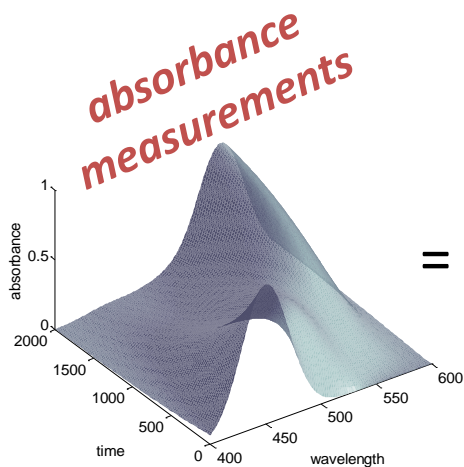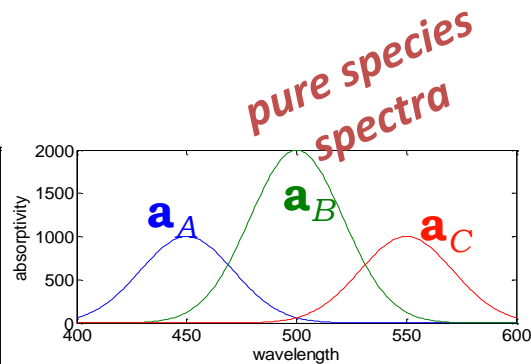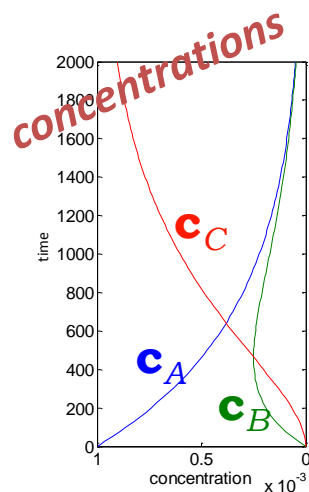# A "not exhaustive" overview of chemometric methods
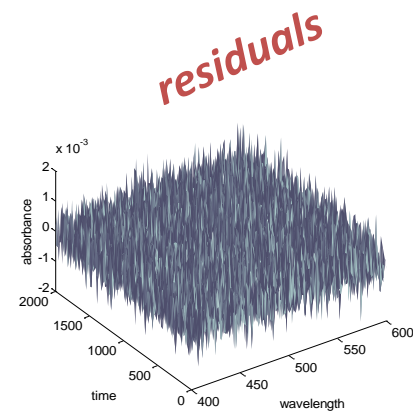## a practical guide for data analysis in chemistry
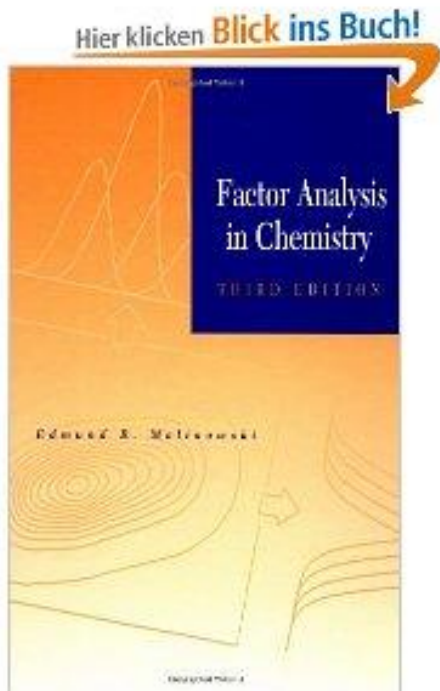
Cap Sébastien
28 November 2014

**Chemometrics** is the science of extracting information from chemical systems by **data-driven means**. It is a highly interfacial discipline, using methods frequently employed in core data-analytic disciplines such as

multivariate statistics,  applied mathematics, and computer science,

in order to address problems in

chemistry, biochemistry, medicine, biology and chemical engineering.

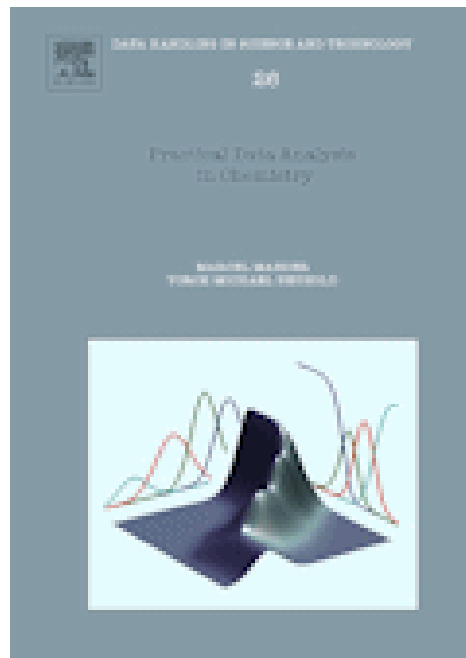In this way, it mirrors several other interfacial '-metrics' such as psychometrics and econometrics.

**Factor Analysis in Chemistry, 3rd Edition**
Edmund R. Malinowski
ISBN: 978-0-471-13479-4
432 pages
March 2002

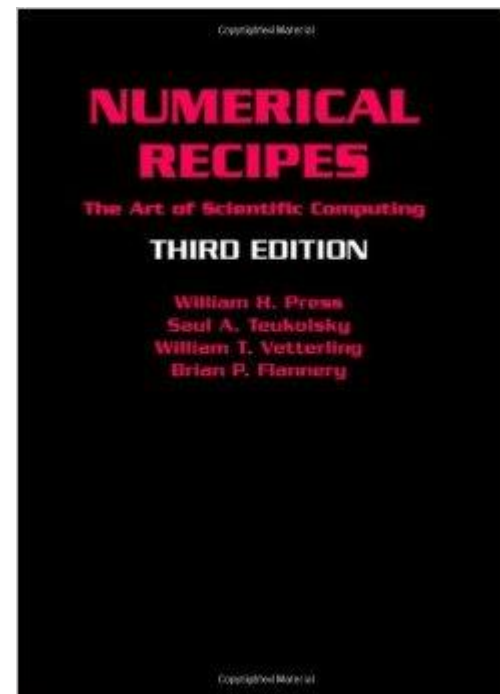**Practical Data Analysis in Chemistry**
**By** Marcel Maeder
Yorck-Michael Neuhold
Published: July 2007
Imprint: Elsevier
ISBN: 978-0-444-53054-7

**Numerical Recipes 3rd Edition: The Art of Scientific Computing Hardcover – September 10, 2007**
ISBN-13: 978-0521880688

# Table of content

## Reactor / Instrumentation

## Chemistry



## Plausible mechanism?



## Data



## Typical Arrhenius representation

Plausible kinetic model is probably wrong…

…after discussion with your gr. leader…new plausible model…

Chemometric will support you for:

- How to design/select the reactor/instrumentation?

- How to do the experimental design?

- How to "find" the correct kinetic model?

- What to "fit" and how ?

- How to determine the rate constant?

- How to "fit" if "C" can not be isolated and unknown?

- What about baseline drift, shift, noise level?

- Finally, are my fitted parameters "correct" to which extends?

