# Data Preprocessing
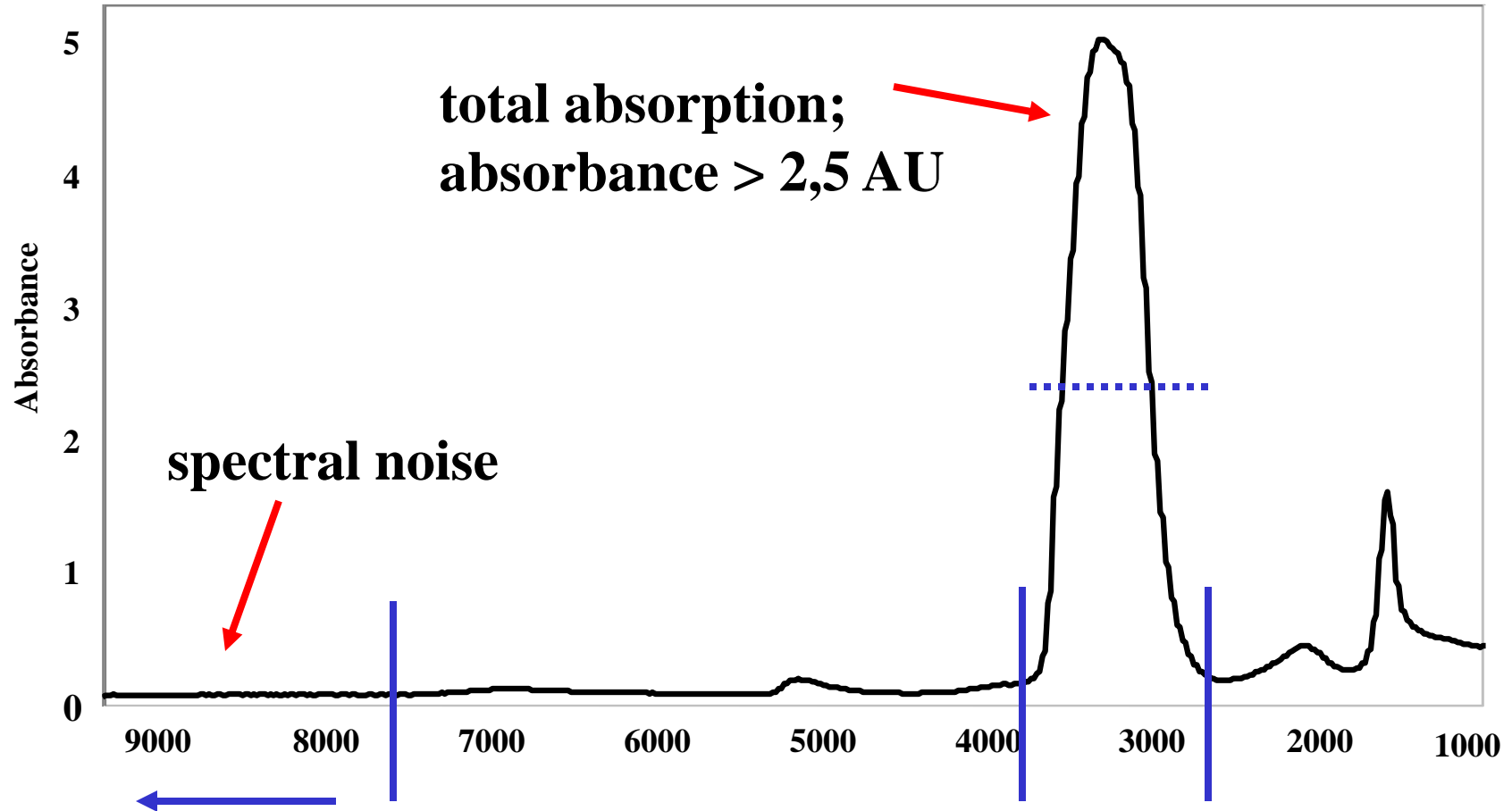
D.N. Rutledge, AgroParisTech

# **Outline**

- Zone selection
- Examining raw data
- The importance of pre-treatment of data

- Common pre-treatment methods

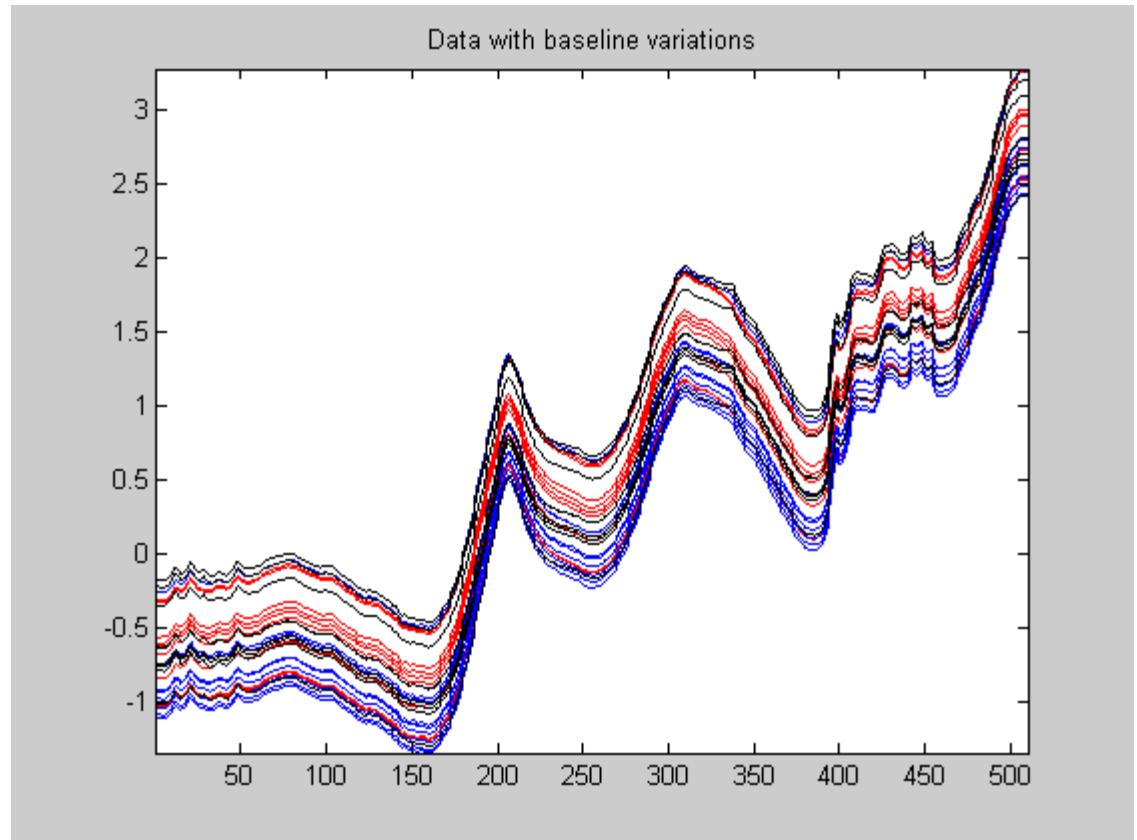# Selection of proper frequency ranges

# Selection of proper frequency ranges



total absorption;
absorbance > 2,5 AU

spectral noise

# Raw data check

The black, red and blue curves indicate different concentration levels
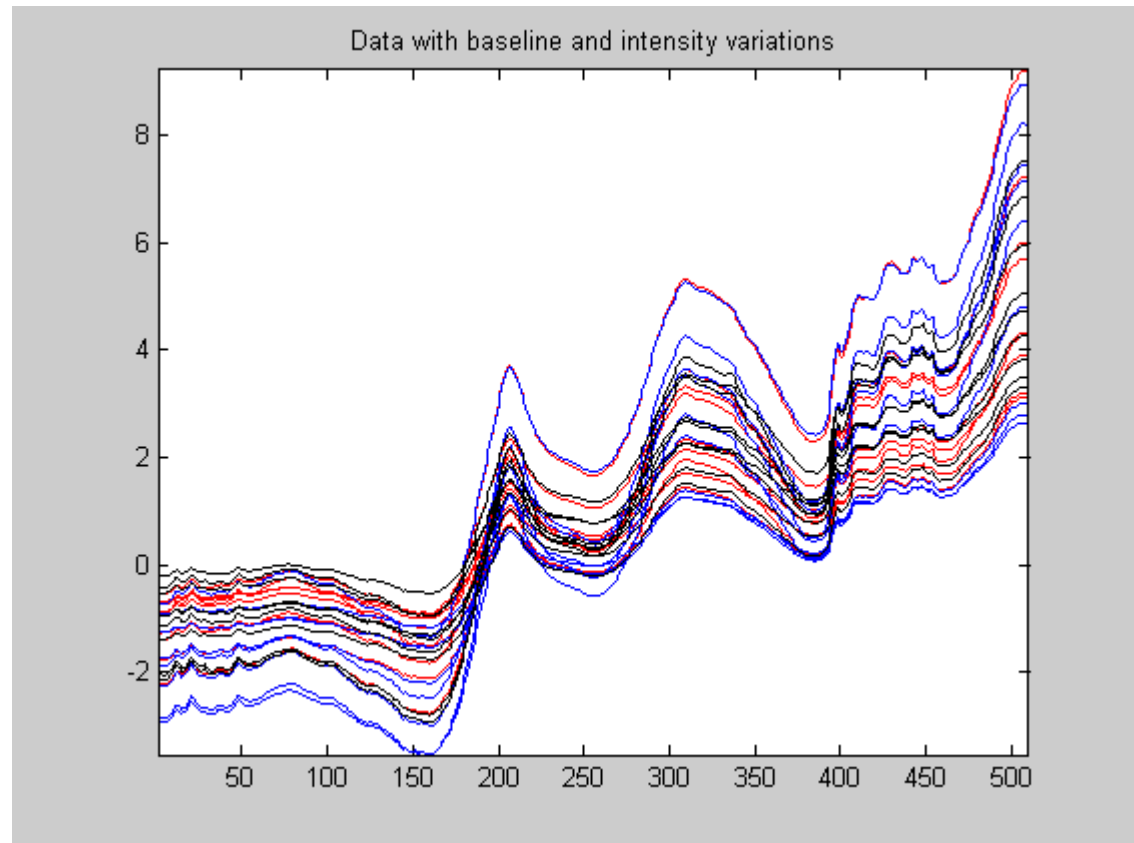
Baseline offset



