

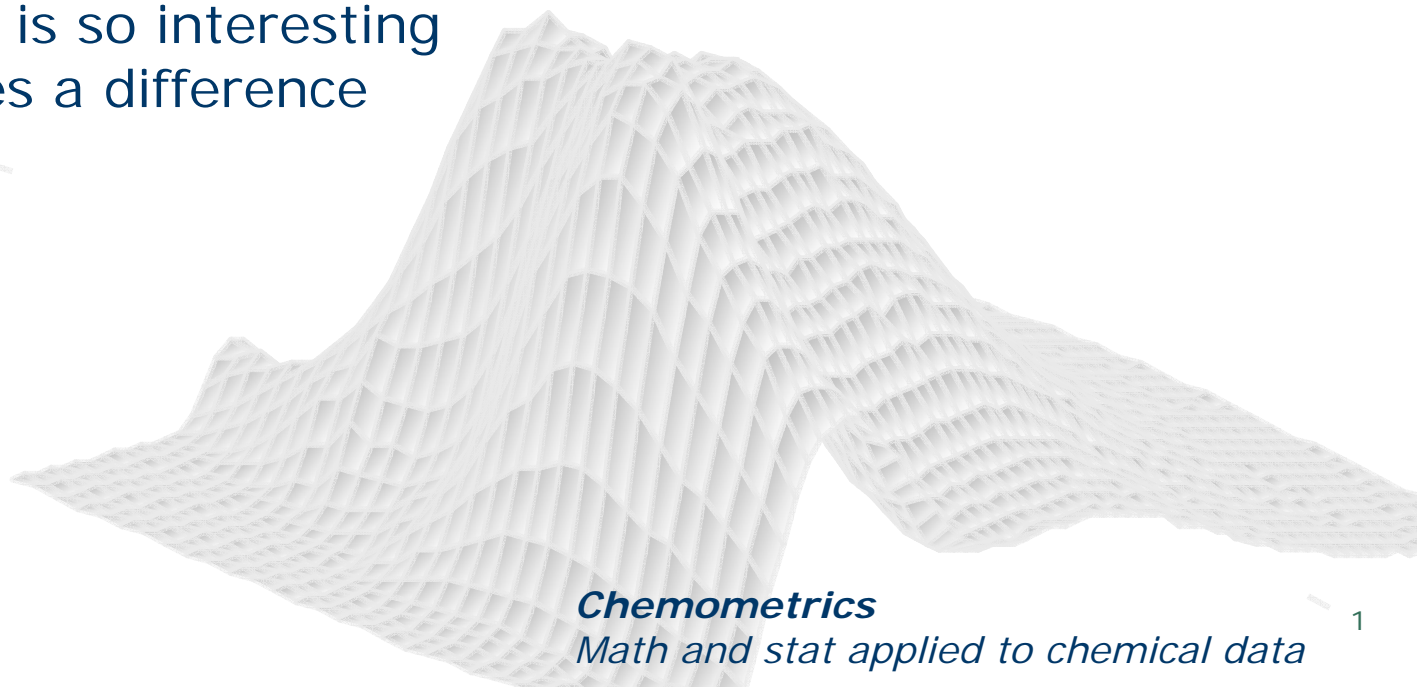
Multi-way analysis of bioinformatic data

Rasmus Bro

Chemometrics Group
Dept. Food Science
Royal Veterinary & Agricultural University (KVL)
rb@kvl.dk

Outline

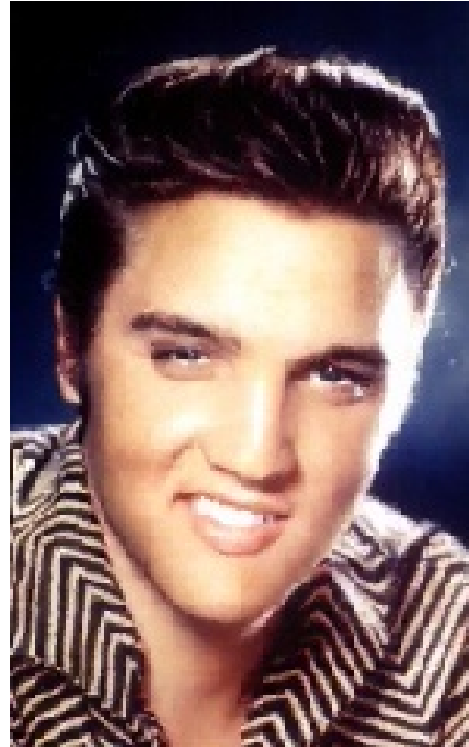
- A little on how chemometrics uses visualization
- Why PARAFAC is so interesting
- Where it makes a difference



Chemometrics

Math and stat applied to chemical data

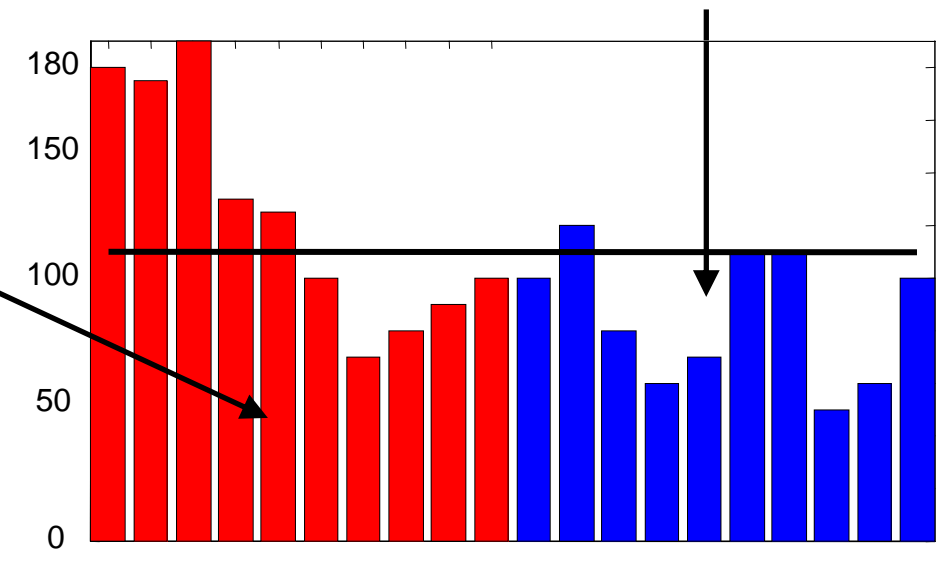
Human pattern recognition uses all available data



