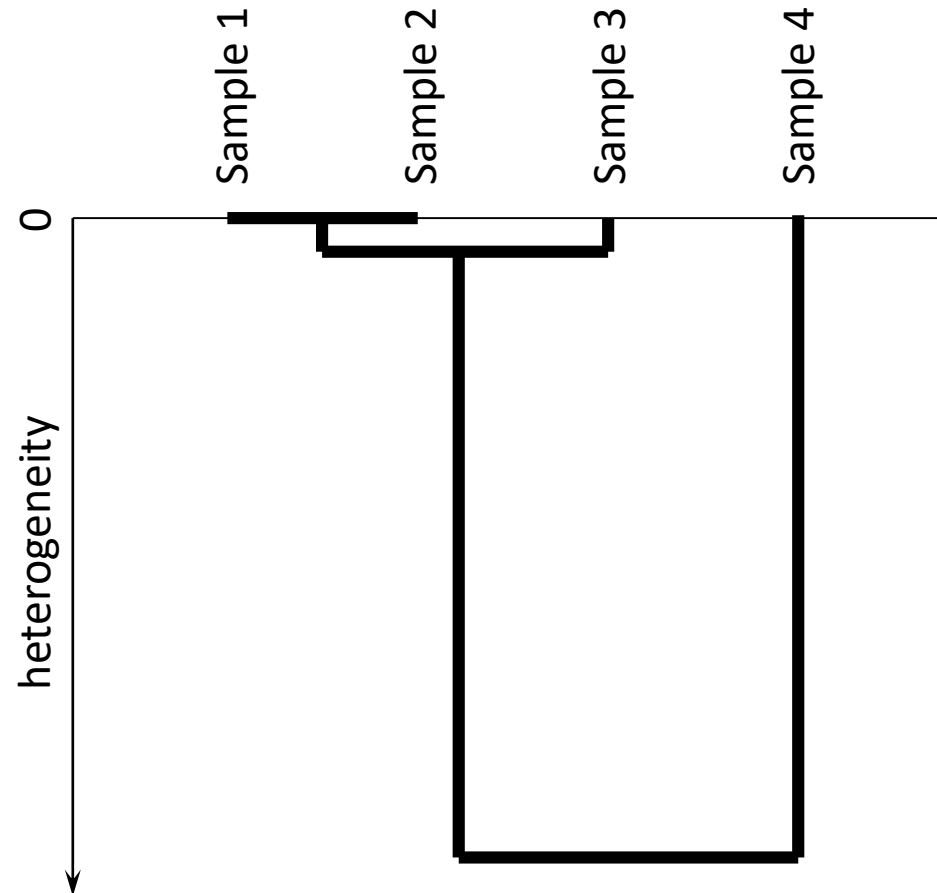
the **chemical** discipline that **uses** mathematical and statistical methods to design or select optimal procedures and experiments, and to provide maximum **chemical** information by analyzing **chemical** data

definition of the Chemometrics Society

Multivariate data analysis

- **Exploratory data analysis** (data mining) – attempts to find the hidden structure in large complex data sets
 - Cluster analysis
 - Principal Component Analysis
- **Regression analysis and Predictive Models** (developing the models from available data and predict desired response)
 - Partial Least Squares Regression
 - Multiplicative Linear Regression
- **Classification Models** (separation of group of object into one or more classes based on distinguished characteristic)
 - Cluster Analysis Test
 - Identity Test
 - SIMCA