Data with baseline variations

# **Raw data check**
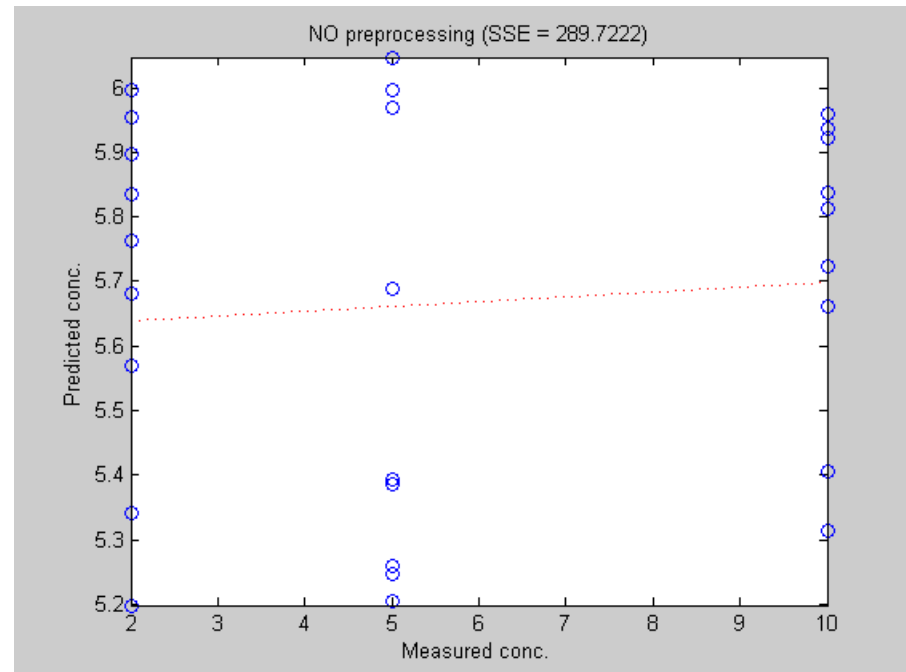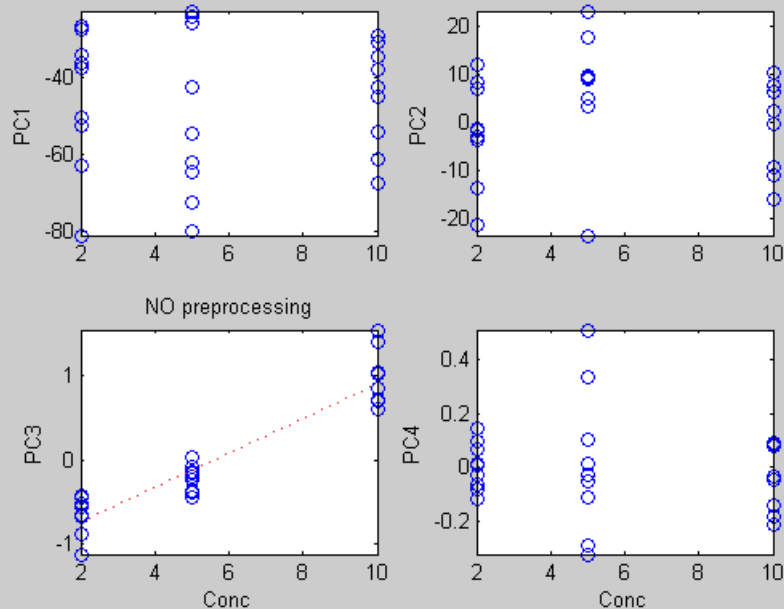
The black, red and blue curves indicate different concentration levels

Baseline offset and global intensity variations


Data with baseline and intensity variations

# PCA & PLS on Raw data

- It is not easy to separate three concentration levels
- Need to correct the spectra

# **Preliminary conclusion**

- Pre-treatment of data is crucial

- But it is not always simple …

# Common pre-treatment methods

- Baseline correction
  - Offset
  - Detrend
  - Spline
  - MSC and EMSC

- Scale correction
  - Standard Normal Variates (snv)
  - MinMax
  - Log

# Common pre-treatment methods

- Data enhancement
  - Centering
  - Standardising
  - $1^{st}$ & $2^{nd}$ order Derivatives
  - Smoothing

- Orthogonalisation
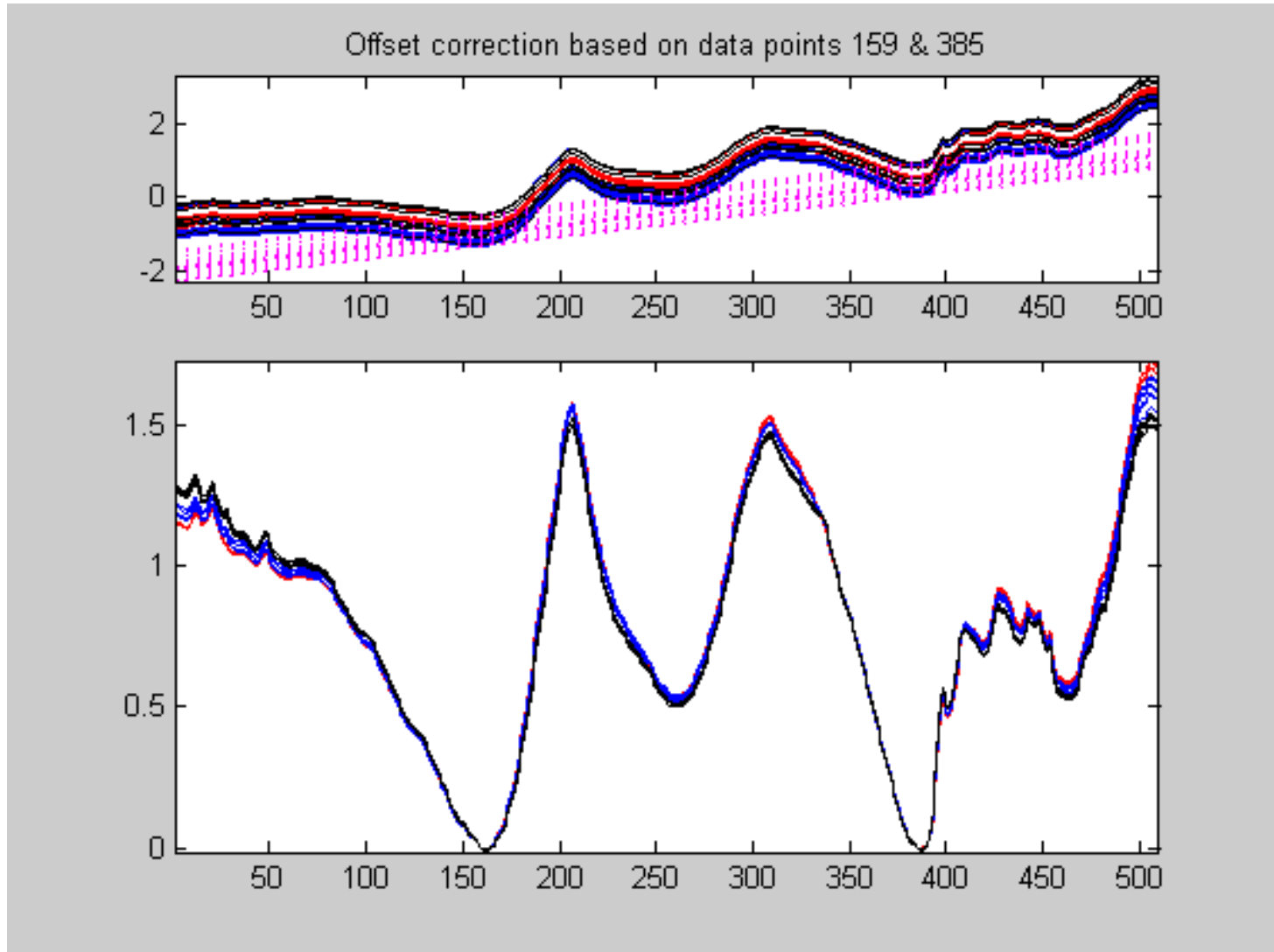  - Direct Othogonalisation
  - O-PLS
  - OSC
  - DOSC
  - …