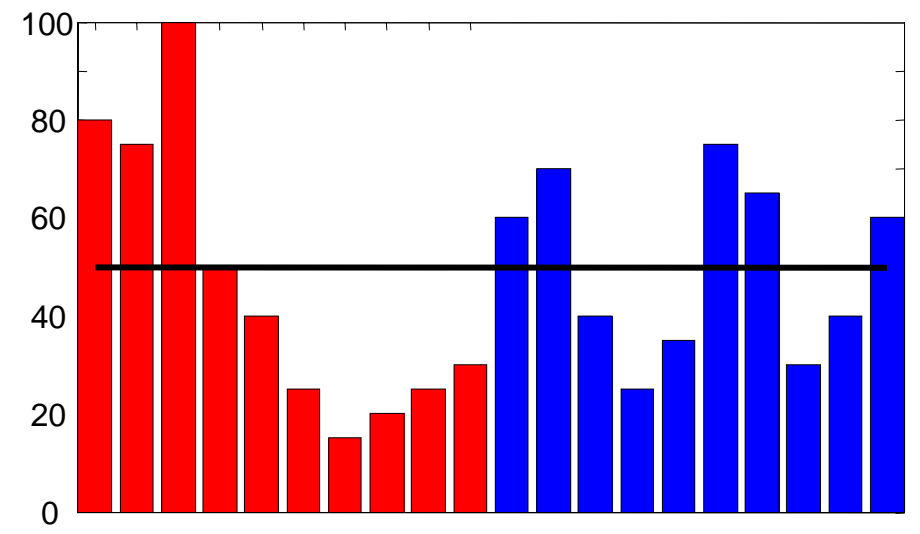
Single measurements/bar-plot

Height humans

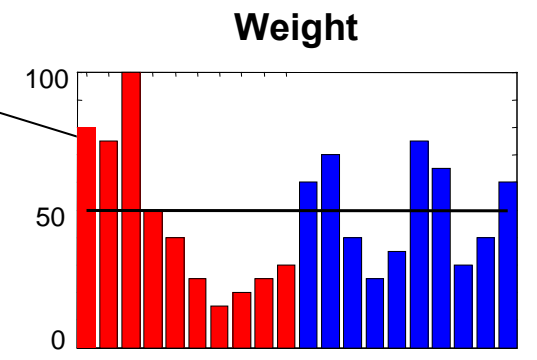
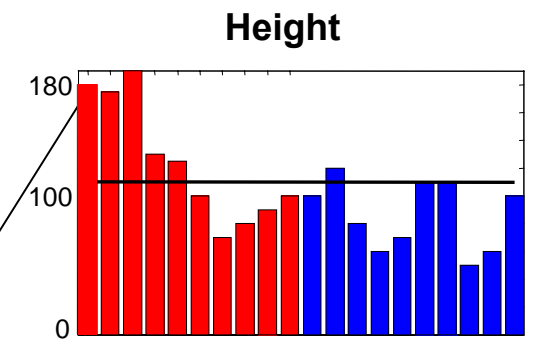
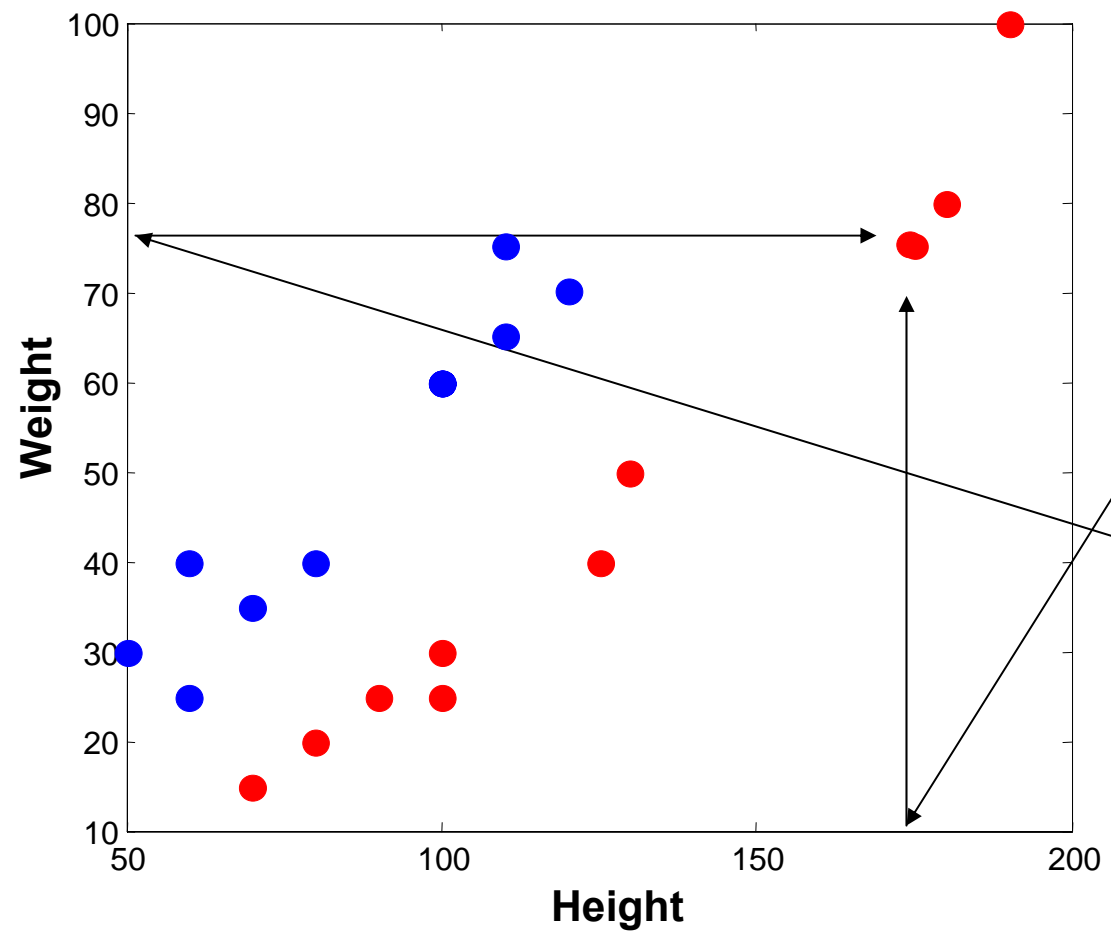
Height monkeys



Ditto weight



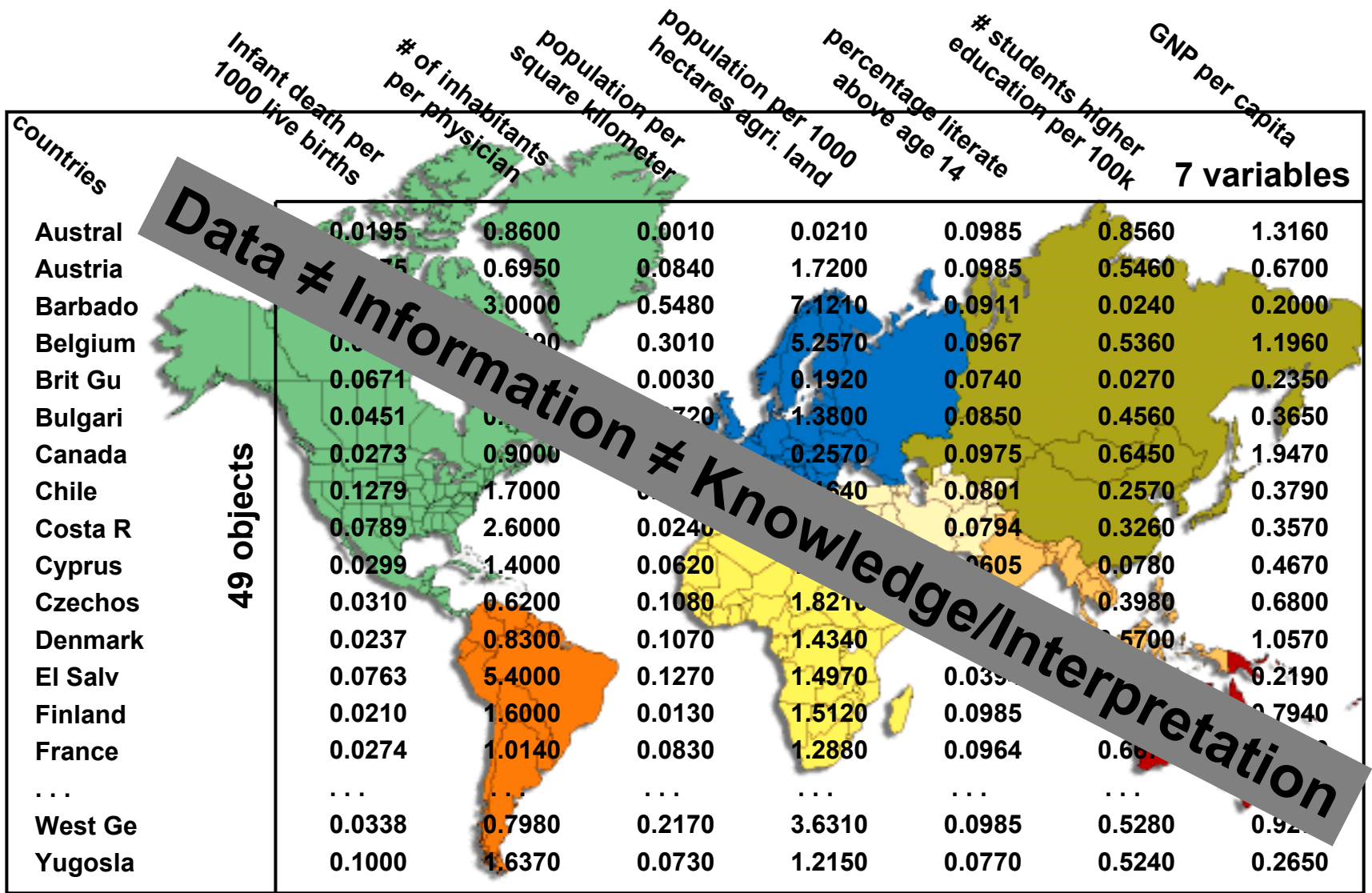
Pattern recognition



Principal component analysis movie

Removed. Find it at www.models.kvl.dk

Data analysis



Data ≠ Information ≠ Knowledge/Interpretation

Incredibly simple questions



What European country is most similar to Japan?