Cluster analysis: CA

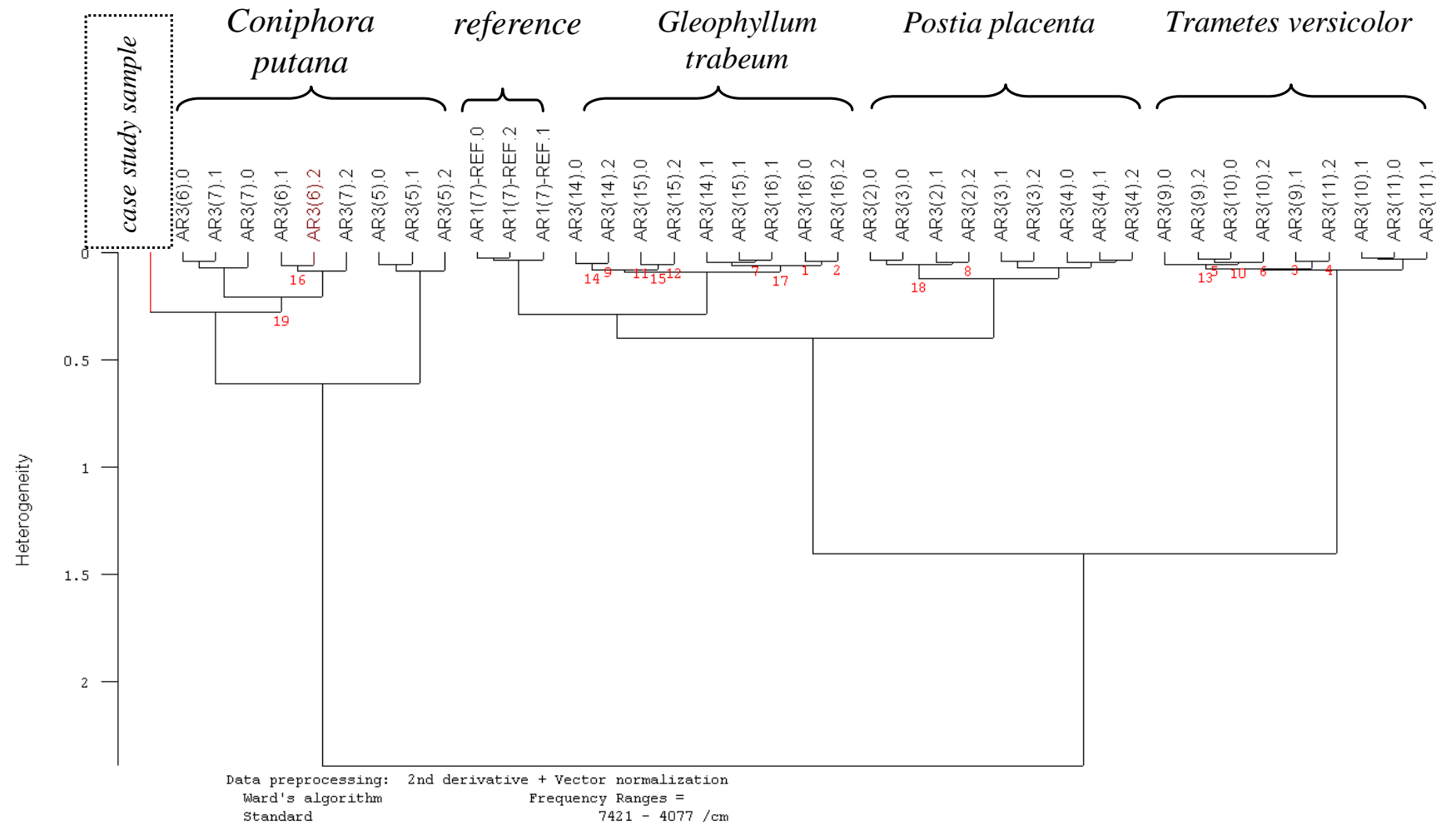


analyze spectra for their similarity, divide most similar spectra into groups, which are called clusters or classes