# Offset correction

- Subtract linear baseline from each signal
  - -Intensity of lowest point
  - -Intensity of a user-chosen point
  - -Intensities calculated between 2 points
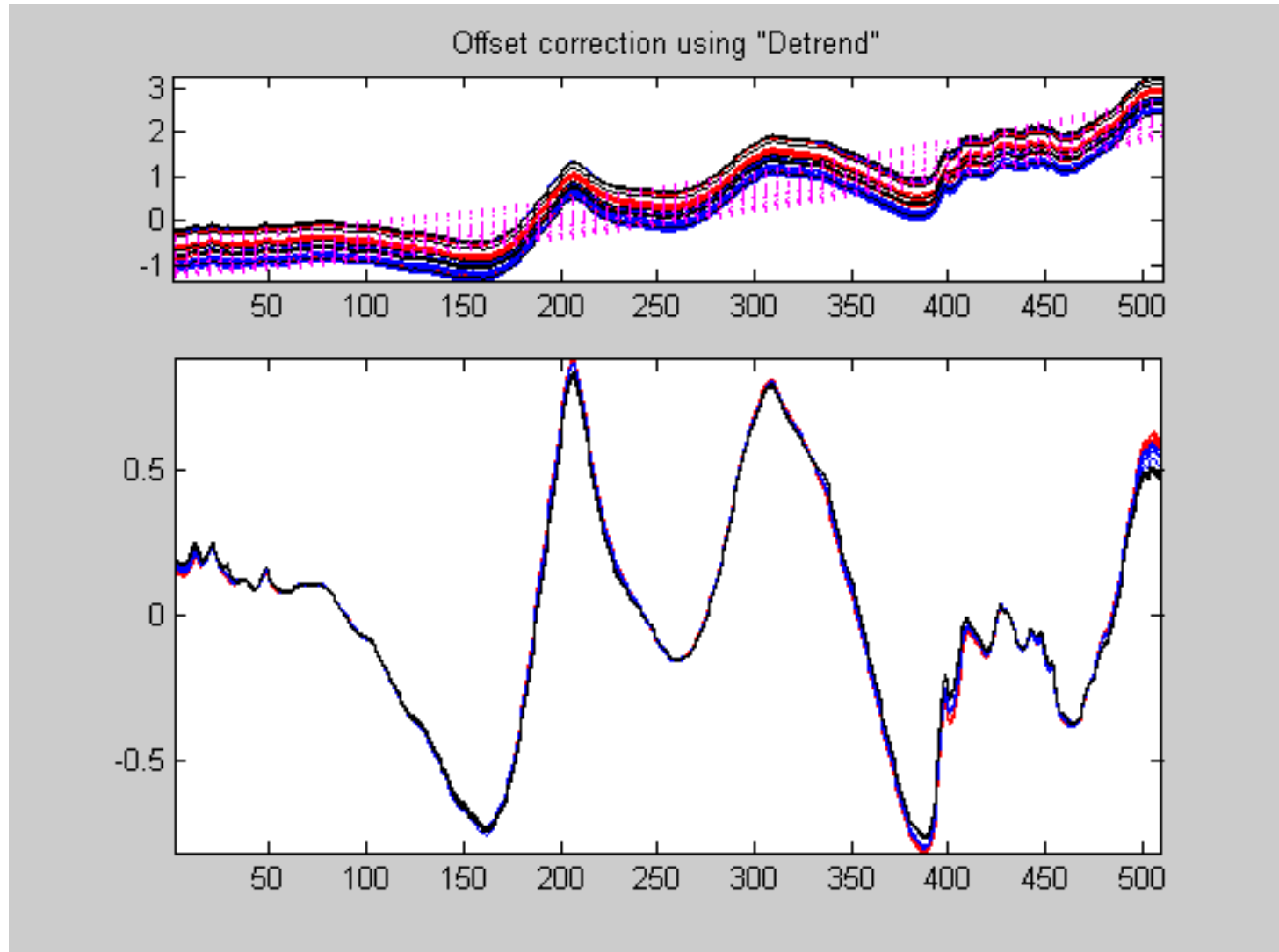
# Baseline correction

# Offset correction



Offset correction based on data points 159 & 385

# Detrend correction

- Subtract 2$^{nd}$ degree polynomial baseline from signals
  - Automatically calculated from data points
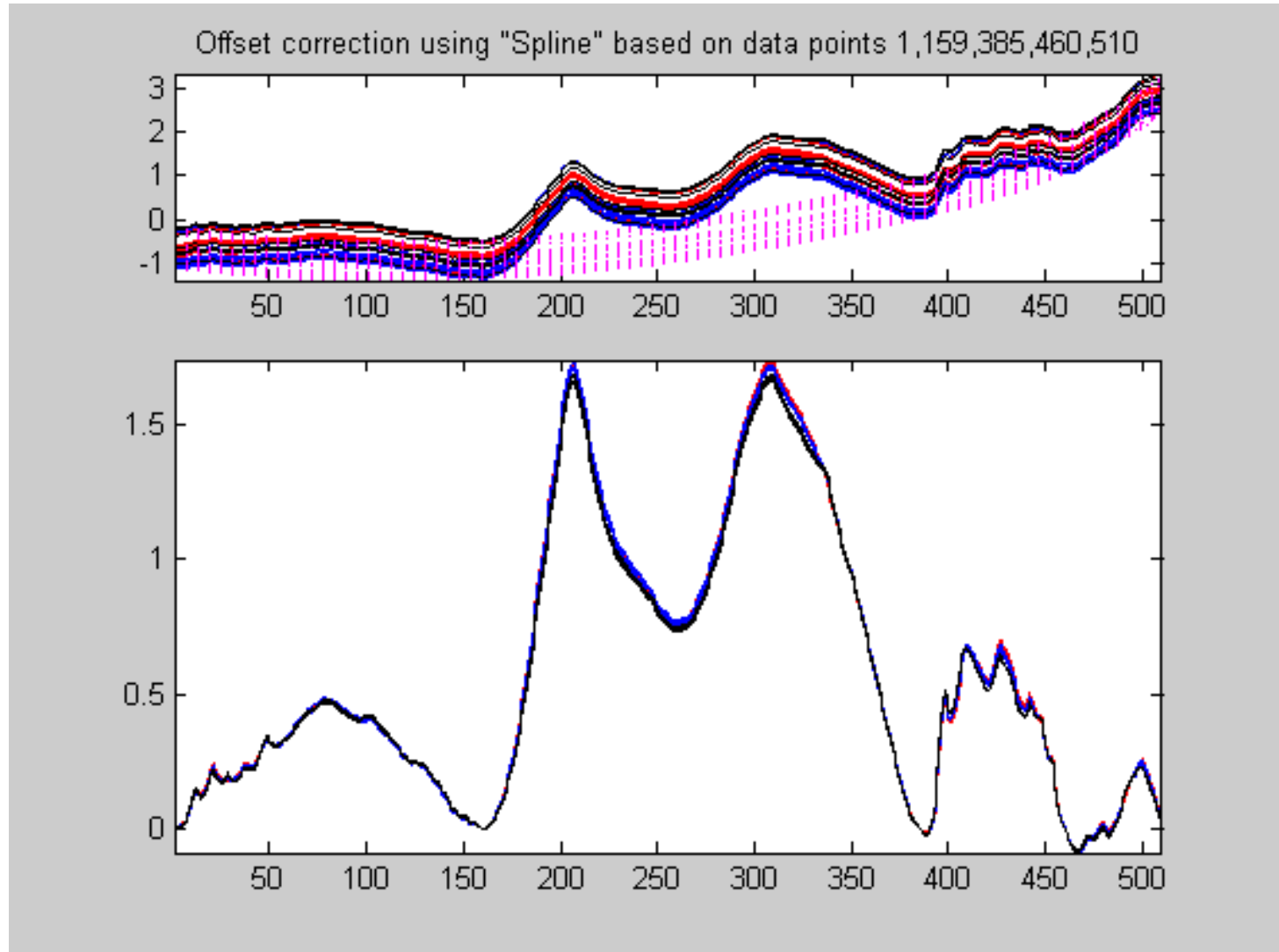
# Detrend correction



Offset correction using "Detrend"