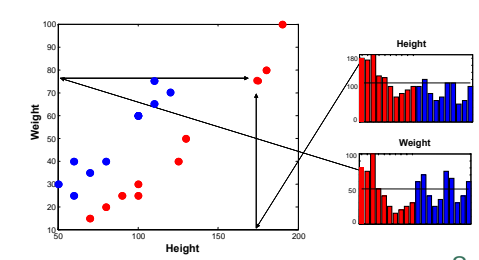
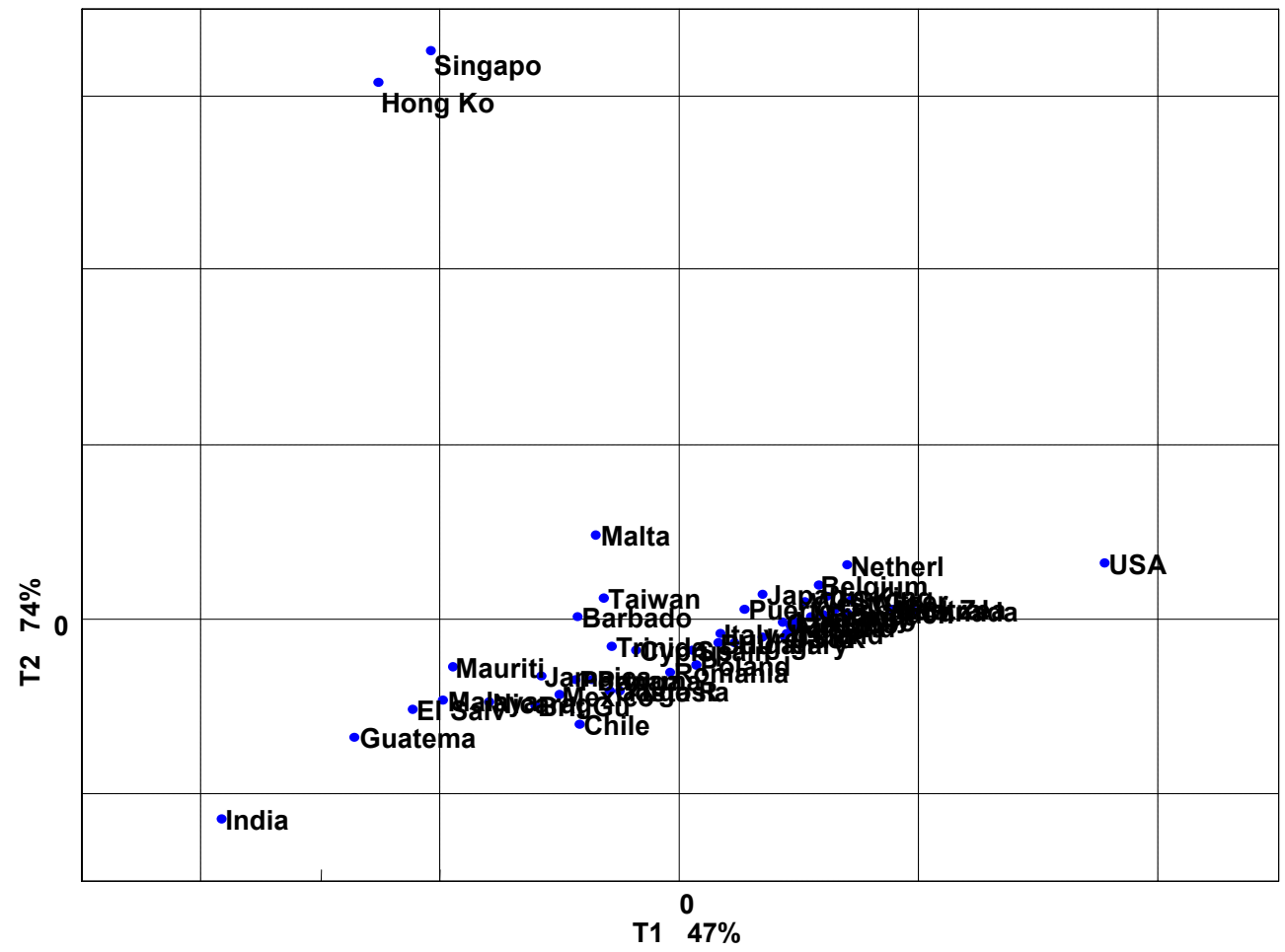
What country is most bizarre?