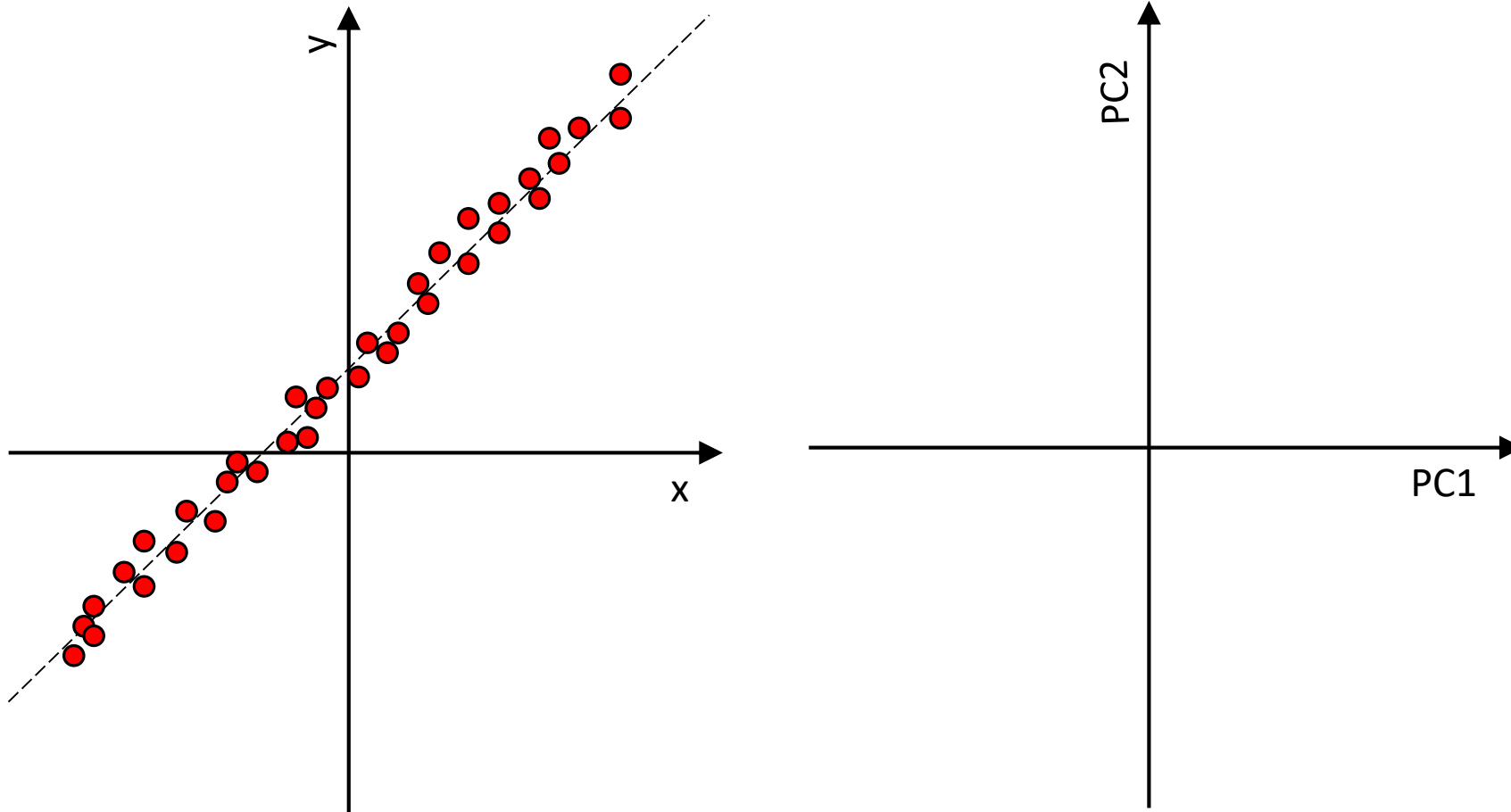
The spectral distance – heterogeneity explains the similarity between the spectra