# Spline correction

- Subtract a cubic piece-wise polynomial baseline from each signal
  - Requires input of a series of spline nodes
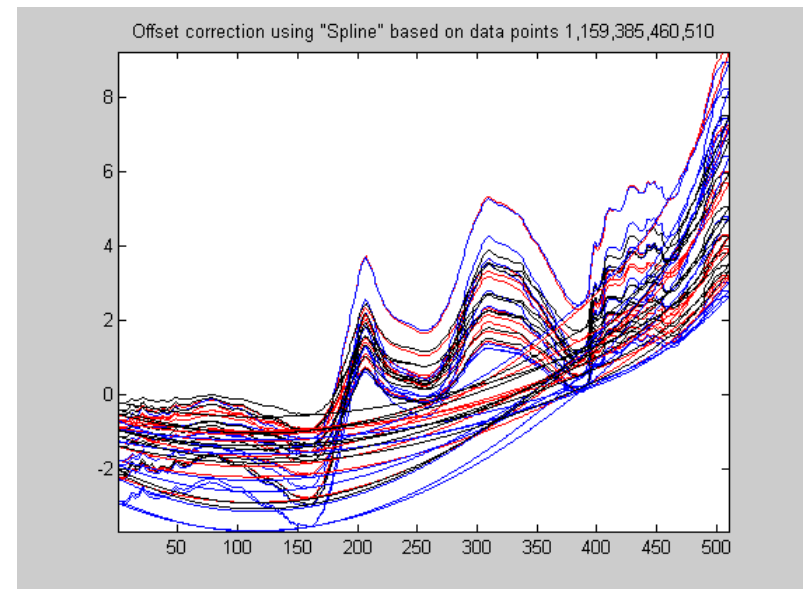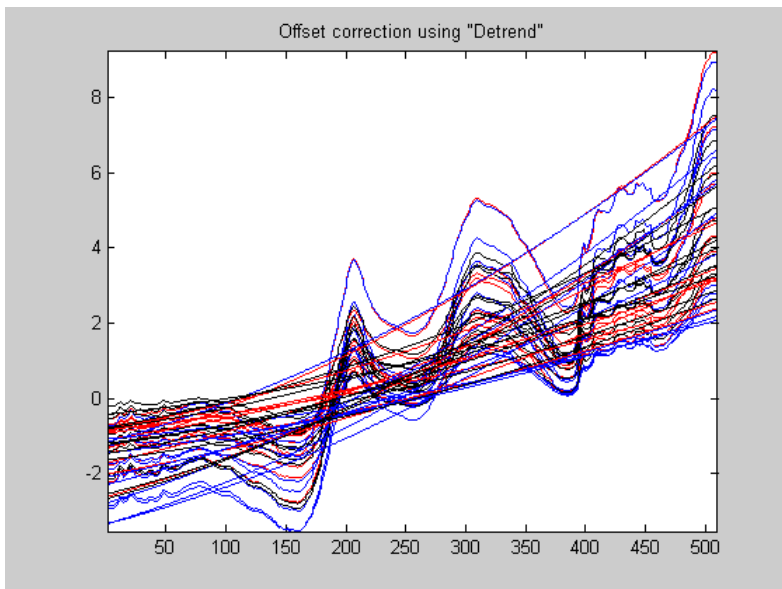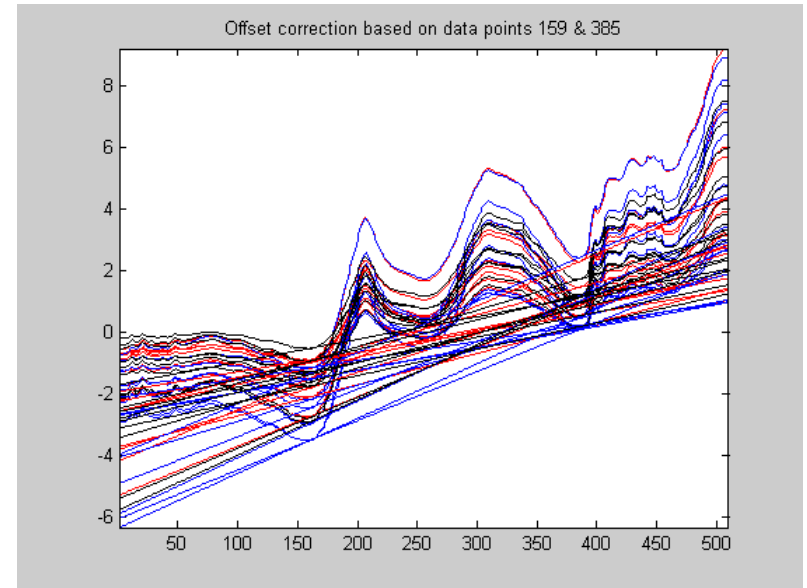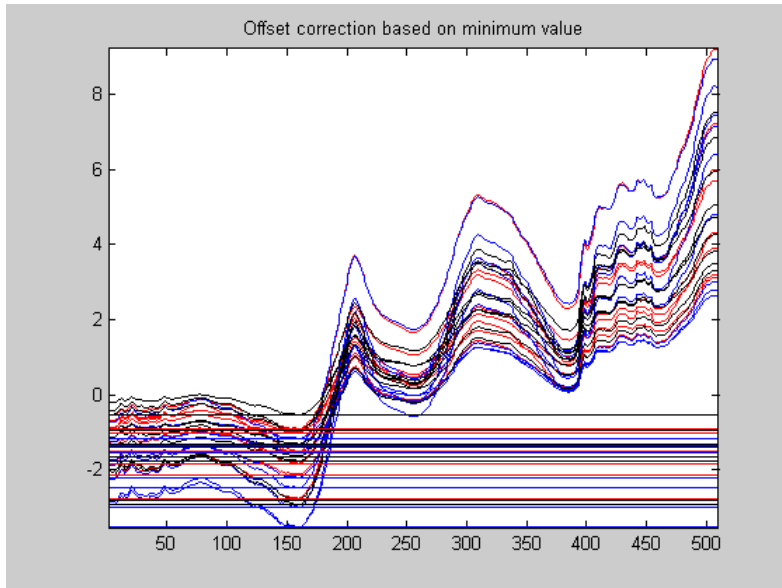  - Delicate choice with important consequences !

# Baseline correction

## Spline correction



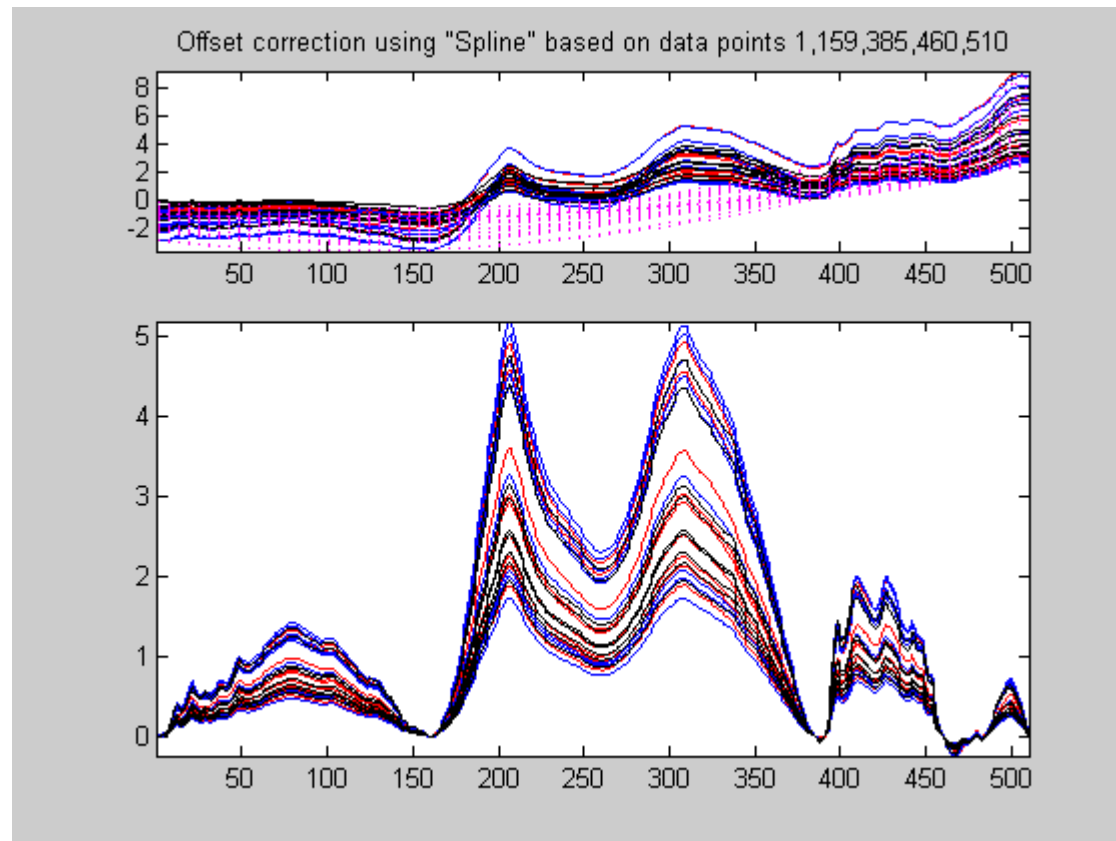Offset correction using "Spline" based on data points 1,159,385,460,510

# Baseline correction

# Baseline correction methods

# Baseline correction

- Only corrects for linear & non-linear baseline shifts
  - Does not correct for global intensity variations



Offset correction using "Spline" based on data points 1,159,385,460,510
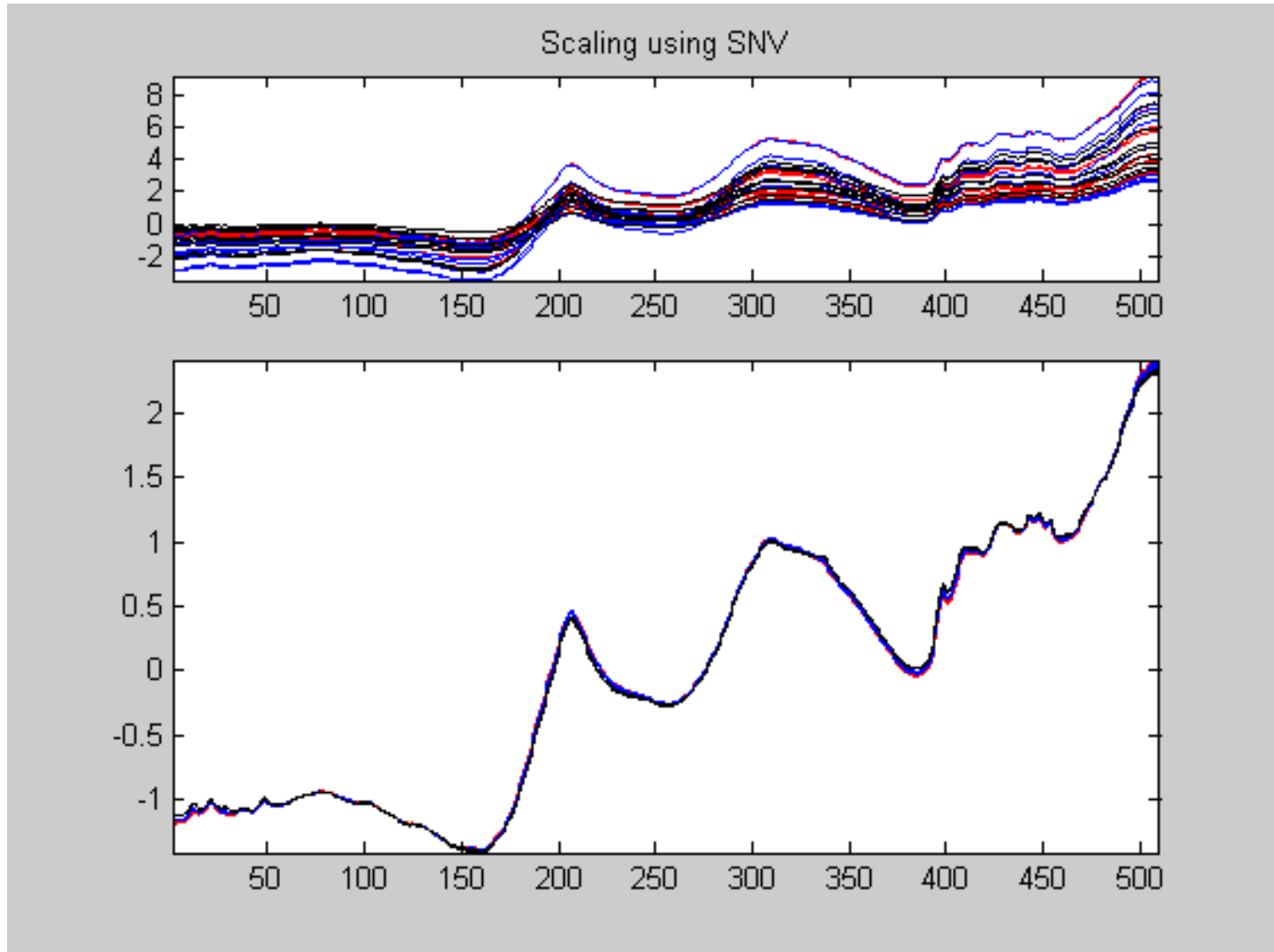
# Standard Normal Variates (SNV)