Principal Component Analysis

Outliers are easily spotted in score scatter plot

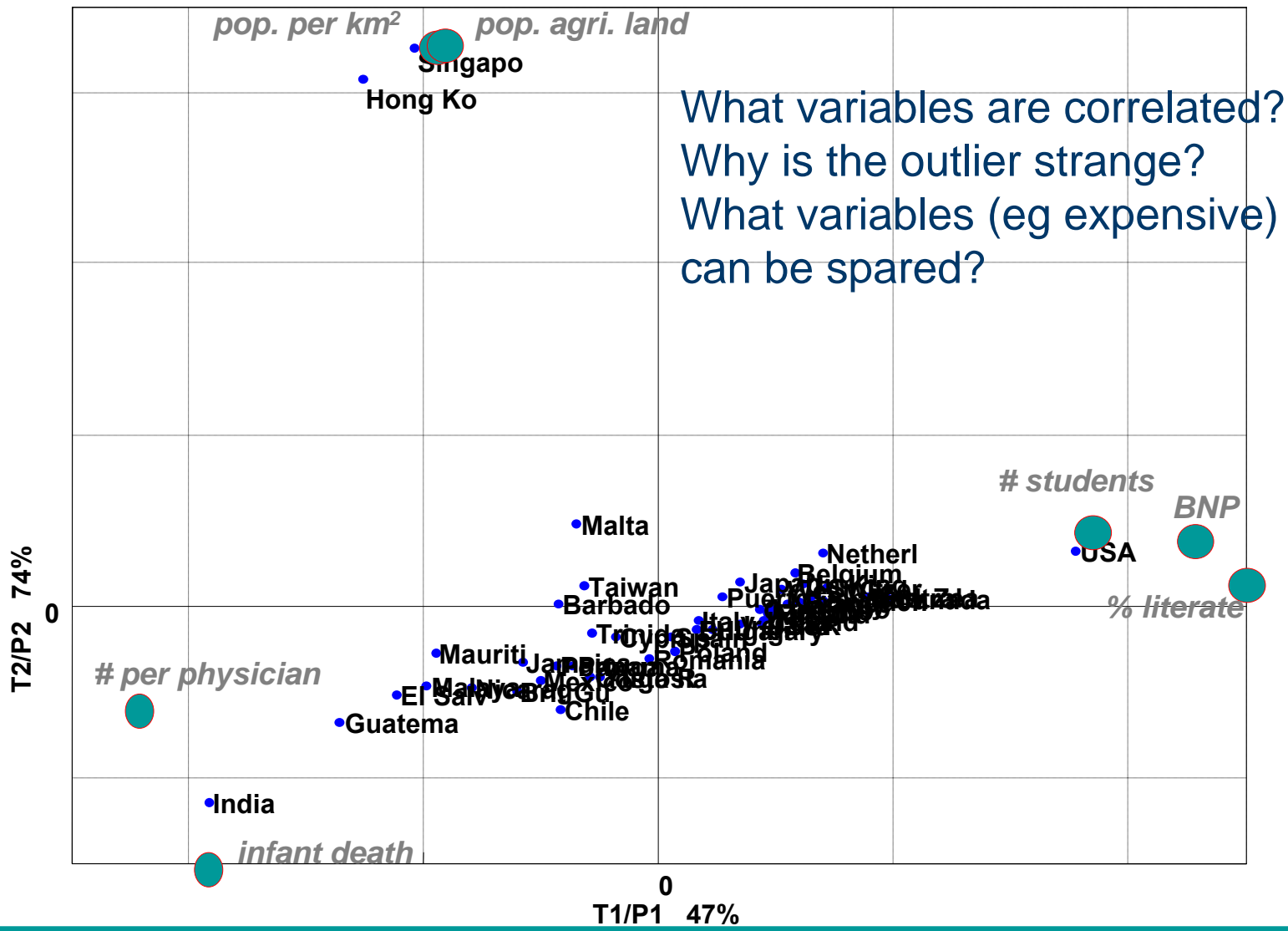
$$X = TP' + E$$

↑ Scores
 ↑ Loadings

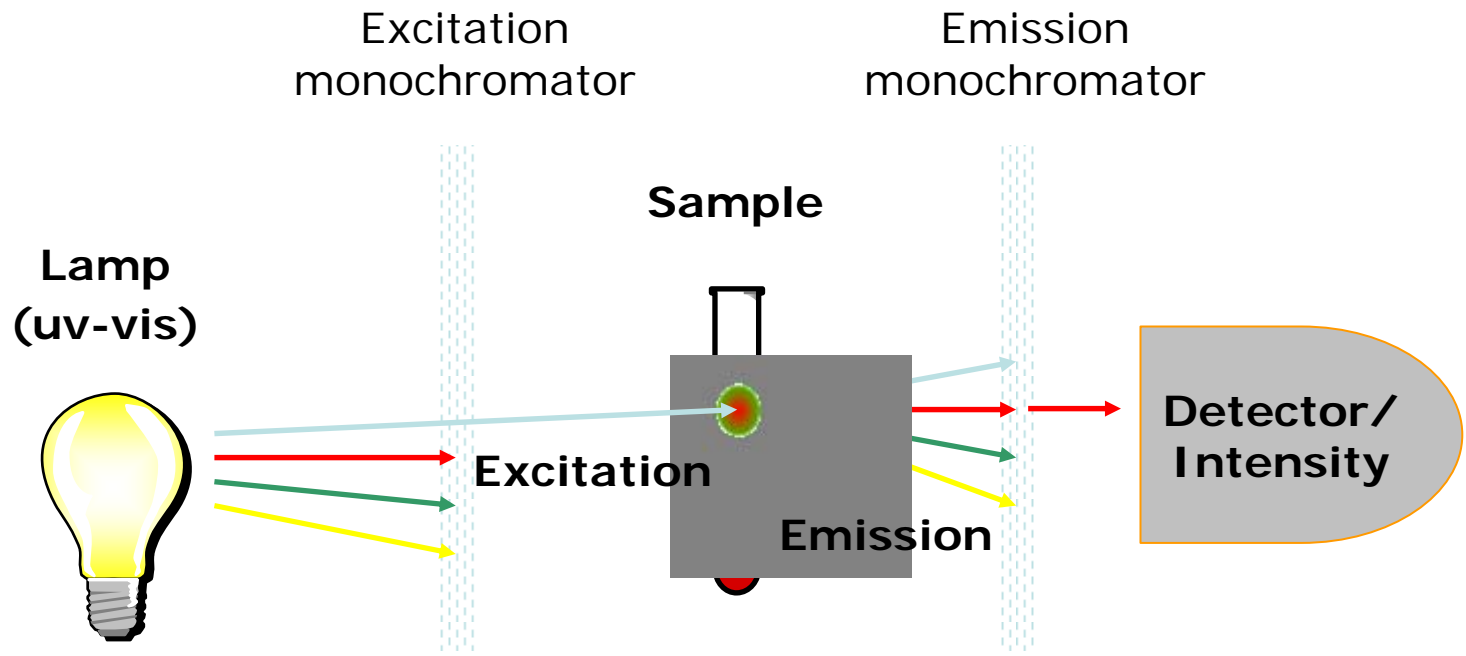


Principal Component Analysis

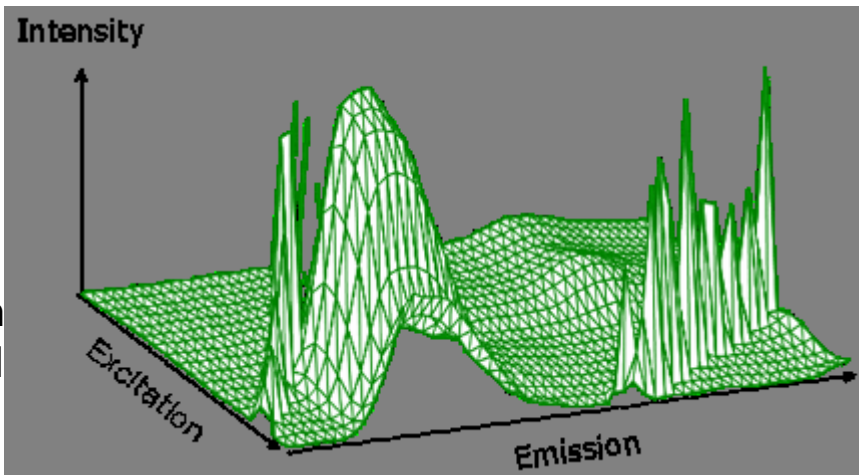
Why – Add loadings – bi-plot?



Fluorescence spectroscopy



Excitation-emission matrix – a chemical fingerprint



PARAFAC movie

Removed. Find it at www.models.kvl.dk

PARAFAC movie

Removed. Find it at www.models.kvl.dk

PARAFAC - algorithm

Efficient ALS algorithm

1. Initialize **B** and **C**

$$2. \mathbf{A} = \left(\sum_{k=1}^K \mathbf{X}_k \mathbf{B} \mathbf{D}_k \right) \{ (\mathbf{B}' \mathbf{B}) * (\mathbf{C}' \mathbf{C}) \}^{-1}$$

$$3. \mathbf{B} = \left(\sum_{k=1}^K \mathbf{X}'_k \mathbf{A} \mathbf{D}_k \right) \{ (\mathbf{A}' \mathbf{A}) * (\mathbf{C}' \mathbf{C}) \}^{-1}$$

$$4. \text{diag} \mathbf{D}_k = \{ (\mathbf{B}' \mathbf{B}) * (\mathbf{A}' \mathbf{A}) \}^{-1} \text{diag}(\mathbf{A}' \mathbf{X}_k \mathbf{B}), k=1, \dots, K$$

5. Step 2 until relative change in fit is small

* Hadamard (elementwise product)

$$x_{ijk} = \sum_{f=1}^F a_{if} b_{jf} c_{kf} + e_{ijk}$$

$$\mathbf{X}_k = \mathbf{A} \mathbf{D}_k \mathbf{B}' + \mathbf{E}_k$$

$$\mathbf{D}_k = \text{diag}(\mathbf{C}(k, :))$$

Why ALS?

Simple

Extends to N-way

Handles missing

Handles ML fitting

Constraints:

- Nonnegativity
- Unimodality
- Orthogonality
- Linear constraints
- Fixed parameters
- Smoothness
- Functional
- etc

Lawton & Sylvestre. Self modeling curve resolution. *Technometrics* 13:617-633, 1971.

Hanson & Lawson. *Solving least squares problems*, Englewood Cliffs: Prentice-Hall, Inc, 1974.

NMF dates back



Why ALS?

Simple

Extends to N-way

Handles missing

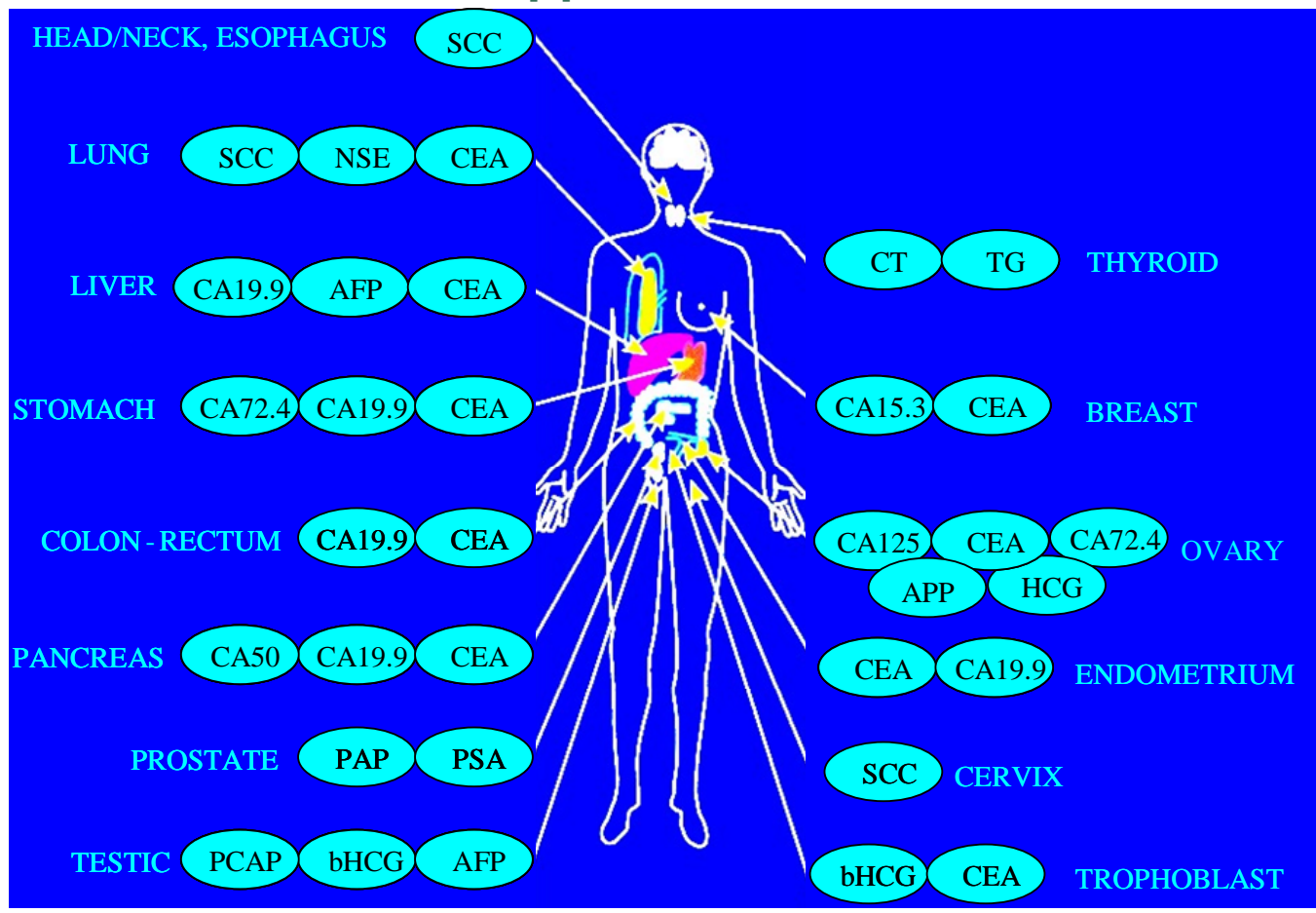
Handles ML fitting

Constraints:

- Nonnegativity
- Unimodality
- Orthogonality
- Linear constraints
- Fixed parameters
- Smoothness
- Functional
- etc