The spectral distance 0 means that the spectra are identical

CA example



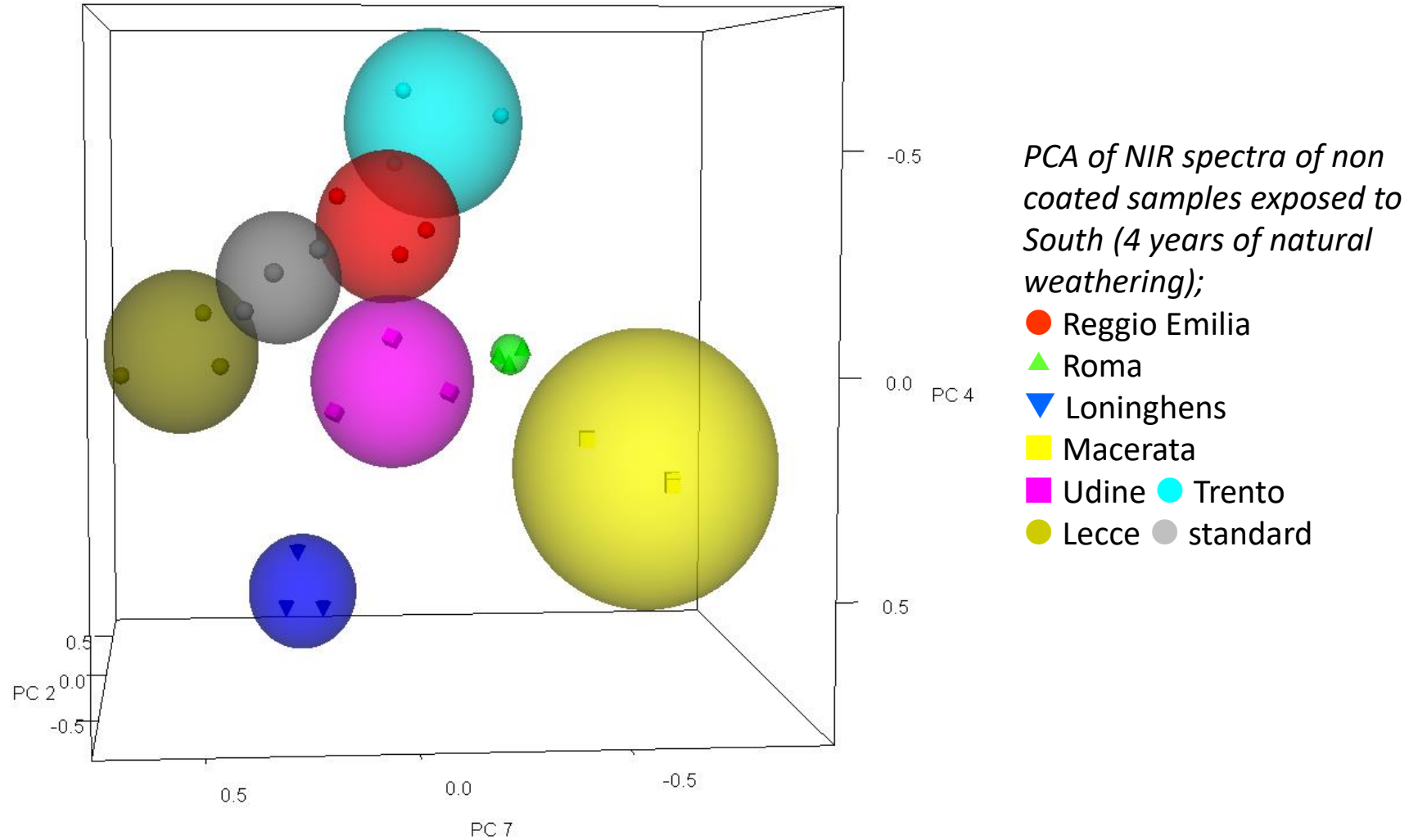
Principal Component Analysis