- Subtract the mean for each spectrum *i*

- Then divide by its standard deviation :

$$x_{ik}^{SNV} = (x_{ik} - m_i) / s_i$$

- SNV is a *baseline* and a *quantity* correction method
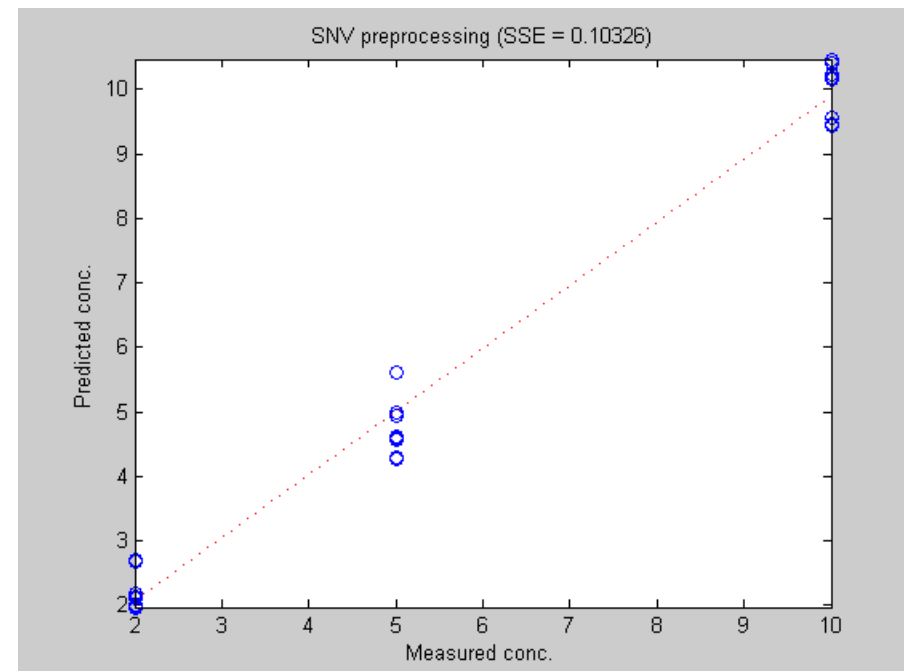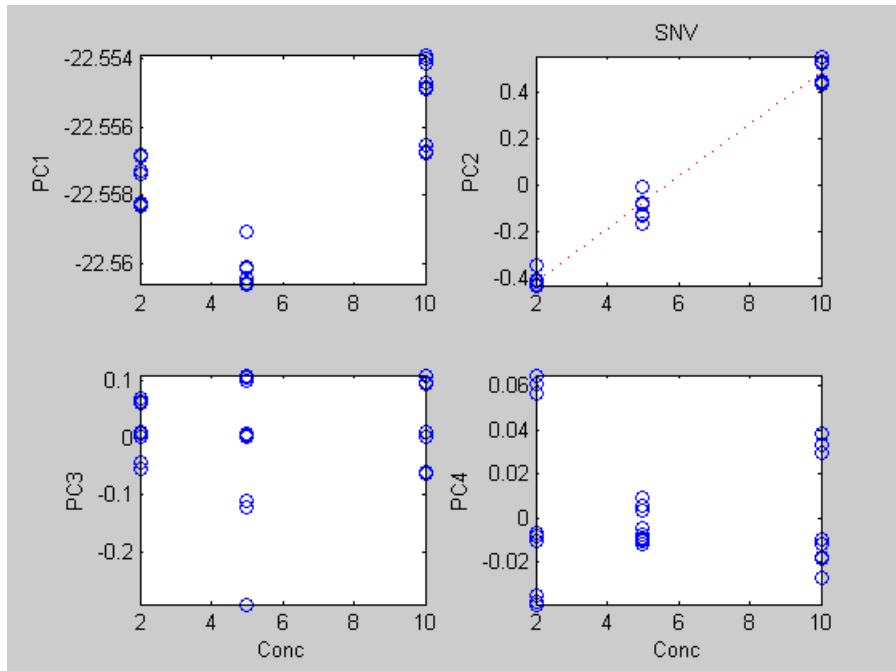
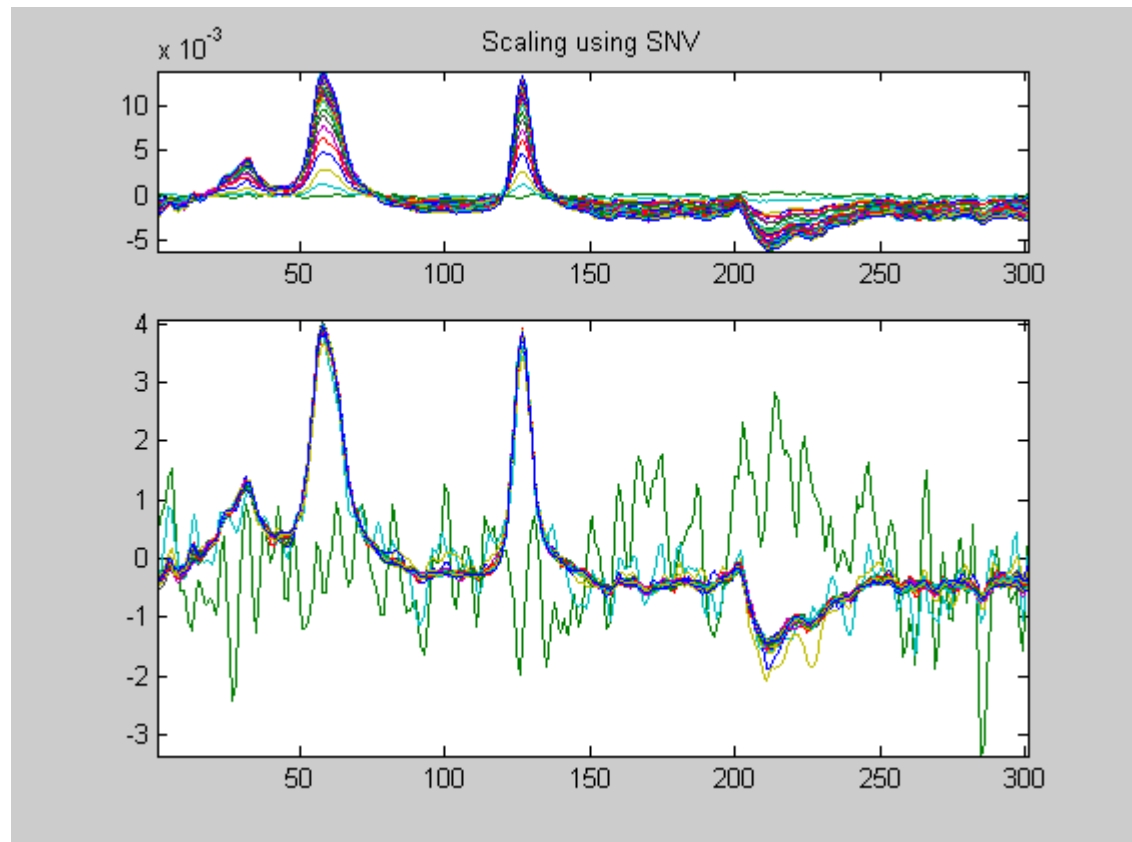**Scale correction**