Cancer diagnostics

Traditional Approach Biomarkers

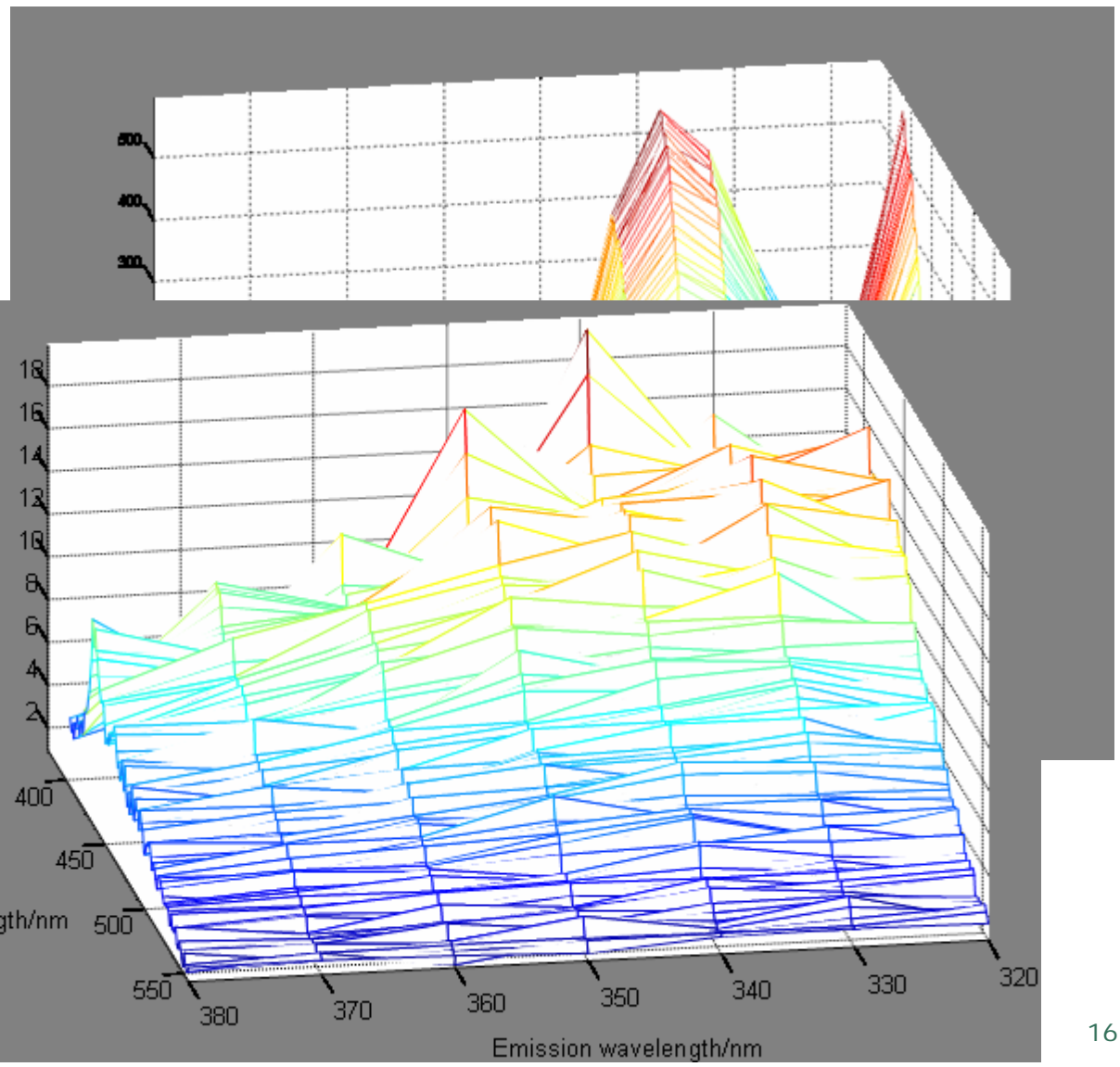


Cancer diagnostics

Alternative
Use fluorescence of
blood samples

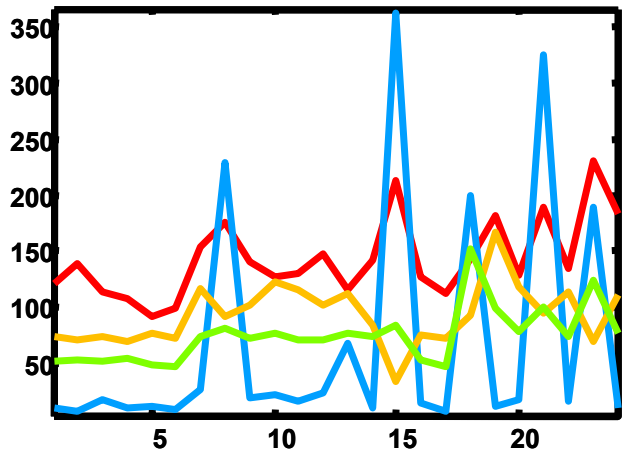
I.e. measure *many*
markers simultaneously