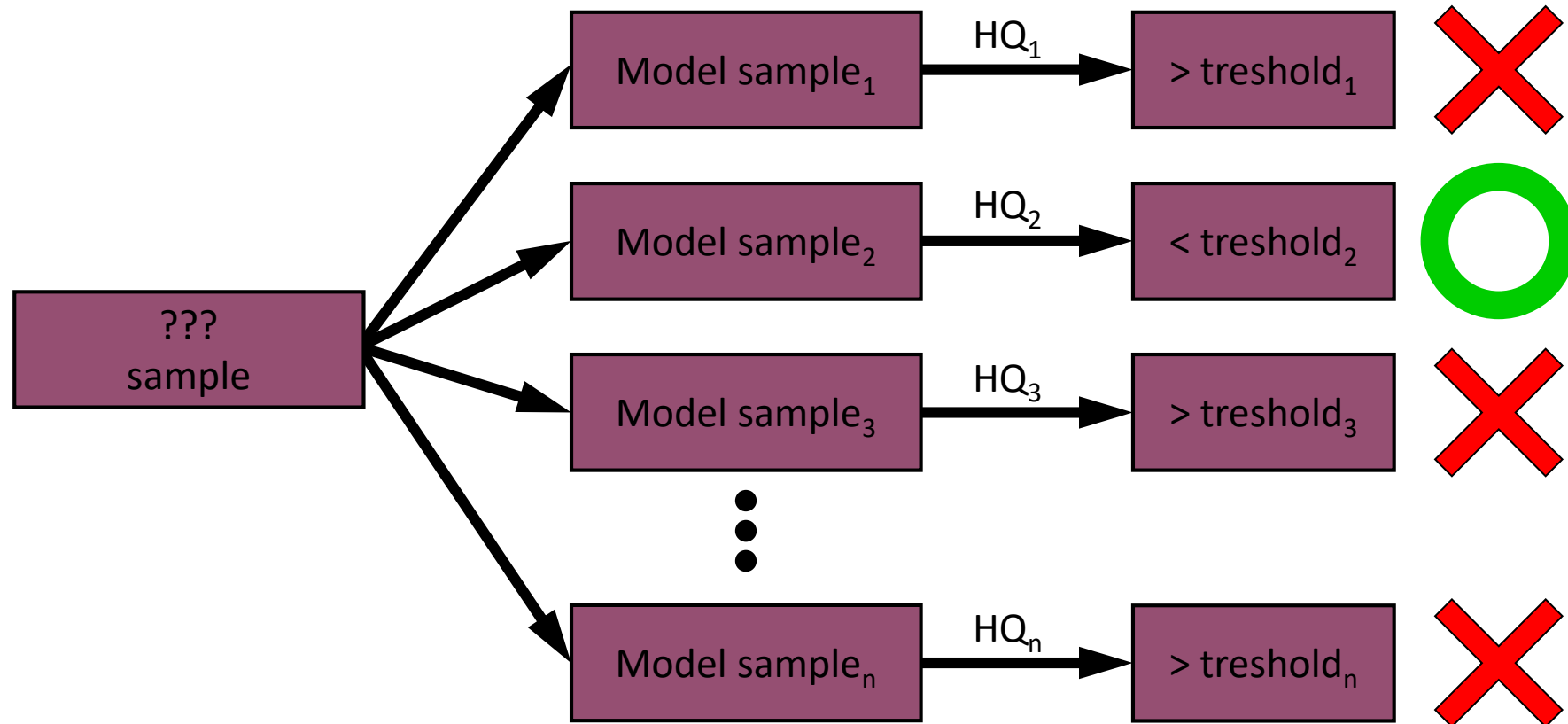
It is used for de-correlation of highly correlated data to reduce multidimensional data set to lower dimensions