# SNV correction

# PCA & PLS on SNV-corrected data

- It is easier to separate three concentration levels

# Problems with SNV

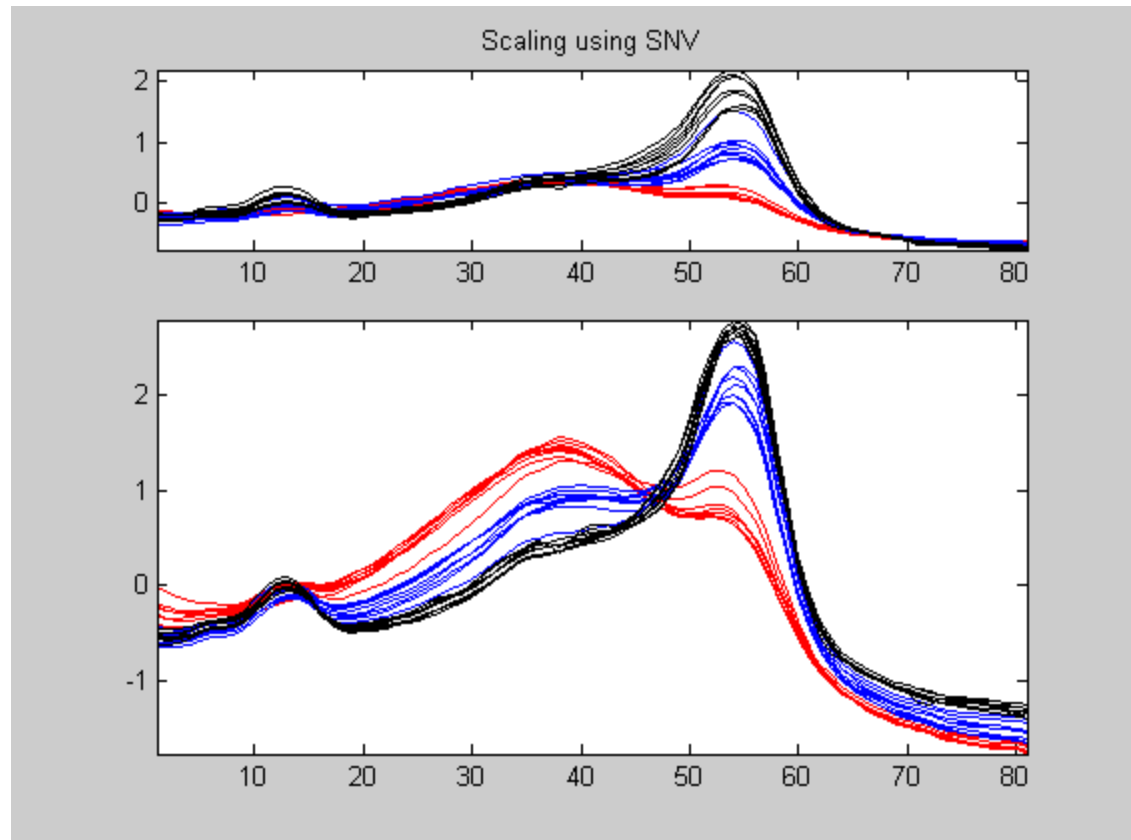- In some cases, global intensity variations *are* interesting !
- SNV enhances noisy signals

# Problems with SNV

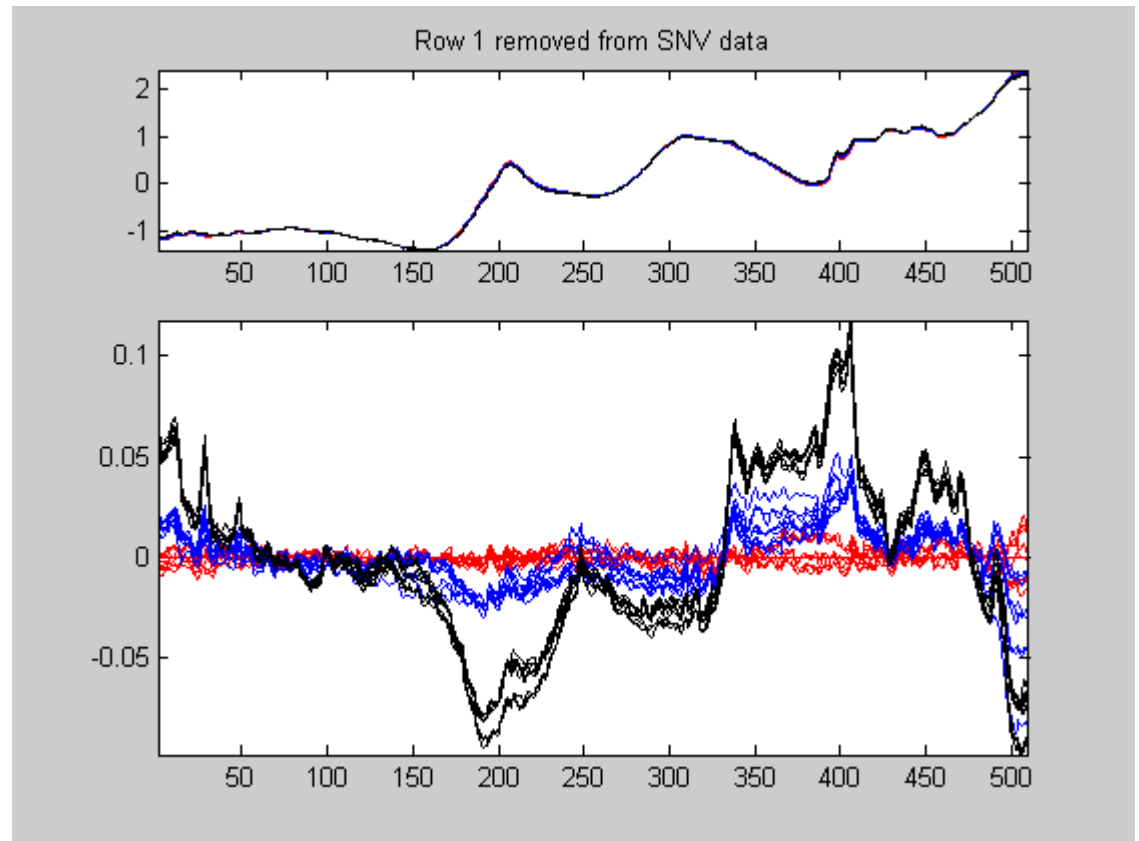- Can change relations between peaks
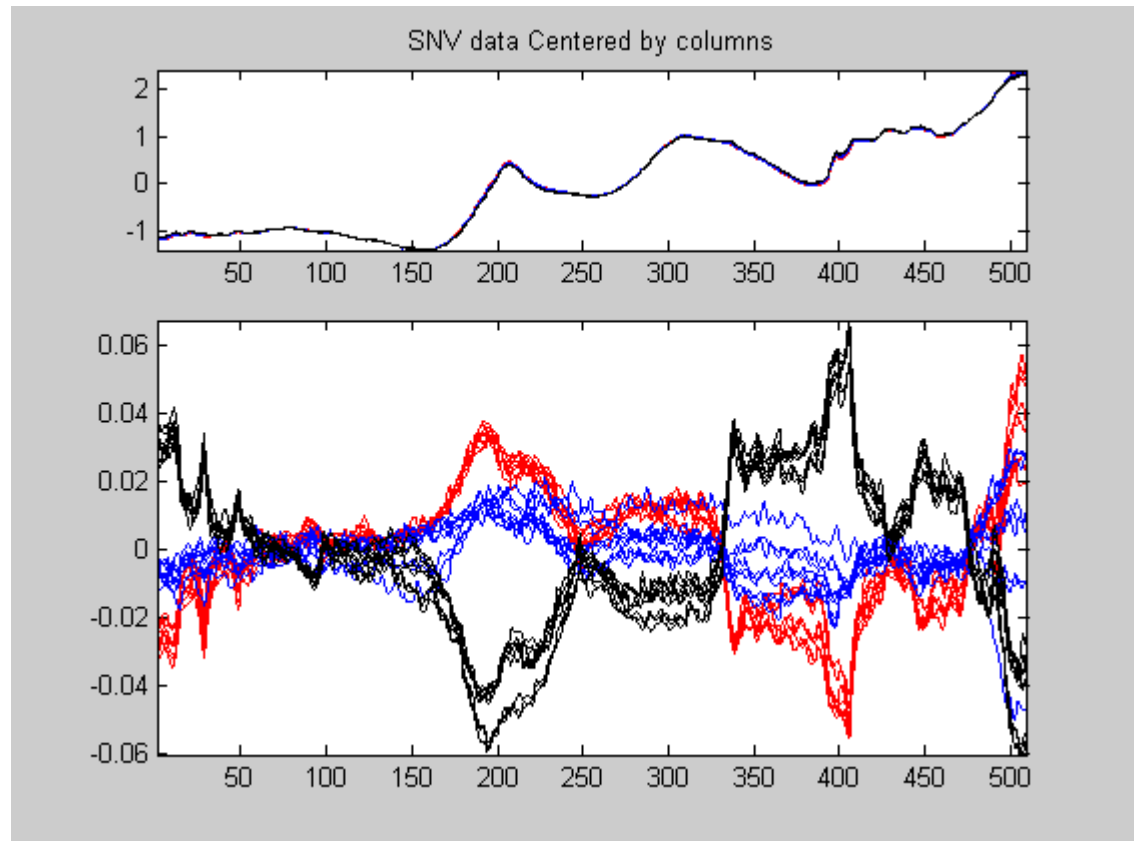- If prior to SNV one peak varies, after SNV all peaks vary