Here we use only a tiny
part of the overall data

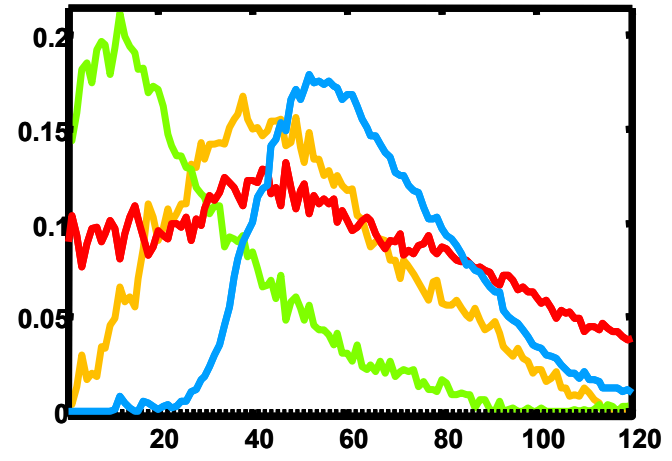


Cancer diagnostics

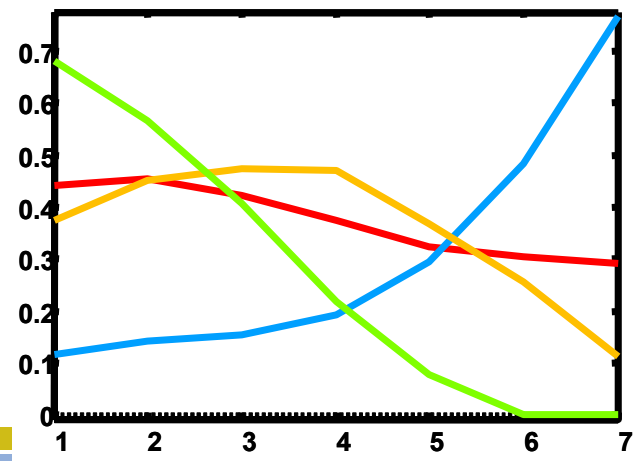
A = Concentrations



B = Emission spectra



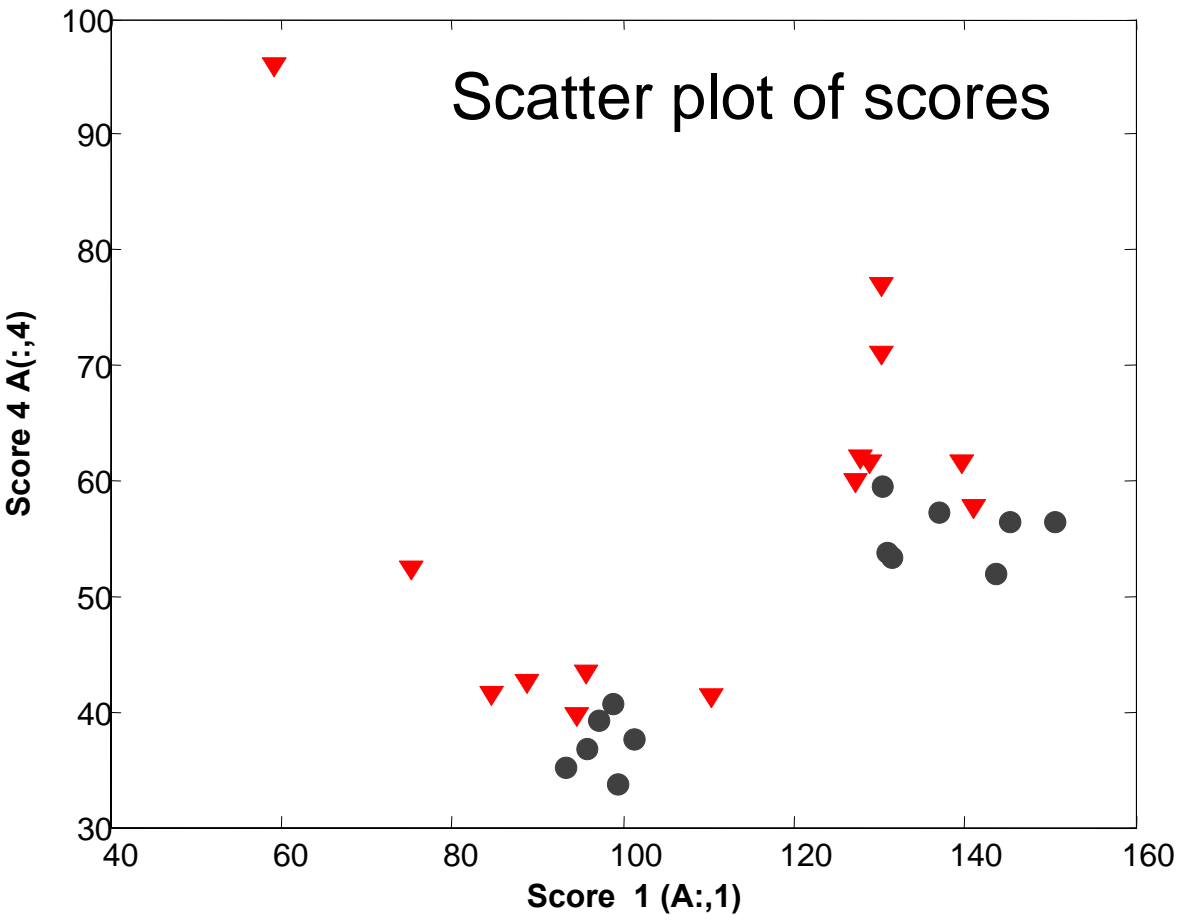
C = Excitation spectra



$$x_{ijk} = \sum_{f=1}^F a_{if} b_{jf} c_{kf} + e_{ijk}$$



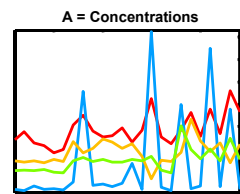
Cancer diagnostics



▼ Tumor
● Control

Grouping according to:
- Disease
- Something else.

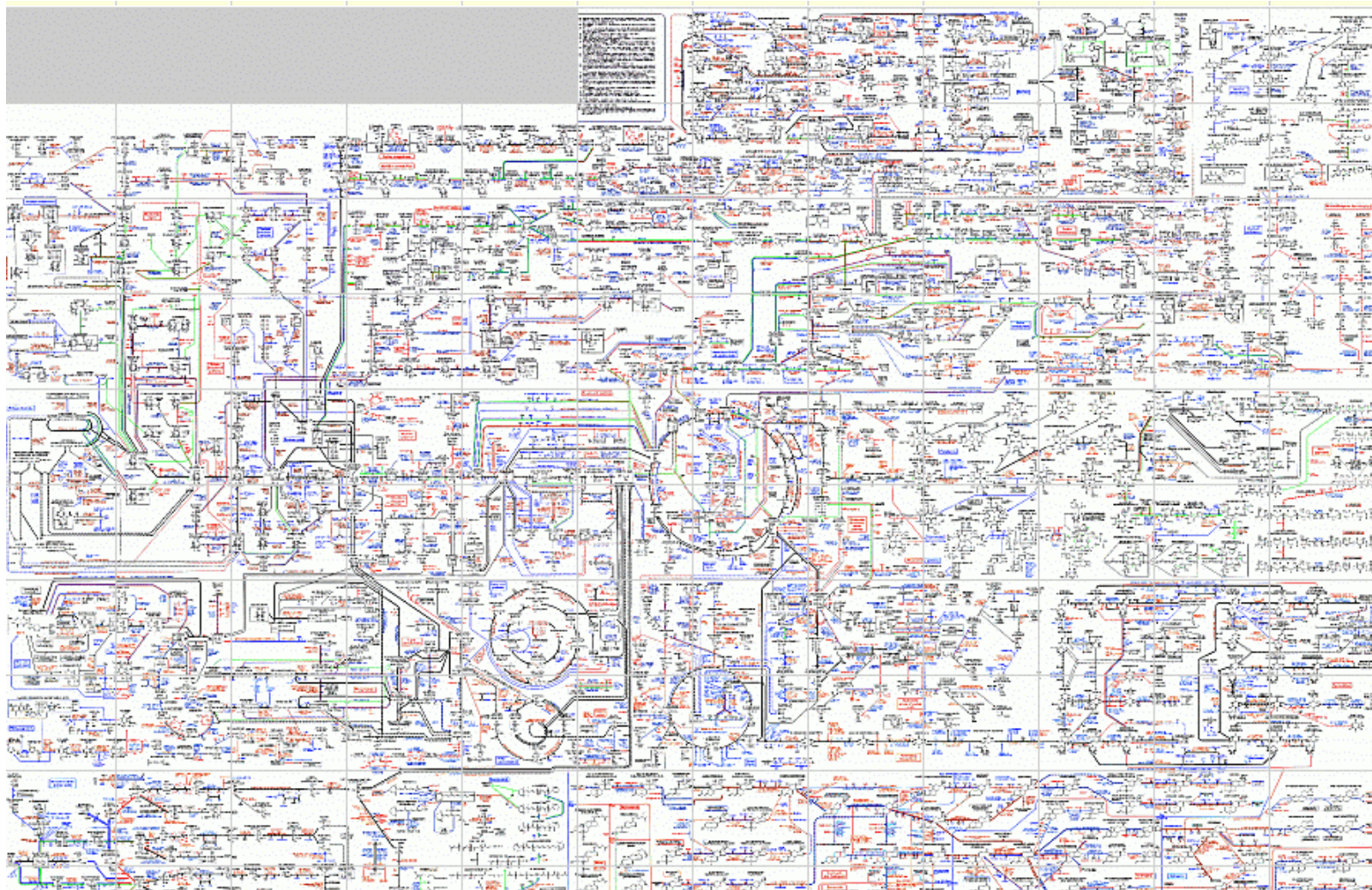
Chemically interpretable on a fluorescence level leading to understanding, biomarkers, etc.



Multi-way analysis in metabonomics

Metabonomics “the systematic study of the unique chemical fingerprints that specific cellular processes leave behind”.

A simplified view on
metabolic pathways



Toxic study

Typical data

- Control (5 rats)
- Low dose (5 rats)
- High dose (5 rats)

ELSEVIER Characterization and identification of metabolites by NMR www.elsevier.com/locate/jchem

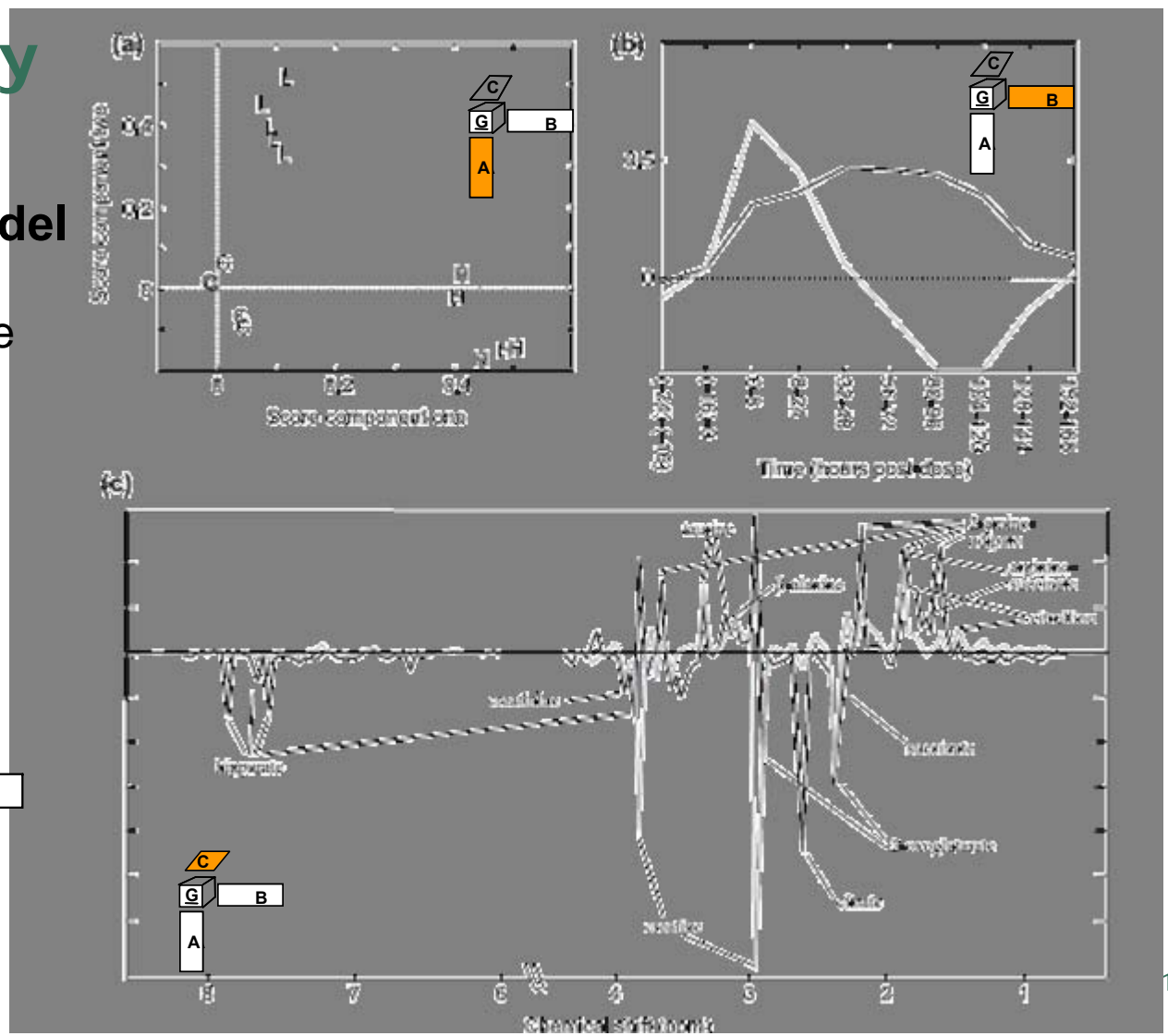
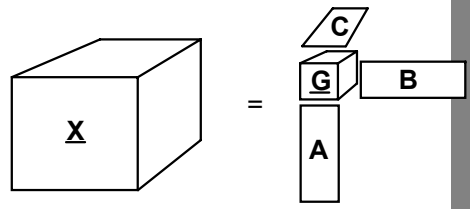
Multway chemometric analysis of the metabolic response to toxins monitored by NMR