- And we have to be quick to do all the above tasks

## Parameters

- Custom designed and versatile
- Adapt lot of probes?
- Time dependent acquisition
- In-situ
- Well controlled
- Use of (SOPs, standard operation protocols)
- Fully automatized
- Fast data acquisition
- Acquire as much data as possible (then only average)

## Reactor and control



CAD

LabView

## Instrumentation

- Adapt an orthogonal instrumental methodology



- Prefer multivariate signal to univariate signal
- Prefer integrated signals (spectroscopy) to differential signals (calorimetry)



Calorimetry

Spectroscopy

## Why a DoE?

➢ Used to reduce the design costs
➢ Speeds up the process design
➢ Reduce late engineering design
➢ Reduce raw and product material
➢ Reduce experimental complexity
➢ Reduce data analysis complexity
➢ Improve the overall robustness of data analysis
➢ Improve the certainty on the fitted parameters



## Experimental purpose

➢ Comparing alternative
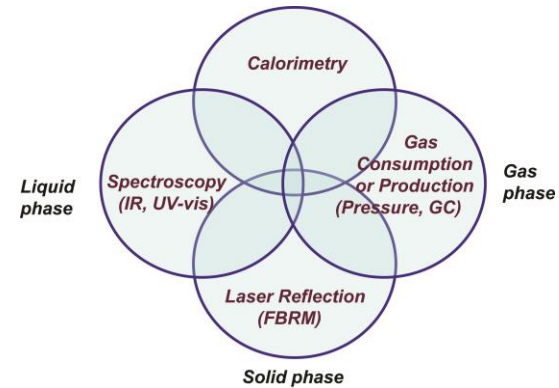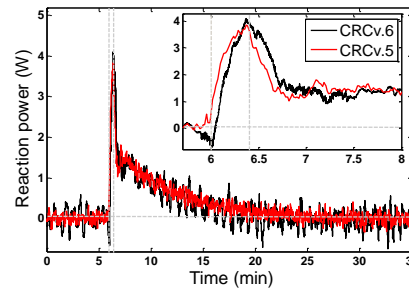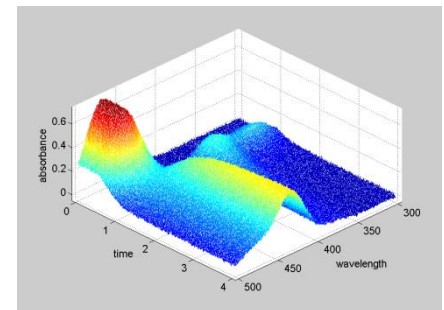  ➢ Egg from different supplier gives the same cake?
➢ Significant inputs (Factors)
  ➢ Sort the relevance of factors: is the amount of sugar more important than the amount of flower for the cake taste?
➢ Optimal process output (Response)
  ➢ What are the "best parameters" to have the tastiest cake.
➢ Reducing variability (sensitivity analysis)
  ➢ Which parameters can change slightly without changing the cake taste?
➢ Target an output
  ➢ Which parameters for, best taste, best color, best consistency?
➢ Balancing tradeoffs
  ➢ How to define optimum parameters to have the best taste at the best price

Nice DoE tutorial and "package " :
I advise the R package RcmdrPlugin.DoE

**Tutorial for designing experiments using the R package RcmdrPlugin.DoE**

**Tutorial for designing experiments using the R package RcmdrPlugin.DoE**

**Tutorial for designing experiments using the R package RcmdrPlugin.DoE**

**Tutorial for designing experiments using the R package RcmdrPlugin.DoE**

**Tutorial for designing experiments using the R package RcmdrPlugin.DoE**



Main effects plot for mean
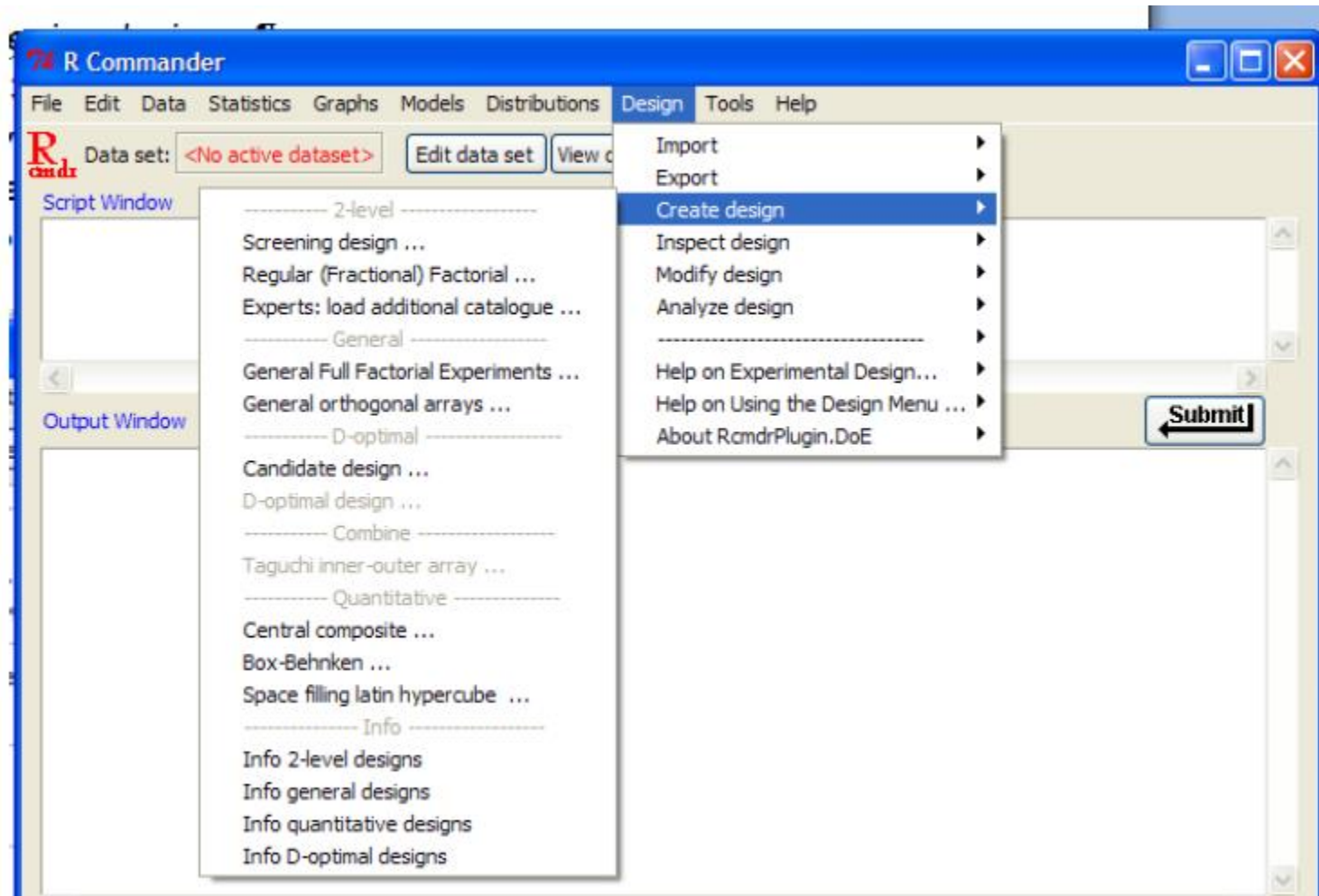
**Effects plots with model uncertainty**

## Interaction plots

- Comparing alternative
  - Egg from different supplier gives the same cake?
- Significant inputs (Factors)
  - Sort the relevance of factors: is the amount of sugar more important than the amount of flower for the cake taste?
- Optimal process output (Response)
  - What are the "best parameters" to have the tastiest cake.
- Reducing variability (sensitivity analysis)
  - Which parameters can change slightly without changing the cake taste?
- Target an output
  - Which parameters for, best taste, best color, best consistency?
- Balancing tradeoffs
  - How to define optimum parameters to have the best taste at the best price

## 3a) Soft modelling versus hard modelling

## 3b) Absorption spectroscopy

- Beer's law in elegant matrix notation ($\mathbf{Y} = \mathbf{C} \times \mathbf{A}$)
- Non-unique factorisation of $\mathbf{Y}$ / rotational ambiguity

## 3c) Principal Component Analysis (**PCA**)

- Abstract Factor analysis (**AFA**) by Singular Value Decomposition (**SVD**)
- Chemical rank of the measurement matrix
- The number of absorbing species

## 3d) Evolving Factor Analysis (**EFA**)

- Evolutionary rank analysis by repeated SVD of sub matrices of $\mathbf{Y}$
- The 'Appearance' & 'Disappearance' of absorbing species

## 3e) Multivariate Curve Resolution by Alternating Least-Squares (**MCR**-**ALS**)

- Model-free iterative decomposition of $\mathbf{Y} = \mathbf{C} \times \mathbf{A} + \mathbf{R}$
- Ideas, principles, limitations