# Subtract first signal

- Highlights *evolution* of signals
- Not often used, but can be very interesting
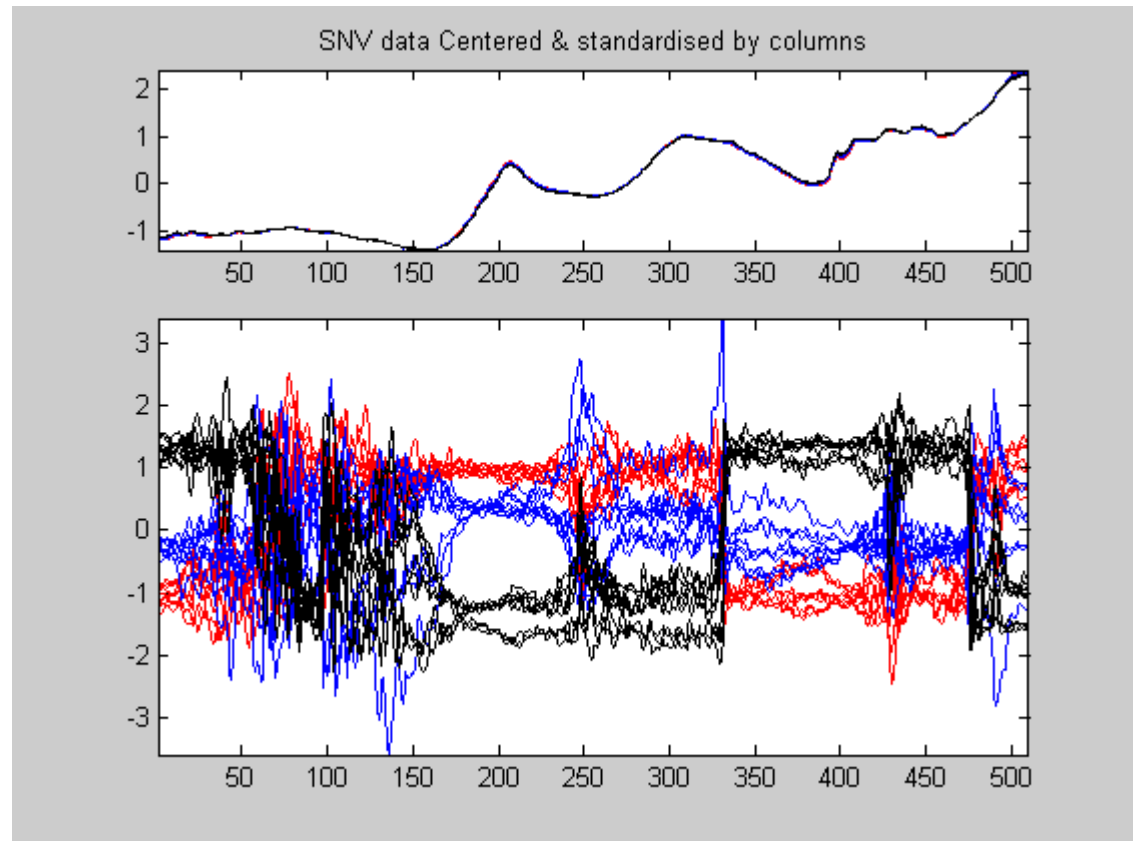- Increases apparent noise level

# Column Centering

- Often used
- Enhances differences among samples
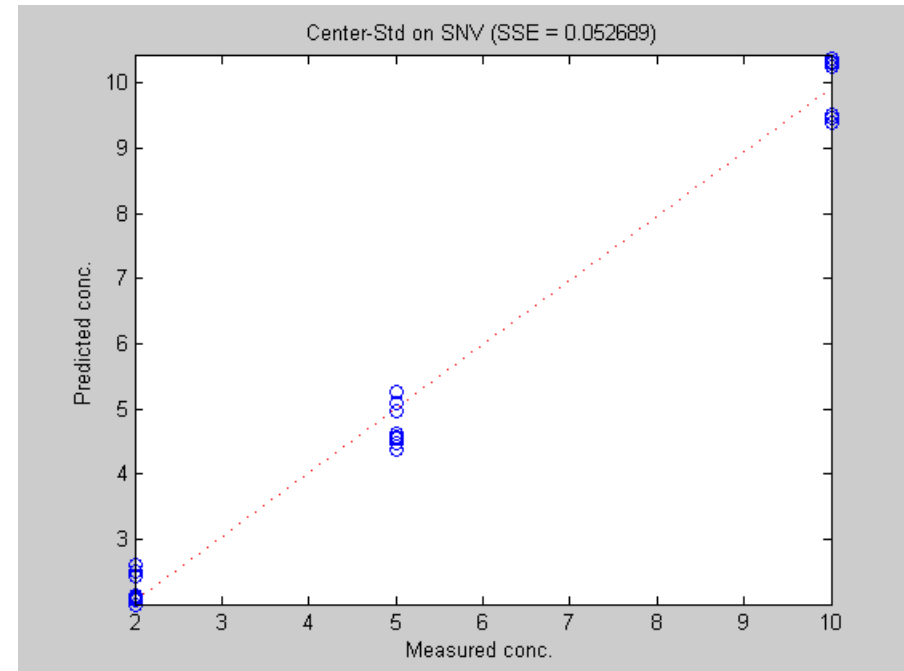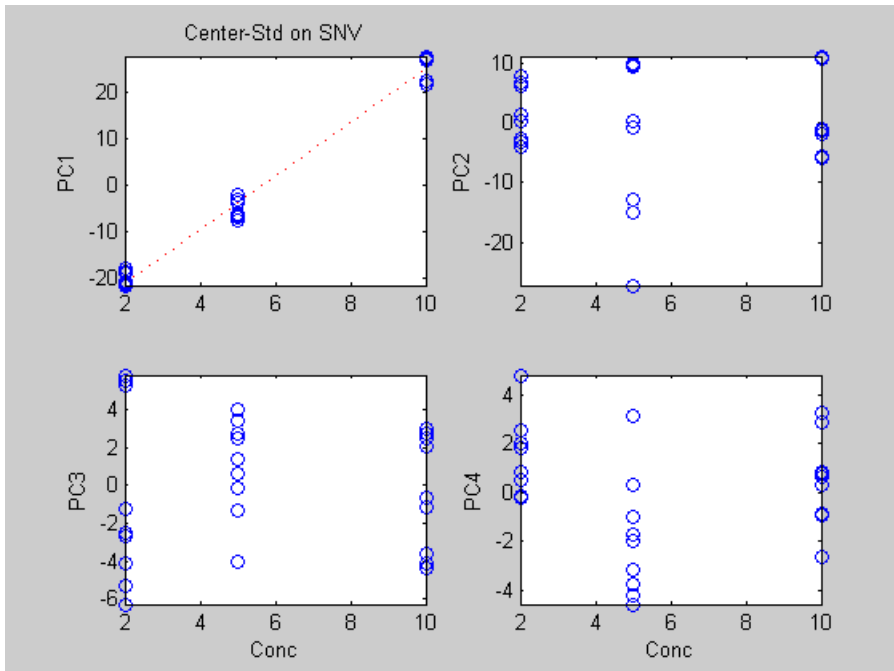- Increases apparent noise level

# Column Centering and Scaling

- Gives equal importance to all parts of signals
  - Both peaks and baseline
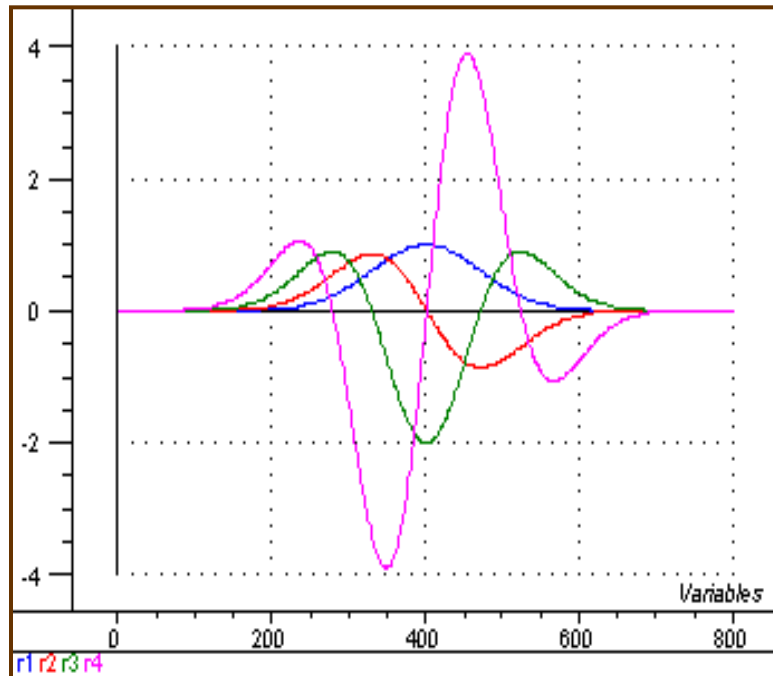  - Makes results difficult to interpret spectrally



SNV data Centered & standardised by columns

# PCA & PLS on centred and on scaled data

- Scaled data noisier
- More difficult to interpret
- But multivariate data analysis results are better
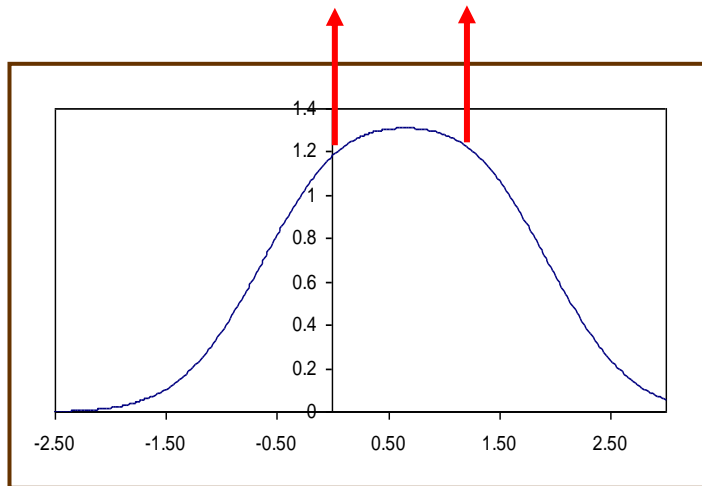
# Derivatives