it can separate set of input data into groups of peculiar similarities

PCA example



Identity test

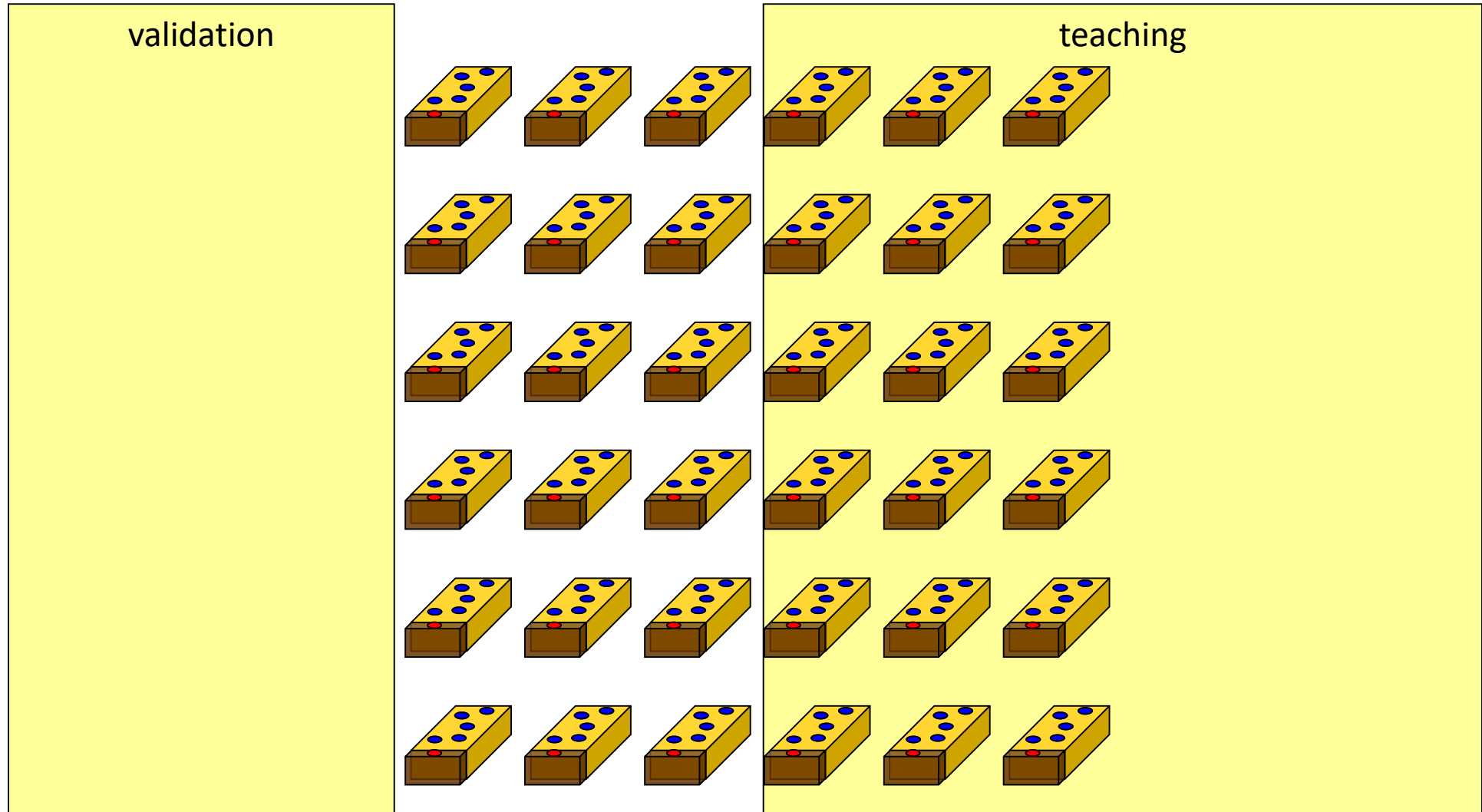


Compare the unknown spectrum with all reference spectra