**Factor model**

**PCA, EFA, MCR-ALS**

**Soft**

**Number of species**
**Approx. conc. & spectra**

**Lambert-Beer**

*Kinetic model*

$$v = k \cdot c_{\text{BuOH}} \cdot c_{\text{AA}} \cdot c_{\text{TMG}}$$

$$\frac{dc_{\text{BuOH}}}{dt} = -v + \frac{f_{in}}{V} \cdot (c_{in,\text{BuOH}} - c_{\text{BuOH}})$$

$$\frac{dc_{\text{AA}}}{dt} = -v - \frac{f_{in}}{V} \cdot c_{\text{AA}}$$

$$\frac{dc_{\text{BuOA}}}{dt} = v - \frac{f_{in}}{V} \cdot c_{\text{BuOA}}$$

$$\frac{dc_{\text{AH}}}{dt} = v - \frac{f_{in}}{V} \cdot c_{\text{AH}}$$

$$\frac{dc_{\text{TMG}}}{dt} = -\frac{f_{in}}{V} \cdot c_{\text{TMG}}$$

$$\frac{dV}{dt} = f_{in}$$

*Reaction Spectra*

*Concentrations*

*Species Spectra*



=

×

**Hard**

**Rate constants**
**Activation energies**

*+ Residuals*

Soft modelling: first insight into the experimental data

**Absorption of UV-vis and IR light:**

input light intensity $I_{0,\lambda}$    cuvette    output light intensity $I_\lambda$

absorbance $y$ at wavelength λ:

$$y_\lambda = -\log\left(I/I_0\right)_\lambda$$

$l$

path length

- absorbance signal $y_\lambda$ [no unit] is linearly dependent on the concentrations $c_s$ [molL$^{-1}$] of $s=1\ldots n_s$ absorbing species, the corresponding coefficients are the molar absorptivities $a_{s,\lambda}$ [Lmol$^{-1}$cm$^{-1}$] and form the pure species spectra

**Beer's Law:**

$$y_\lambda = \left(c_1 a_{1,\lambda} + \ldots + c_s a_{s,\lambda} + \ldots + c_{n_s} a_{n_s,\lambda}\right) \times l$$

$$= \sum_{s=1}^{n_s} c_s a_{s,\lambda} \times l$$

Often:
path length $l = 1$cm

● But it is time to introduce some useful conventions for a scalar/vector/matrix notation:

  – matrices are given in **boldface capital** characters, e.g. $\mathbf{Y}$

  – vectors are given in **boldface lower case** characters, e.g. $\mathbf{y}$, this includes row or column vectors of a matrix $\mathbf{Y}$, e.g. $\mathbf{y}_{2,:}$ (2nd row of $\mathbf{Y}$), or $\mathbf{y}_{:,3}$ (3rd column of $\mathbf{Y}$)

  – scalars are given in *italic* characters, this includes the elements of a vector $\mathbf{y}$ or matrix $\mathbf{Y}$, e.g. $y_i$ or $y_{i,j}$

● In some cases, it is very illustrative to write vector/matrix equations in a "line/box" notation, e.g.

$$\left| \quad \right| = \left| \; \square \; \right| \; \left| + \right|$$

$$\mathbf{y} \;=\; \mathbf{F} \times \mathbf{a} + \mathbf{r}$$

- For time and/or wavelength resolved kinetic data, Beer's law can be written in elegant vector (single wavelength) or matrix notation (multi wavelengths):



at $i=1\ldots n_t$ times

$$y_i = \sum_{s=1}^{n_s} c_{i,s} a_s + r_i$$

at $j=1\ldots n_\lambda$ wavelengths

$$y_{i,j} = \sum_{s=1}^{n_s} c_{i,s} a_{s,j} + r_{i,j}$$

for example: a reaction $A \rightarrow B \rightarrow C$

e.g. A→B→C



$$\mathbf{Y} \quad = \quad \mathbf{C} \quad \times \quad \mathbf{A} \quad + \quad \mathbf{R}$$

**Goal:** Find concentration profiles **C** and species spectra **A** such that the residuals **R**=**Y**-**CA** become small only using a 'soft model', i.e. by linear factorisation

**Problem:** Factorisation is not unique (rotational ambiguity)

- By using appropriate 'soft' restrictions on **C** and **A**, e.g. non-negativity, windows of existence, closure, unimodality, known spectra, the number of possible solutions can be reduced, sometimes can even lead to a unique solution for **C** & **A**
- There are 2 major classes

1) Factor Analysis (AFA) based

$$\mathbf{Y} = \overline{\mathbf{U}}\,\overline{\mathbf{S}}\,\overline{\mathbf{V}} = \overline{\mathbf{U}}\mathbf{T}\mathbf{T}^{-1}\overline{\mathbf{S}}\,\overline{\mathbf{V}} = \mathbf{C}\mathbf{A}$$

Find **T** such that

$$\mathbf{C} = \overline{\mathbf{U}}\mathbf{T}, \text{ and } \mathbf{A} = \mathbf{T}^{-1}\overline{\mathbf{S}}\,\overline{\mathbf{V}}$$

2) Alternating Least Squares (ALS) based

Start from some guessed **C**,

then recalculate **A** and **C** until satisfied:

$$\mathbf{A} = (\mathbf{C}^t\mathbf{C})^{-1}\mathbf{C}^t\mathbf{Y} = \mathbf{C}^+\mathbf{Y}$$

$$\mathbf{C} = \mathbf{Y}\mathbf{A}^t(\mathbf{A}\mathbf{A}^t)^{-1} = \mathbf{Y}\mathbf{A}^+$$

## 3a) Soft modelling versus hard modelling

## 3b) Absorption spectroscopy

- Beer's law in elegant matrix notation ($Y = C \times A$)
- Non-unique factorisation of $Y$ / rotational ambiguity

## 3c) Principal Component Analysis (PCA)

- Abstract Factor analysis (AFA) by Singular Value Decomposition (SVD)
- Chemical rank of the measurement matrix
- The number of absorbing species

## 3d) Evolving Factor Analysis (EFA)

- Evolutionary rank analysis by repeated SVD of sub matrices of $Y$
- The 'Appearance' & 'Disappearance' of absorbing species

## 3e) Multivariate Curve Resolution by Alternating Least-Squares (MCR-ALS)

- Model-free iterative decomposition of $Y = C \times A + R$
- Ideas, principles, limitations

- One very well defined solution is the one received from Abstract Factor Analysis (AFA) using Singular Value Decomposition (SVD)



$$X = USV^T$$

**S** is a diagonal matrix with the square root of their eigenvalues

columns of **U** (rows of **V**) are eigenvectors of $YY^t$ ($Y^tY$)

**U** and $V^t$ are orthonormal

in Matlab:

```
[U,S,Vt]=svd(Y,0);
```

Indeed very efficient tool for data compression and noise filtration (orthogonal noise is removed)

- Eigenvectors in **U** (columns) and **V** (rows) are arranged in decreasing order of magnitude of their corresponding singular values in **S**

- Many of them just represent 'noise' and can be neglected; the significant 'factors', the Principal Components, are retained in $\bar{\mathbf{U}}$ and $\bar{\mathbf{V}}$ and form 'abstract' concentration profiles and spectra

- The diagonal elements of $\bar{\bar{\mathbf{S}}}$, the singular values, can be seen as normalisation coefficients for $\bar{\mathbf{U}}$ or $\bar{\mathbf{V}}$

$$N_t \begin{bmatrix} \bar{\mathbf{Y}} \\ {}^{N_\lambda} \end{bmatrix} = \bar{\mathbf{U}} \,|\, \text{noise} \quad \begin{bmatrix} \bar{\bar{\mathbf{S}}} \\ \text{noise} \end{bmatrix} \begin{bmatrix} \bar{\mathbf{V}} \\ \text{noise} \end{bmatrix} N_e = \bar{\mathbf{U}}\bar{\bar{\mathbf{S}}} \,|\, \bar{\mathbf{V}} \,|\, N_e \approx \mathbf{C} \,|\, \mathbf{A}$$

**U**     **S**     **V**

• The number of significant singular(eigen) values and –vectors is the chemical rank of **Y** and a 1st estimate on the number of absorbing species

e.g. A→B→C

$\mathbf{Y}$

$= \mathbf{C}$

$\mathbf{A}$

$+ \mathbf{R}$

$\overline{\mathbf{U}} \quad \overline{\mathbf{S}}$

$\overline{\mathbf{V}}$

$+$

$\mathbf{R}_{PCA}$

→ rank($\mathbf{Y}$) = 3
→ 3 absorbing species

The **rank** of a matrix *A* is the size of the largest collection of linearly independent columns or rows of *A*.