Mona Dyby^a, Dennis Burggraf^a, Pärsons Bro^a, Søren Eklöv Egeiser^{a,b}

Toxic study

Exploratory Tucker model

Unlike unfolding/flattening:
Shows recovery of low-dose rats as well as an early response to taurine and creatine.



Toxic study

Multilevel component analysis of time-resolved metabolic fingerprinting data

Jeroen J. Jansen^a, Hub C.J. Hoefsloot^a, Jan van der Greef^{b,c},
 Marieke E. Timmerman^d, Age K. Smilde^{a,b,*}

Further analysis

Use ANOVA-SCA (simultaneous component analysis)

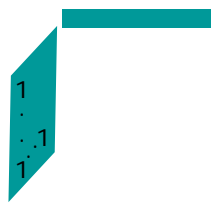
ANOVA

$$X_{rat_{dose}, time, nmr} = \mu_{nmr} + \alpha_{time, nmr} + (\alpha\beta)_{dose, time, nmr} + e_{rat_{dose}, time, nmr}$$

$X_{rat_{dose}, time, nmr}$



μ_{nmr}



$\alpha_{time, nmr}$



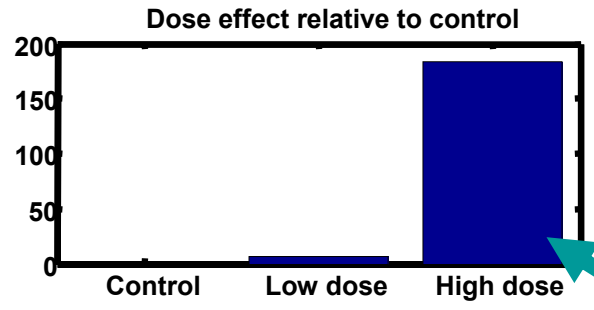
$(\alpha\beta)_{dose, time, nmr}$



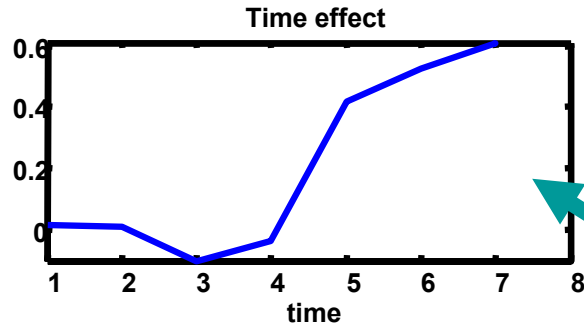
PARAFAC on dose/time-effect

$$(\alpha\beta)_{dose,time,nmr}$$

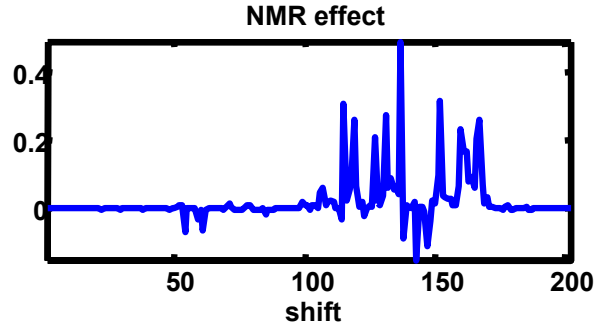
Separate the effect into a shock



High-dose on

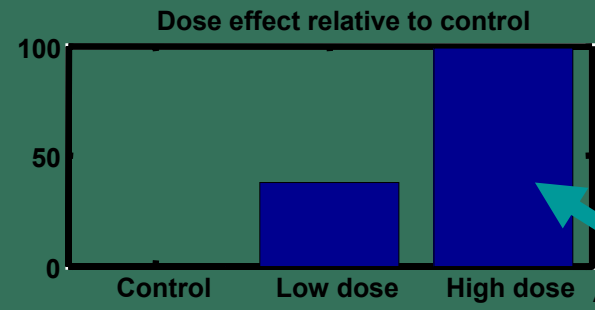


Irreversible

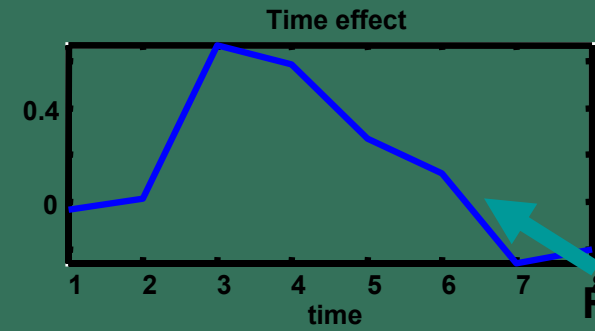


PARAFAC factor 1

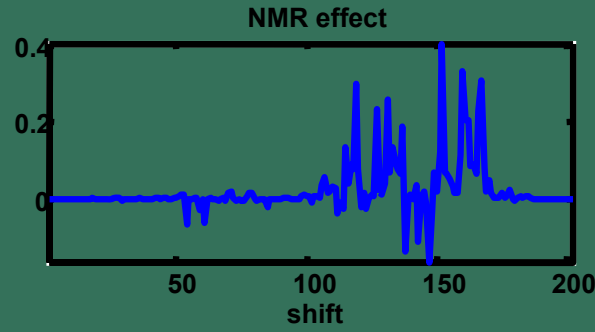
and a reversible effect



All-dose linear



Reversible



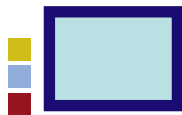
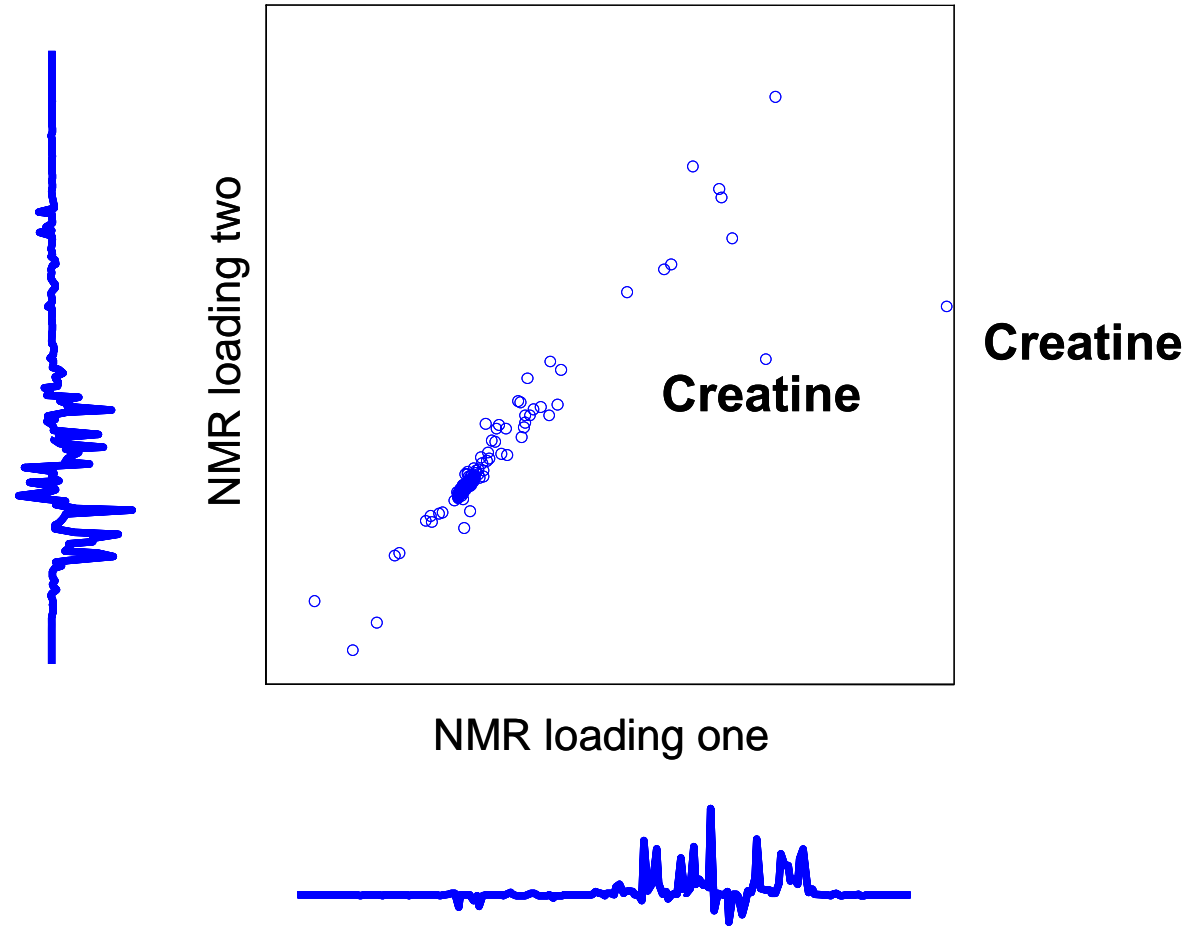
PARAFAC factor 2