The result of comparison between two spectra is the spectral distance called hit quality

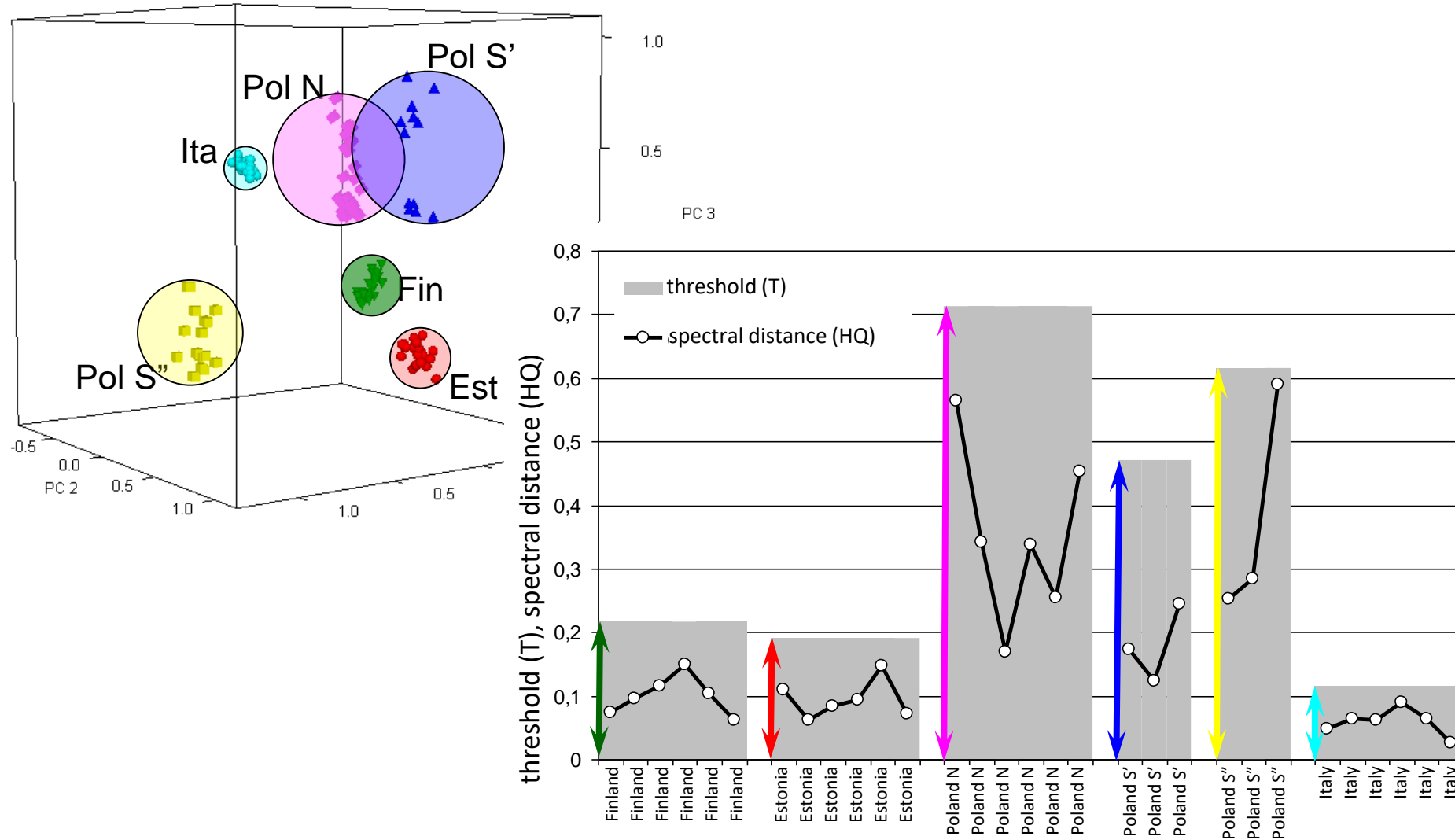
The better spectra match the smaller is spectral distance; HQ for identical spectra is 0