$\mathbf{Y}$ ($N_t \times N_\lambda$)                    $\log(s_{i,i})$ vs $i$



The noise level in the data matrix **Y** determines the drop in the magnitude from significant to insignificant singular values

## 3a) Soft modelling versus hard modelling

## 3b) Absorption spectroscopy

- Beer's law in elegant matrix notation ($\mathbf{Y} = \mathbf{C} \times \mathbf{A}$)
- Non-unique factorisation of $\mathbf{Y}$ / rotational ambiguity

## 3c) Principal Component Analysis (**PCA**)

- Abstract Factor analysis (**AFA**) by Singular Value Decomposition (**SVD**)
- Chemical rank of the measurement matrix
- The number of absorbing species
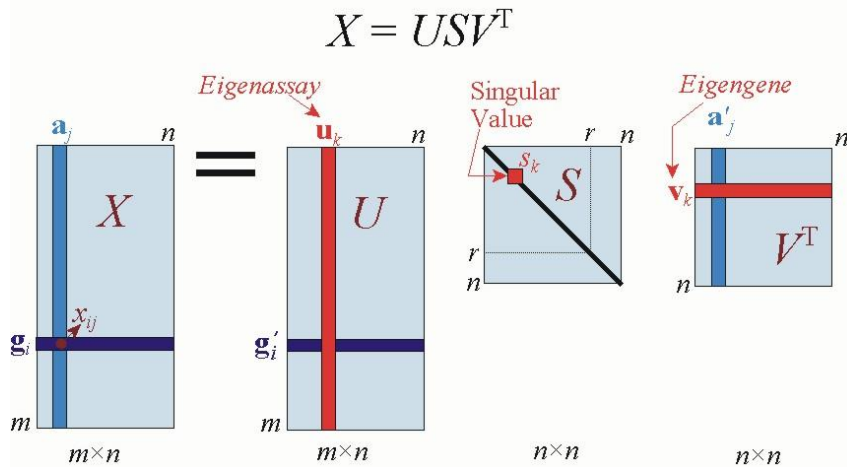
## 3d) Evolving Factor Analysis (**EFA**)

- Evolutionary rank analysis by repeated SVD of sub matrices of $\mathbf{Y}$
- The 'Appearance' & 'Disappearance' of absorbing species

## 3e) Multivariate Curve Resolution by Alternating Least-Squares (**MCR**-**ALS**)

- Model-free iterative decomposition of $\mathbf{Y} = \mathbf{C} \times \mathbf{A} + \mathbf{R}$
- Ideas, principles, limitations

- the chemical rank of the spectral data matrix **Y** is determined by the number of its significant singular vectors

- the number of significant singular vectors of **Y** is determined by the number of linearly independent columns or rows in the matrix of pure species spectra (**A**) and corresponding concentration profiles (**C**)

- linear dependencies in **C** due to the kinetic model are common and sometimes difficult to predict (e.g. $A+B \rightarrow C$)

- linear dependencies in **A** are less common

- **sequential rank analysis** of the data matrix along its time domain by repeated SVD
- can be performed in a **forward and backward** way
- indicates the **rise of new singular vectors** and thus gives an estimate for the appearance & disappearance of new absorbing species
- **ideally** designed to follow **chromatography** experiments
  - species appear & disappear **sequentially**
- capable of **roughly** following **kinetic** profiles
  - species can appear & dissappear **simultaneously**

Forward EFA



SVD $\bar{\mathbf{s}}_{i,:}$

$N_e$

$1$

$i$

$N_t$

• Repeated rank analysis by SVD in forward direction
• The appearance of a new 'species' is indicated by a gradual rise of a new singular value

Backward EFA



• Repeated rank analysis by SVD in backward direction

• A 'disappearing species' is indicated by a gradual rise of a new singular value

Forward and backward EFA



Combined forward/backward EFA results can be used as reasonable initial guesses of concentration profiles for subsequent iterative refinement e.g. by ALS

## 3a) Soft modelling versus hard modelling

## 3b) Absorption spectroscopy

- Beer's law in elegant matrix notation ($\mathbf{Y} = \mathbf{C} \times \mathbf{A}$)
- Non-unique factorisation of $\mathbf{Y}$ / rotational ambiguity

## 3c) Principal Component Analysis (**PCA**)

- Abstract Factor analysis (**AFA**) by Singular Value Decomposition (**SVD**)
- Chemical rank of the measurement matrix
- The number of absorbing species

## 3d) Evolving Factor Analysis (**EFA**)

- Evolutionary rank analysis by repeated SVD of sub matrices of $\mathbf{Y}$
- The 'Appearance' & 'Disappearance' of absorbing species

## 3e) Multivariate Curve Resolution by Alternating Least-Squares (**MCR**-**ALS**)

- Model-free iterative decomposition of $\mathbf{Y} = \mathbf{C} \times \mathbf{A} + \mathbf{R}$
- Ideas, principles, limitations

- Conceptually very simple



http://www.mcrals.info/

http://www.mcrals.info/

http://www.mcrals.info/

- **Advantages**
  - No prior knowledge on the chemical system required
  - Estimation of the number of linearly dependent absorbing species and their approximate evolution from PCA, EFA & ALS
  - Info for the development of a 'hard' model
  - 'Better than nothing'

- **Drawbacks**
  - No physical model
  - No predictions for other exp. conditions possible
  - Uniqueness of the result is rarely given and difficult to validate

MCR-ALS is a very nice software, I strongly recommend it. http://www.mcrals.info

MCR-ALS becomes a very powerful method when multiple datasets are used.

**Factor model**

PCA, EFA, MCR-ALS

**Soft** — Number of species / Approx. conc. & spectra

Lambert-Beer

*Kinetic model*

$$\upsilon = k \cdot c_{BuOH} \cdot c_{AA} \cdot c_{TMG}$$

$$\frac{dc_{BuOH}}{dt} = -\upsilon + \frac{f_{in}}{V} \cdot (c_{in,BuOH} - c_{BuOH})$$

$$\frac{dc_{AA}}{dt} = -\upsilon - \frac{f_{in}}{V} \cdot c_{AA}$$

$$\frac{dc_{BuOA}}{dt} = \upsilon - \frac{f_{in}}{V} \cdot c_{BuOA}$$

$$\frac{dc_{AH}}{dt} = \upsilon - \frac{f_{in}}{V} \cdot c_{AH}$$

$$\frac{dc_{TMG}}{dt} = -\frac{f_{in}}{V} \cdot c_{TMG}$$

$$\frac{dV}{dt} = f_{in}$$

*Reaction Spectra* = *Concentrations* × *Species Spectra*

**Hard** — Rate constants / Activation energies

+ *Residuals*

Hard modelling: Use a mathematical model to "explain" the experimental data

## The objective of this section

– to understand the underlying ideas and principles of the calibration-free modelling (fitting) of spectro-kinetic absorbance data