Toxic study

$(\alpha\beta)_{dose,time,nmr}$

Difference between reversible and irreversible effect
Creatine indicating chronic kidney damage



Concluding remarks

Many interesting solutions using tensor approaches

Uniqueness

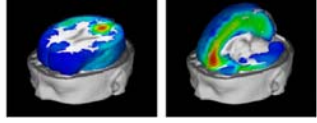
- Pure spectra
- Pure profiles
- Pure concentrations
- Pure magic!!!

=

- Mathematical chromatography

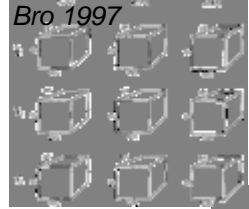
EEG

Miwakeichi et al 2004



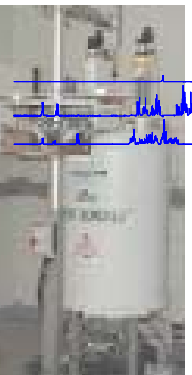
Other Examples

5-way analysis

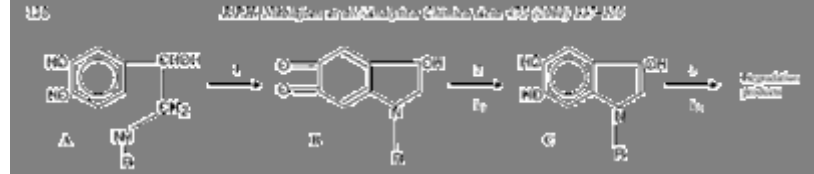


DOSY

Toft et al 2004

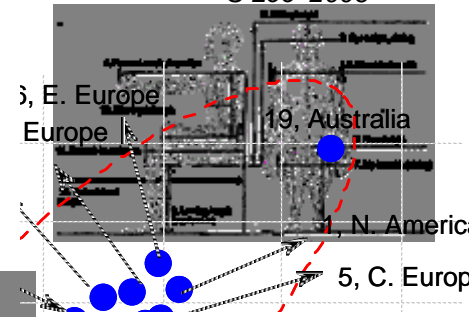


Quantify catecholamine in urine

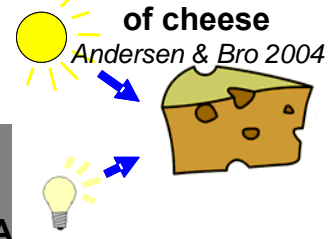


Anthropometry

S Lee 2006

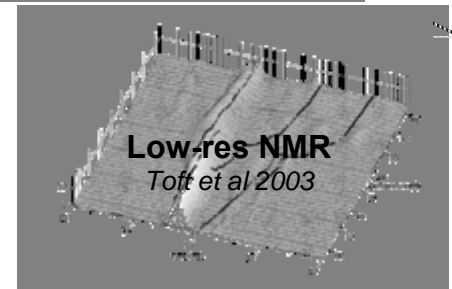


Light-induced oxidation of cheese



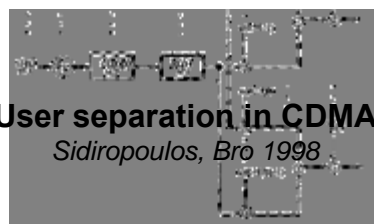
Low-res NMR

Toft et al 2003



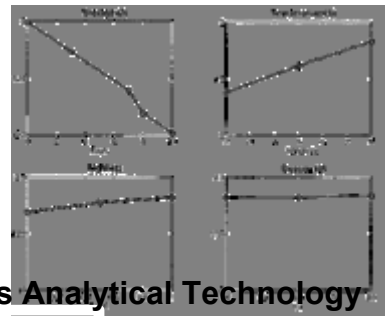
User separation in CDMA

Sidiropoulos, Bro 1998



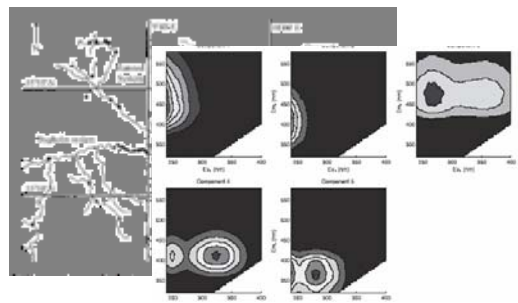
Generalized ANOVA

Bro & Jakobsen 1996, 2002



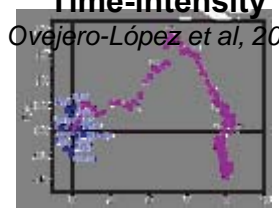
Tracing DOM

Stedmon, Markager, Bro 2003



Time-intensity

Ovejero-López et al, 2004



Process Analytical Technology

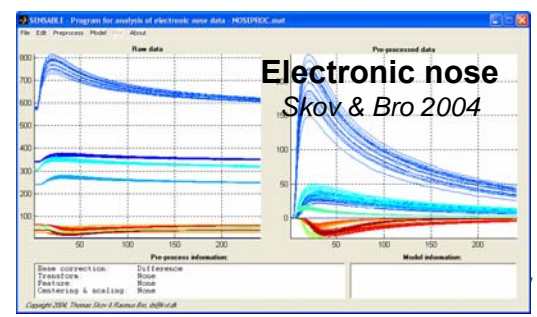
Datamining

Bro 1998



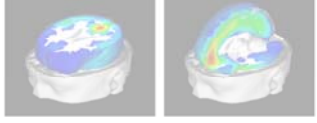
Electronic nose

Skov & Bro 2004

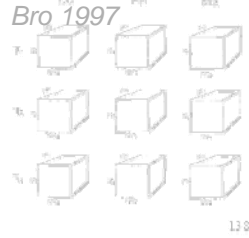


EEG

Miwakeichi et al 2004

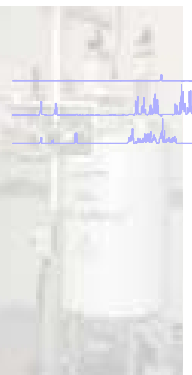


5-way analysis



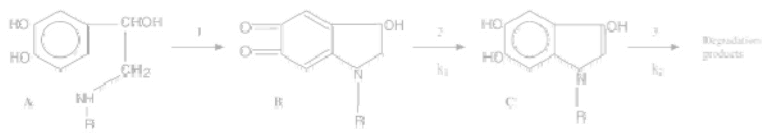
DOSY

Toft et al 2004



Quantify catecholamine in urine

R.P.H. Nijboer et al., Analytica Chimica Acta 475 (2002) 237-250



Light-induced oxidation of cheese

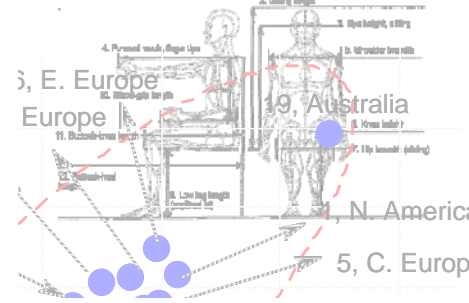
Andersen & Bro 2004



Other Examples

Anthropometry

S Lee 2006



Low-res-NMR

Toft et al 2004



www.models.kvl.dk

Time-intensity

Ovejero-López et al, 2004