Identity test (IT)

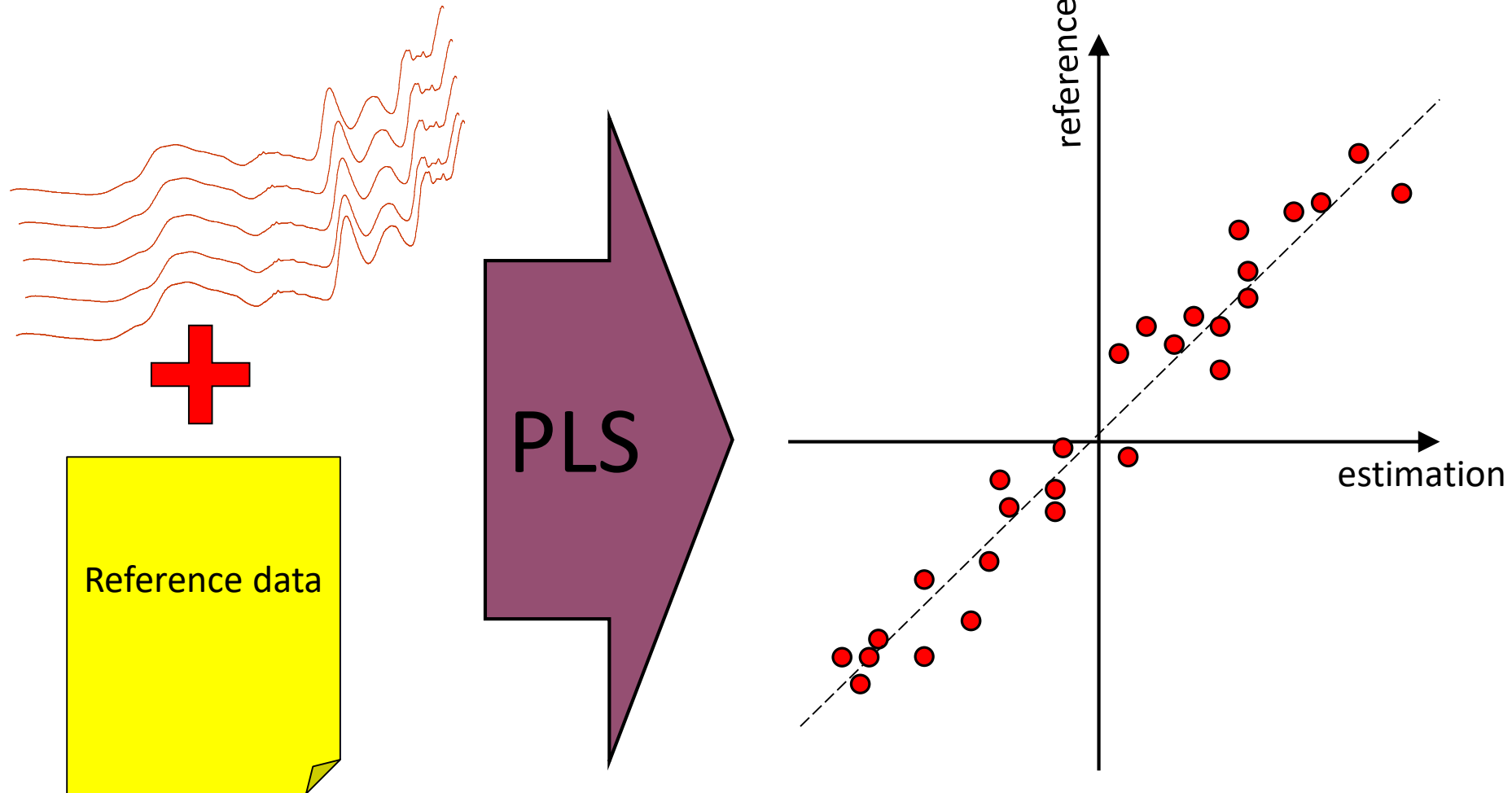


20% for validation (80% for building model)

Identity test (IT)



Partial Least Squares



How it works?

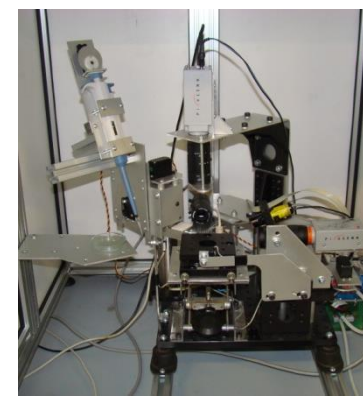
spectra



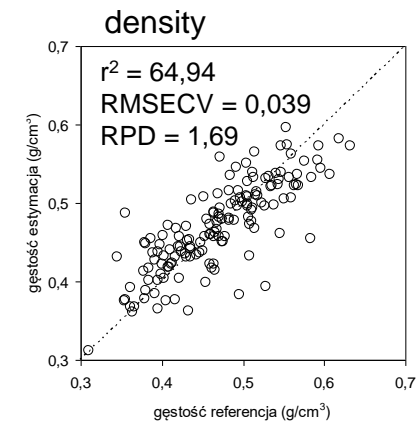
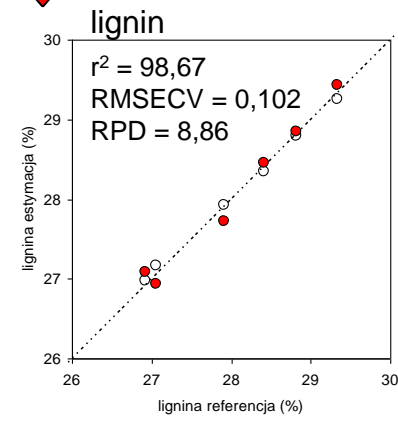
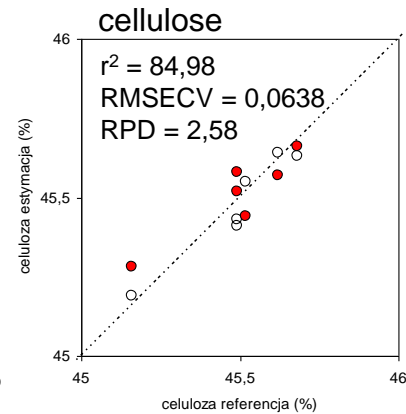
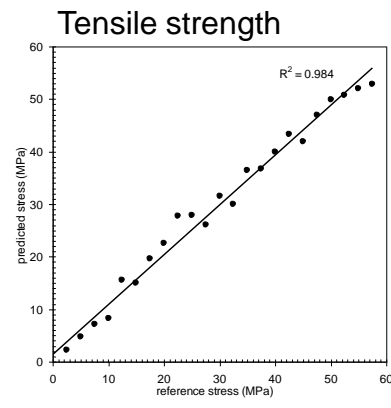
+



reference data



calibration (PLS)



PLS in practice