– to apply these principles to kinetic problems using Matlab (or R)

$$Y_{Measured}$$

4c) $$\min_{k} \left\| Y_{Measured} - C_{modelled} A_{modelled} \right\|^2$$

4b) $$A_{modelled} \qquad A_{modelled} = C_{modelled}^{+} Y_{Measured}$$

4a) $$C_{Modelled}$$

$$
\begin{cases}
\frac{dC_{Ac_2O}}{dt} &= r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,Ac_2O} - C_{Ac_2O}) \\
\frac{dC_{H_2O}}{dt} &= r(k,t) - \frac{f_{dos}}{V_r(t)} \cdot C_{H_2O} \\
\frac{dC_{AcOH}}{dt} &= -2 \cdot r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,AcOH} - C_{AcOH}) \\
\frac{dV_r}{dt} &= f_{dos}
\end{cases}
$$

$$C = f \text{ (kinetic model, parameters)}$$

rate law

- initial concentrations
- rate constants
- activation energies
- temperature

The rate of a molecular reaction is defined by the derivative of the concentration of the reactants with respect to time normalised by the corresponding stoichiometric coefficient

e.g.
$$aA + bB \xrightarrow{k} cC$$

$$rate = -\frac{d[A]_t}{dt \cdot a} = -\frac{d[B]_t}{dt \cdot b} = \frac{d[C]_t}{dt \cdot c} = k[A]_t^a[B]_t^b$$

$$A + B \xrightarrow{k_1} C$$

$$2C \xrightarrow{k_2} D$$

*Rate laws*

$$rate_1 = k_1 \cdot [A]_t \cdot [B]_t$$

$$rate_2 = k_2 \cdot [C]_t^2$$

$k_1,\ k_2 :\ rate\ constants\ [L/mol/s]$

$[A]_t,\ [B]_t,\ [C]_t :\ conc.\ [mol/L]\ of\ species\ at\ time\ t$

*Derivatives*

$$\frac{d[A]_t}{dt} = \frac{d[B]_t}{dt} = -rate_1$$

$$\frac{d[C]_t}{dt} = rate_1 - 2rate_2$$

$$\frac{d[D]_t}{dt} = rate_2$$

Except for some very simple cases there is no explicit solution for systems of kinetic ODEs ➜ **Numerical integration**

$$A + B \xrightarrow{k_1} C$$

$$2C \xrightarrow{k_2} D$$

*When species are dosed, the set of ODEs can be modified accordingly :*

**Rate laws**

$$rate_1 = k_1 \cdot \left[ A \right]_t \cdot \left[ B \right]_t$$

$$rate_2 = k_2 \cdot \left[ C \right]_t^2$$

**Derivatives**

$$\frac{d\left[ A \right]_t}{dt} = -rate_1 + \frac{F_t}{V_t} \cdot \left( \left[ A \right]^{dos} - \left[ A \right]_t \right)$$

$$\frac{d\left[ B \right]_t}{dt} = -rate_1 + \frac{F_t}{V_t} \cdot \left( \left[ B \right]^{dos} - \left[ B \right]_t \right)$$

$$\frac{d\left[ C \right]_t}{dt} = rate_1 - 2 \cdot rate_2 + \frac{F_t}{V_t} \cdot \left( \left[ C \right]^{dos} - \left[ C \right]_t \right)$$

$$\frac{d\left[ D \right]_t}{dt} = rate_2 + \frac{F_t}{V_t} \cdot \left( \left[ D \right]^{dos} - \left[ D \right]_t \right)$$

$$\frac{dV_t}{dt} = F_t$$

$F_t$ : *flow rate* $[mL / s]$ *at time t*

$V_t$ : *volume* $[mL]$ *at time t*

$[A]^{dos}$ : *conc.* $[mol / L]$ *of dosed species*

*dosed material:* $\quad + \dfrac{F_t}{V_t} \cdot [A]^{dos}$

*dilution:* $\quad - \dfrac{F_t}{V_t} \cdot [A]_t$

Dosing requires some modifications to the ODEs for all species and the inclusion of an ODE for the change of volume due to the flow-rate

- Crude approach : Euler's method (truncated Taylor series)

$$c_i\left(t + \Delta t\right) \approx c_i\left(t\right) + \left(\frac{dc_i}{dt}\right)_t \cdot \Delta t$$

Applied to our specific example *without dosing*

$$\left[C\right]_{t+\Delta t} \approx \left[C\right]_t + \left(\frac{d\left[C\right]_t}{dt}\right)_t \cdot \Delta t$$

$$\approx \left[C\right]_t + \left(k_1 \cdot \left[A\right]_t \cdot \left[B\right]_t - 2k_2 \cdot \left[C\right]_t^2\right)\Delta t$$

$$A + B \xrightarrow{k_1} C$$

$$2C \xrightarrow{k_2} D$$

- Much more sophisticated integration methods exist, e.g. Matlab's '`ode45`', a 4th order Runge-Kutta with *automatic stepsize control*

- In stepsize controlled ODE solvers, the stepsize is adjusted at each step to meet the user-specified accuracy

- The accuracy is measured with absolute (Matlab: `AbsTol`) and relative (Matlab: `RelTol`) tolerance's values.

- For some kinetic models, the concentration profiles change on dramatically different scales (stiff problem) and a **stiff ODE solver** (eg. Matlab's `'ode15s'`) is required

```
% Mechanism : A+B>C, 2C>D (batch)
t         = [0:0.01:2]';          % time vector (column)
c0        = [1 1 0 0];            % initial concentrations (row)
k         = [10; 5];              % rate constants (column)

% call ODE solver for numerical integration of dC/dt
odeoptions = odeset('RelTol',1e-10,'AbsTol',1e-12);
[tout,C] = ode45('ode_ApBtoC_2CtoD_batch',t,c0',odeoptions,k);
```

$$A + B \xrightarrow{k_1} C$$

$$2C \xrightarrow{k_2} D$$

batch conditions

```
function cdot = ode_ApBtoC_2CtoD_batch(ti,ci,flag,k)
% Mechanism : A+B>C, 2C>D (batch)

% rate law
rates(1)   = k(1)*ci(1)*ci(2);
rates(2)   = k(2)*ci(3)*ci(3);

% System of ODEs
cdot(1,1)  = -rates(1);                % dA/dt
cdot(2,1)  = -rates(1);                % dB/dt
cdot(3,1)  =  rates(1) - 2*rates(2);   % dC/dt
cdot(4,1)  =  rates(2);                % dD/dt
```

```
% Mechanism : A+B>C, 2C>D (semi batch)
t      = [0:0.01:4]'; % time vector (column)
c0     = [0 1 0 0];   % initial concentrations (row)
k      = [10;5];      % rate constants (column)
V0     =  0.020;      % initial volume [L]
tdos   = [0.5 1];     % start and stop times of dosing (row)
cdos   = [3 0 0 0];   % concentrations of dosed species (row)
frate  =  0.01;       % flow-rate [L/(unit of time)]
F = flow(t,tdos,frate);   % generate vector F of flow-rates

% call ODE solver for numerical integration of dC/t and dV/t
odeoptions = odeset('RelTol',1e-10,'AbsTol',1e-12);
[tout,CV]  = ode45('ode_ApBtoC_2CtoD_semi',t,[c0 V0]',odeoptions,k,F,t,cdos);
% extract volume and concentration profiles from columns of matrix CV
C = CV(:,1:end-1); V = CV(:,end);
```



```
function cdot = ode_ApBtoC_2CtoD_semi(ti,civi,flag,k,F,t,cdos)
% Mechanism : A+B>C, 2C>D (semi batch)
ci       = civi(1:end-1);            % concentration at time ti
Vi       = civi(end);                % volume at time ti
fi       = F(find(t<=ti,1,'last')); % flow-rate at time ti
% rate law
rates(1)  = k(1)*ci(1)*ci(2);
rates(2)  = k(2)*ci(3)*ci(3);
% System of ODEs
cdot(1,1) = -rates(1)              + (fi/Vi)*(cdos(1)-ci(1)); % dA/dt
cdot(2,1) = -rates(1)              + (fi/Vi)*(cdos(2)-ci(2)); % dB/dt
cdot(3,1) =  rates(1) - 2*rates(2) + (fi/Vi)*(cdos(3)-ci(3)); % dC/dt
cdot(4,1) =  rates(2)              + (fi/Vi)*(cdos(4)-ci(4)); % dD/dt
cdot(5,1) =  fi;                                              % dV/dt
```

$$A + B \xrightarrow{k_1} C$$

$$2C \xrightarrow{k_2} D$$

semi-batch
(A dosed)

## The objective of this section

– to understand the underlying ideas and principles of the calibration-free modelling (fitting) of spectro-kinetic absorbance data

– to apply these principles to kinetic problems using Matlab (or R)



4c) $\min_{k} \left\| \mathbf{Y_{Measured}} - \mathbf{C_{modelled}A_{modelled}} \right\|^2$

4b) $\mathbf{A_{modelled}}$ ← $\mathbf{A_{modelled}} = \mathbf{C_{modelled}^+ Y_{Measured}}$

4a) $\mathbf{C_{Modelled}}$ ←

$$\begin{cases} \frac{dC_{Ac_2O}}{dt} &= r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,Ac_2O} - C_{Ac_2O}) \\ \frac{dC_{H_2O}}{dt} &= r(k,t) - \frac{f_{dos}}{V_r(t)} \cdot C_{H_2O} \\ \frac{dC_{AcOH}}{dt} &= -2 \cdot r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,AcOH} - C_{AcOH}) \\ \frac{dV_r}{dt} &= f_{dos} \end{cases}$$

Decision tree for the selection of an appropriate method for the kinetic hard-modelling of spectroscopic data

**Target:** Find the least-squares minimum as a function of the rate constant(s), the non-linear parameters

for 'concentration measurements', $\mathbf{C}$

$$ssq(\mathbf{k}) = \sum_{i=1}^{n_t} \sum_{s=1}^{n_s} r_{i,s}^2(\mathbf{k}), \qquad \mathbf{R}(\mathbf{k}) = \mathbf{C} - \mathbf{C_{calc}}(\mathbf{k})$$

for spectral absorbance measurements, $\mathbf{Y}$

$$ssq(\mathbf{k}) = \sum_{i=1}^{n_t} \sum_{j=1}^{n_\lambda} r_{i,j}^2(\mathbf{k}), \qquad \mathbf{R}(\mathbf{k}) = \mathbf{Y} - \mathbf{Y_{calc}}(\mathbf{k}) = \mathbf{Y} - \mathbf{C_{calc}}(\mathbf{k}) \cdot \mathbf{A_{calc}}(\mathbf{k})$$

$$\partial ssq / \partial \mathbf{k} = 0$$

No explicit solution !!!
→ problem to be solved iteratively

- The most widely used data fitting is the linear regression of a data vector **y** to a straight line (1st order polynomial)



$$y_i = a_1 + a_2 x_i + r_i$$

$$i = 1 \ldots m, \; m > 2$$

● As with all other regression methods the task is to minimise the least-squares sum $ssq$ of all residuals $r_i$ between the data $y_i$ and the underlying model $a_1 + a_2 x_i$ by optimising the linear parameters defining the model, here slope $a_2$ and intercept $a_1$

$$\min_{a_1, a_2} |ssq| \qquad \text{where} \qquad ssq(a_1, a_2) = \sum_{i=1}^{m} r_i^2 = \sum_{i=1}^{m} \left( \underbrace{y_i - (\underbrace{a_1 + a_2 x_i}_{y_{calc_i}})} \right)^2$$

- In order to find the "best" parameters $a_1$ & $a_2$, i.e. that lead to a minimal $ssq$ the following two derivatives must be zero:

$$\frac{\partial ssq}{\partial a_1} = \frac{\partial ssq}{\partial a_2} = 0$$

- For linear parameters, such as slope and intercept there is a non-iterative solution, i.e. for $m$ data points, $a_1$ and $a_2$ can be calculated explicitly:

$$a_1 = \frac{\Sigma x_i^2 \Sigma y_i - \Sigma x_i \Sigma x_i y_i}{m \Sigma x_i^2 - (\Sigma x_i)^2}$$

$$a_2 = \frac{-\Sigma x_i \Sigma y_i + m \Sigma x_i y_i}{m \Sigma x_i^2 - (\Sigma x_i)^2}$$

- But linear regression is much more than just straight line fitting. It also includes the fitting to polynomials of any other order or any other function that depends on linear parameters only.

- Linear relationships can be written much simpler in an elegant vector or matrix notation. The straight line regression problem is then denoted by:

$$y_1 = a_1 + a_2 x_1 + r_1$$
$$y_2 = a_1 + a_2 x_2 + r_2$$
$$\vdots$$
$$y_i = a_1 + a_2 x_i + r_i$$
$$\vdots$$
$$y_m = a_1 + a_2 x_m + r_m$$

$$\longrightarrow \quad \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_i \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} + \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_i \\ \vdots \\ r_m \end{bmatrix}$$

$$\mathbf{y} = \mathbf{F}(\mathbf{x}) \cdot \mathbf{a} + \mathbf{r}$$

$\mathbf{F}(\mathbf{x})$ is the design matrix

- And the generalisation for any polynomial

$$y_i = a_1 + a_2 x_i + a_3 x_i^2 + \ldots + a_j x_i^{j-1} + \ldots + a_{n_a} x_i^{n_a-1} = \sum_{j=1}^{n_a} a_j x_i^{j-1}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{j-1} & \cdots & x_1^{n_a-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{j-1} & \cdots & x_2^{n_a-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \cdots & \vdots \\ 1 & x_i & x_i^2 & \cdots & x_i^{j-1} & \cdots & x_i^{n_a-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \cdots & x_m^{j-1} & \cdots & x_m^{n_a-1} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_j \\ \vdots \\ a_{n_a} \end{bmatrix} + \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_i \\ \vdots \\ r_m \end{bmatrix}$$

$$\mathbf{y} = \mathbf{F(x)} \quad \mathbf{a} + \mathbf{r}$$

*Calibration with pure spectra and fitting* **C**:

● Analogously, if the pure species spectra **A** are known the corresponding concentration profiles **C** can also be determined by **multivariate linear regression**:

$$\mathbf{Y} = \mathbf{C} \cdot \mathbf{A} + \mathbf{R}$$

such that $ssq = \sum\limits_{i=1}^{n_t} \sum\limits_{j=1}^{n_\lambda} r_{i,j}^2$ is minimal

$$\mathbf{C} = \mathbf{Y}\mathbf{A}^t(\mathbf{A}\mathbf{A}^t)^{-1}$$

$\mathbf{A}^t(\mathbf{A}\mathbf{A}^t)^{-1}$ is called the right pseudoinverse, $\mathbf{A}^+$, of $\mathbf{A}$

*Calibration free and fitting* **Y**:

• For a given multi wavelength kinetic measurement, **Y**, if the concentration profiles **C** are known the corresponding pure species spectra **A** can be determined by **multivariate linear regression**:

$$\mathbf{Y} = \mathbf{C} \cdot \mathbf{A} + \mathbf{R}$$

such that $ssq = \sum\limits_{i=1}^{n_t} \sum\limits_{j=1}^{n_\lambda} r_{i,j}^2$ is minimal

$$\mathbf{A} = (\mathbf{C}^t\mathbf{C})^{-1}\mathbf{C}^t\mathbf{Y}$$

$(\mathbf{C}^t\mathbf{C})^{-1}\mathbf{C}^t$ is called the left pseudoinverse, $\mathbf{C}^+$, of $\mathbf{C}$

*Beer's law in matrix notation*



**With calibration**

$$A_{Calibration} = \text{function }(\boldsymbol{\theta})$$

$$\mathbf{C} = (\mathbf{Y}\mathbf{A}^{+}_{calibration})$$

**Calibration free**

$$\mathbf{C}_{modelled} = \text{function }(\mathbf{k})$$

$$\mathbf{A} = \mathbf{C}^{+}_{modelled}\mathbf{Y} \qquad [1]$$

*Objective function:* $\quad \min_{k} \left\| \mathbf{C}_{Measured} - \mathbf{C}_{modelled} \right\|^{2} \quad \min_{k} \left\| \mathbf{Y}_{Measured} - \mathbf{C}_{modelled}\mathbf{A} \right\|^{2}$

$$[1] \ \mathbf{C}^{+} = \left(\mathbf{C}^{T}\mathbf{C}\right)^{-1}\mathbf{C}^{T}$$

## The objective of this section

– to understand the underlying ideas and principles of the calibration-free modelling (fitting) of spectro-kinetic absorbance data

– to apply these principles to kinetic problems using Matlab (or R)

4c) $$\min_{k} \left\| \mathbf{Y}_{\textbf{Measured}} - \mathbf{C}_{\textbf{modelled}}\mathbf{A}_{\textbf{modelled}} \right\|^{2}$$

$\mathbf{Y}_{\textbf{Measured}}$

4b) $\mathbf{A}_{\textbf{modelled}}$ $\quad\longleftarrow\quad$ $\mathbf{A}_{\textbf{modelled}} = \mathbf{C}_{\textbf{modelled}}^{+}\mathbf{Y}_{\textbf{Measured}}$

4a) $\mathbf{C}_{\textbf{Modelled}}$

$$\begin{cases} \frac{dC_{Ac_2O}}{dt} &= r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,Ac_2O} - C_{Ac_2O}) \\ \frac{dC_{H_2O}}{dt} &= r(k,t) - \frac{f_{dos}}{V_r(t)} \cdot C_{H_2O} \\ \frac{dC_{AcOH}}{dt} &= -2 \cdot r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,AcOH} - C_{AcOH}) \\ \frac{dV_r}{dt} &= f_{dos} \end{cases}$$

Sphere:

Eggholder:

Goldstein-Price:

HölderTable:

"Real case":

Easom:

Global optimizer (Heuristic algoritm):

➢ Evolutionary algorithm
➢ Genetic algorithms (GA)
➢ Simulated annealing (SA)
➢ Tabu search
➢ Memetic algorithm
➢ Particle swarm optimization (PSO)

Iterative (gradient based methods)

➢ Sequential quadratic programming (SQP, Hessian ev.)
➢ Quasi-Newton methods
➢ Gradient descent
➢ Simplex
➢ Simultaneous perturbation stochastic approximation (SPSA)

Global optimizer Ex. Particle Swarm Optimisation (PSO)



https://www.youtube.com/watch?v=3CR5y8qZf0Y

Iterative (gradient based methods)



Response surface

Taylor's approximation (parabola)

Minimum of the parabola

Steepest direction

$ssq(\mathbf{k}_0)$

Steepest direction

Newton-Gauss
Steepest descent

$ssq$

Steepest direction

$ssq(\mathbf{k}_0+\Delta\mathbf{k})$

$\mathbf{k}$

Minimum $ssq$
optimum $\mathbf{k}$

- It is possible to circumvent the calculation of **Hessian**

$$ssq\left(\mathbf{k}+\Delta\mathbf{k}\right)\approx ssq\left(\mathbf{k}\right)+\left(\frac{\partial ssq}{\partial\mathbf{k}}\right)\cdot\Delta\mathbf{k}+\frac{1}{2}\left(\frac{\partial^2 ssq}{\partial\mathbf{k}\partial\mathbf{k}}\right)\cdot\Delta\mathbf{k}^2$$

and to develop the Taylor series for the residuals $\mathbf{R}$

- To do so it is convenient to first vectorise the matrices of residuals $\mathbf{R}$ = $\mathbf{C}$-$\mathbf{C}_{calc}$ (fitting concentration data $\mathbf{C}$)
$\mathbf{R}=\mathbf{Y}-\mathbf{Y}_{calc}=\mathbf{Y}-\mathbf{C}_{calc}\mathbf{A}_{calc}$ (fitting spectral data $\mathbf{Y}$)



**VECTORISATION :**
$\mathbf{R}$ is unfolded into a ´long´ column vector $\mathbf{r}$
__in Matlab:__ $\mathbf{r}$=R(:)

- **Computationally easier and (almost) equivalent:** The residuals are approximated by a Taylor series expansion truncated after the first derivative

**Jacobian J**

$$\mathbf{r}(\mathbf{k}+\Delta\mathbf{k}) = \mathbf{r}(\mathbf{k}) + \frac{\partial \mathbf{r}(\mathbf{k})}{\partial \mathbf{k}} \cdot \Delta\mathbf{k}$$

rearrange for $\mathbf{r}(\mathbf{k})$

$$\mathbf{r}(\mathbf{k}) = -\mathbf{J} \cdot \Delta\mathbf{k} + \mathbf{r}(\mathbf{k}+\Delta\mathbf{k})$$

same structure as before in
$\mathbf{y} = \mathbf{C} \times \mathbf{a} + \mathbf{r}$, with $\mathbf{a} = \mathbf{C}^{+}\mathbf{y} = (\mathbf{C}^t\mathbf{C})^{-1}\mathbf{C}^t\mathbf{y}$

*thus: linear regression to minimize* $\mathbf{r}(\mathbf{k}+\Delta\mathbf{k})$ *by optimising* $\Delta\mathbf{k}$ *for a given* $\mathbf{J}$

$$\Delta\mathbf{k} = -(\mathbf{J}^t\mathbf{J})^{-1}\mathbf{J}^t \cdot \mathbf{r}(\mathbf{k}) = -\mathbf{J}^{+} \cdot \mathbf{r}(\mathbf{k})$$

The **SHIFT VECTOR** $\otimes\mathbf{k}$ is added to $\mathbf{k}$ for the next iteration

**approx. Hessian H**

```
Matlab: Δk = -J \ r
```

The Newton-Gauss algorithm
requires the calculation of the Jacobian

$$\mathbf{J} = \frac{\partial \mathbf{r}(\mathbf{k})}{\partial \mathbf{k}}$$

The NG algorithm **CONVERGES** if the Taylor series expansion is a good approximation for the new residuals $\mathbf{r}(\mathbf{k}+\Delta\mathbf{k})$, and the shift vector

$$\Delta\mathbf{k} = -(\mathbf{J}^t\mathbf{J})^{-1}\mathbf{J}^t \cdot \mathbf{r}(\mathbf{k}) = -\mathbf{J}^+ \cdot \mathbf{r}(\mathbf{k})$$

leads to a better (smaller) least-squares sum $ssq(\mathbf{k}+\Delta\mathbf{k})$ than $ssq(\mathbf{k})$ of the previous iteration or with the initial guess of $\mathbf{k}$.

**Convergence criterion:** "small" relative change of $ssq$ (e.g. $\mu = 10^{-4}$)

$$\frac{ssq(\mathbf{k}) - ssq(\mathbf{k}+\Delta\mathbf{k})}{ssq(\mathbf{k})} > \mu \qquad \text{positive,} \rightarrow \text{convergence}$$

$$\frac{ssq(\mathbf{k}) - ssq(\mathbf{k}+\Delta\mathbf{k})}{ssq(\mathbf{k})} < -\mu \qquad \text{negative,} \rightarrow \text{divergence}$$

$$abs\left(\frac{ssq(\mathbf{k}) - ssq(\mathbf{k}+\Delta\mathbf{k})}{ssq(\mathbf{k})}\right) \leq \mu \qquad \text{approx. equal} \rightarrow \text{minimum reached}$$

**PROBLEM** :    The NG algorithm **DIVERGES** if the Taylor series expansion

is not a good approximation for the residuals function

(e.g. with poor initial guesses of $\mathbf{k}$)

**SOLUTION** :    Do not use a Taylor series expansion but move in the
direction of steepest descent

$$\Delta \mathbf{k} = -\mathbf{H}^{-1} \cdot \mathbf{J}^{t} \cdot \mathbf{r}(\mathbf{k})$$

**Inverse Hessian method**
(Newton-Gauss)

*Is there a way to switch progressively
from one method to the other ?*

$$\Delta \mathbf{k} = -\left(\mathbf{H} + mp \cdot \mathbf{I}\right)^{-1} \cdot \mathbf{J}^{t} \cdot \mathbf{r}(\mathbf{k})$$

**Levenberg-Marquardt**
modification

The Marquardt parameter ($mp$) is a scalar added
to the diagonal elements of $\mathbf{H}$ to decrease its
influence on $\otimes\mathbf{k}$ and shorten
the magnitude of $\otimes\mathbf{k}$

$$\Delta \mathbf{k} = -\mathbf{J}^{t} \cdot \mathbf{r}(\mathbf{k})$$

**steepest descent**

- With $\mathbf{H} = \mathbf{J}^t \mathbf{J}$ the shift vector Δ$\mathbf{k}$ can be written as:

$$\Delta\mathbf{k} = -\mathbf{H}^{-1}\mathbf{J}^t\mathbf{r}(\mathbf{k})$$

The Hessian $\mathbf{H}$ is a square matrix ($n_k \times n_k$). It's inverse is an estimate for the variance/covariance matrix for $\mathbf{k}$ !

The diagonal element(s) of the inverted Hessian allow the calculation of the **standard error(s)** $\sigma_\mathbf{k}$ for the rate constant(s) $\mathbf{k}$:

$$\sigma_\mathbf{k} = \sigma_\mathbf{r} \cdot \sqrt{diag\left(\mathbf{H}^{-1}\right)} \quad with \quad \sigma_\mathbf{r} = \sqrt{\frac{ssq}{df}} \approx \sigma_\mathbf{y}$$

$\sigma_\mathbf{r}$: standard deviation of the residuals $\mathbf{r}$ ($\mathbf{R}$)

$\sigma_\mathbf{y}$: ´true´ standard deviation of the measurement $\mathbf{Y}$ (or $\mathbf{C}$)

$df$: degree of freedom, $df = n_t n_s - n_k$ ($\mathbf{C}$ fitted), or $df = n_t n_\lambda - n_k - n_s n_\lambda$ ($\mathbf{Y}$ fitted)

## The objective of this section

– to understand the underlying ideas and principles of the calibration-free modelling (fitting) of spectro-kinetic absorbance data

– to apply these principles to kinetic problems using Matlab (or R)



$$\textbf{Y}_{\textbf{Measured}}$$

4c) $$\min_{k} \left\| \textbf{Y}_{\textbf{Measured}} - \textbf{C}_{\textbf{modelled}}\textbf{A}_{\textbf{modelled}} \right\|^2$$

4b) $$\textbf{A}_{\textbf{modelled}} \qquad \textbf{A}_{\textbf{modelled}} = \textbf{C}^{+}_{\textbf{modelled}}\textbf{Y}_{\textbf{Measured}}$$

4a) $$\textbf{C}_{\textbf{Modelled}}$$

$$\begin{cases} \frac{dC_{Ac_2O}}{dt} &= r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,Ac_2O} - C_{Ac_2O}) \\ \frac{dC_{H_2O}}{dt} &= r(k,t) - \frac{f_{dos}}{V_r(t)} \cdot C_{H_2O} \\ \frac{dC_{AcOH}}{dt} &= -2 \cdot r(k,t) + \frac{f_{dos}}{V_r(t)} \cdot (C_{dos,AcOH} - C_{AcOH}) \\ \frac{dV_r}{dt} &= f_{dos} \end{cases}$$

```
% load spectro-kinetic experimental data (Y,t,w,c0,V0,F,cdos)
load('data_fit_CY.mat');
% Initial guess for the rate constant(s)
k0 = [7; 3]
% Fitting of Y (optimization of k0)
[k, ssq, Hessian, Ccalc, Acalc] =
              nglm_Y('r_Ycalc', k0, t, c0, Y, cdos, V0, F);
% Standard deviation on k
df     = prod(size(Y)) - length(k0) - prod(size(Acalc));
sig_r = sqrt(ssq/df)
sig_k = sig_r*sqrt(diag(inv(Hessian)))
```

$$A + B \xrightarrow{k_1} C \quad \text{semi-batch}$$
$$2C \xrightarrow{k_2} D \quad (A \text{ dosed})$$

Matlab output

```
k0 = 7 3

it=0, k(1)=7,       k(2)=3,       ssq=8.1585
it=1, k(1)=9.51204, k(2)=4.41167, ssq=7.97518
it=2, k(1)=10.404,  k(2)=4.84956, ssq=7.96041
it=3, k(1)=10.463,  k(2)=4.83658, ssq=7.9604

sig_r = 0.01
sig_k = 0.2664  0.1312
```

```
function [r, ssq, Ccalc, Acalc] = r_Ycalc(k, t, c0, Y, cdos, V0, F)
odeoptions  = odeset('RelTol', 1e-10, 'AbsTol', 1e-12);
[tdummy,CV] = ode45('ode_ApBtoC_2CtoD_semi', t, [c0 V0]',
                    odeoptions, k, F, t, cdos);
% Extraction of the concentration profiles
Ccalc       = CV(:,1:end-1);
Vcalc       = CV(:,end);
% Calculation of Acalc
Acalc       = Ccalc\Y;
% Calculation of the residuals
R           = Y - Ccalc*Acalc;
% Vectorization
r           = R(:);
% Calculation of the sum of squares
ssq         = sum(r.^2);
```

elimination of $A$



fitted vs 'measured'
absorbances at selected
wavelengths

- Occam´s razor (*lex parsimoniae*)

  The principle states that among competing hypotheses, the one with the fewest assumptions should be selected. Other, more complicated solutions may ultimately prove correct, but—in the absence of certainty—the fewer assumptions that are made, the better.

- Akaike information criterion (AIC)

  The AIC is a measure of the relative quality of a statistical model for a given set of data. As such, AIC provides a means for model selection.

  AIC = 2k − 2ln(L)

  where $k$ is the number of parameters in the model, and $L$ is the maximized value of the likelihood function for the model.

  Models with small AIC should be preferred.

- Bayesian information criterion (BIC)…and more

- We "only" considered irreversible homogeneous solution chemistry & absorbance measurements but the principles can be extended to also deal with

  - instantaneous & kinetically observable equilibria
  - multiphasic transitions (e.g. surface catalysis, gas formation)
  - change of pH, ionic strength (activities), temperature
  - other data types (e.g. heat, pressure, pH, conductivity, particle size)

- Generally, this "only" requires an adaptation of the differential equations defining the mass transfer and a reassignment of the signal(s) to be fitted and/or parameters to be optimised

Chemometric will support you for:

- How to design/select the reactor/instrumentation?

- How to do the experimental design?

- How to "find" the correct kinetic model?

- What to "fit" and how ?

- How to determine the rate constant?

- How to "fit" if "C" can not be isolated and unknown?

- What about baseline drift, shift, noise level?

- Finally, are my fitted parameters "correct" to which extends?

- And we have to be quick to do all the above tasks

- **Factor Analysis in Chemistry**

  *E.R. Malinowski, 3rd ed., Wiley, New York 2002*

- **Practical Data Analysis in Chemistry**

  *M. Maeder, Y.M. Neuhold, Elsevier, Amsterdam 2007*

- **Practical Guide to Chemometrics**

  *P. Gemperline (editor), 2nd ed., CRC Press, Boca Raton 2006*

- **The Investigation of Organic Reactions and their Mechanisms**

  *H. Maskill (editor), Blackwell Publishing, Oxford 2006*

- **Evolving factor analysis for the resolution of overlapping chromatographic peaks**

  *M. Maeder, Anal. Chem. 59 (1987), 527-530*

- **Nonlinear Least-Squares Fitting of Multivariate Absorption Data**

  *M. Maeder, A. Zuberbühler. Anal. Chem. 62 (1990), 2220-2224*

- **Analyses of 3-way data from equilibrium and kinetic investigations**

  *R. Dyson , M. Maeder, Y.M. Neuhold, G. Puxty. Anal. Chim. Act. 490 (2003), 99-108*

- **Empirical kinetic modelling of on-line simultaneous infrared and calorimetric measurement using a pareto optimal approach and multi-objective genetic algorithm**

  *S.I. Gianoli, G. Puxty, U. Fischer, M. Maeder, K. Hungerbühler. Chemom. Int. Lab. Syst. 85 (2007), 47-62*

- **Tutorial on the Fitting of Kinetic Models to Multivariate Spectroscopic Measurements with Non-Linear Regression**

  *G. Puxty, M. Maeder, K. Hungerbühler. Chemom. Int. Lab. Syst. 81 (2006), 149-164*

- **Data Oriented Process Development: Determination of Reaction Parameters by Small-Scale Calorimetry with in situ Spectroscopy**

  *G. Puxty, U. Fischer, M. Jecklin, K. Hungerbühler. Chimia 60 (2006), 605-610*