- Computing derivatives of various orders is a classical technique widely used for spectroscopic applications

- Information in a spectrum may be more easily revealed when working on a 1st or 2nd order derivative
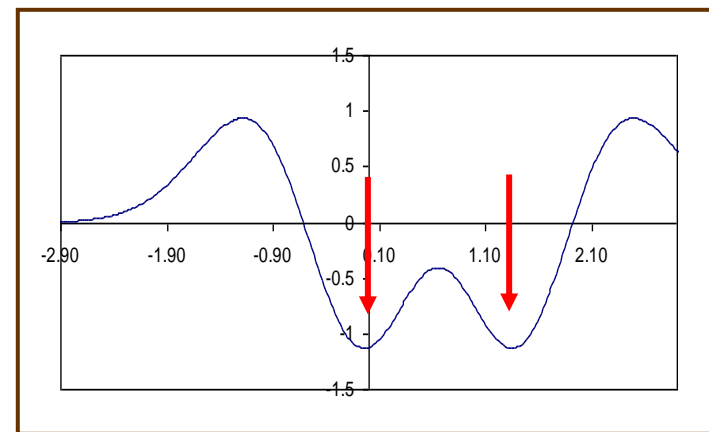
# 2nd Derivative is preferred

- 2nd derivatives is the most common preprocessing
- Removes background drift due to scattering
- Can help resolve nearby peaks
- Peak positions are at the same place as in the original spectra.
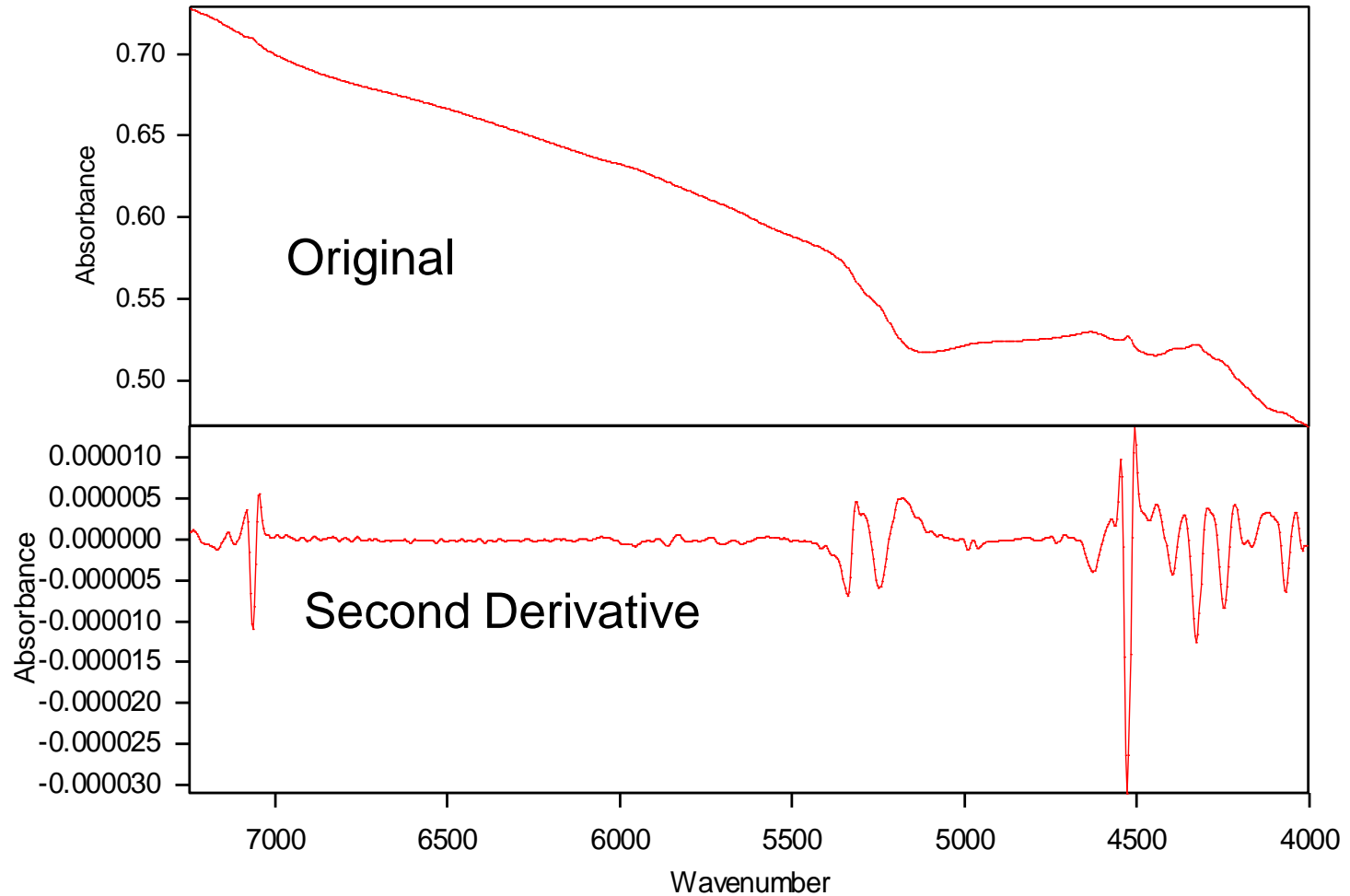- Can improve spectral resolution:



Invisible Peaks at 0.0 and 1.3           2nd derivative
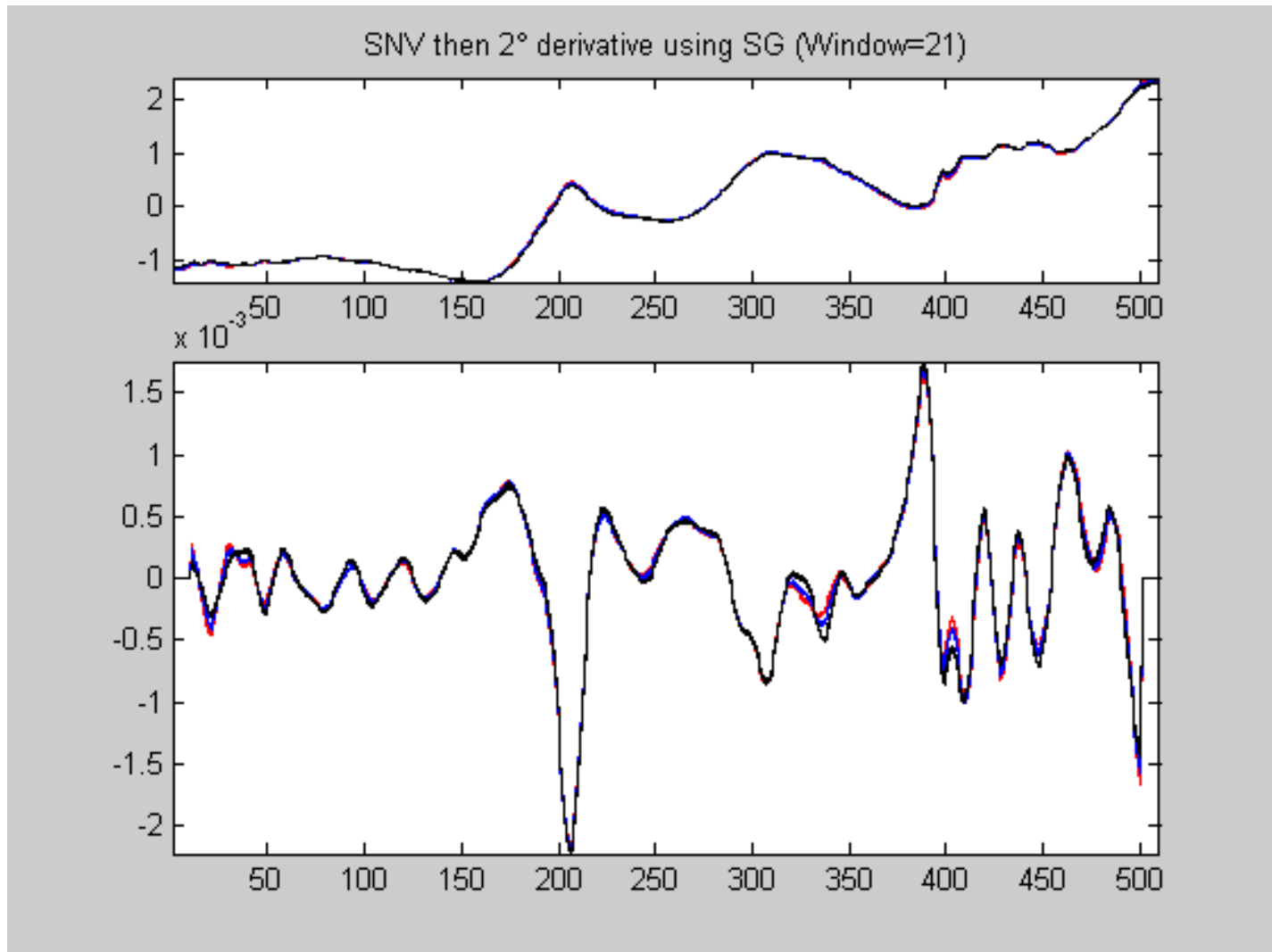
# Signal Enhancement

## Example : NIR Spectrum of Coal

# Savitsky-Golay Derivatives

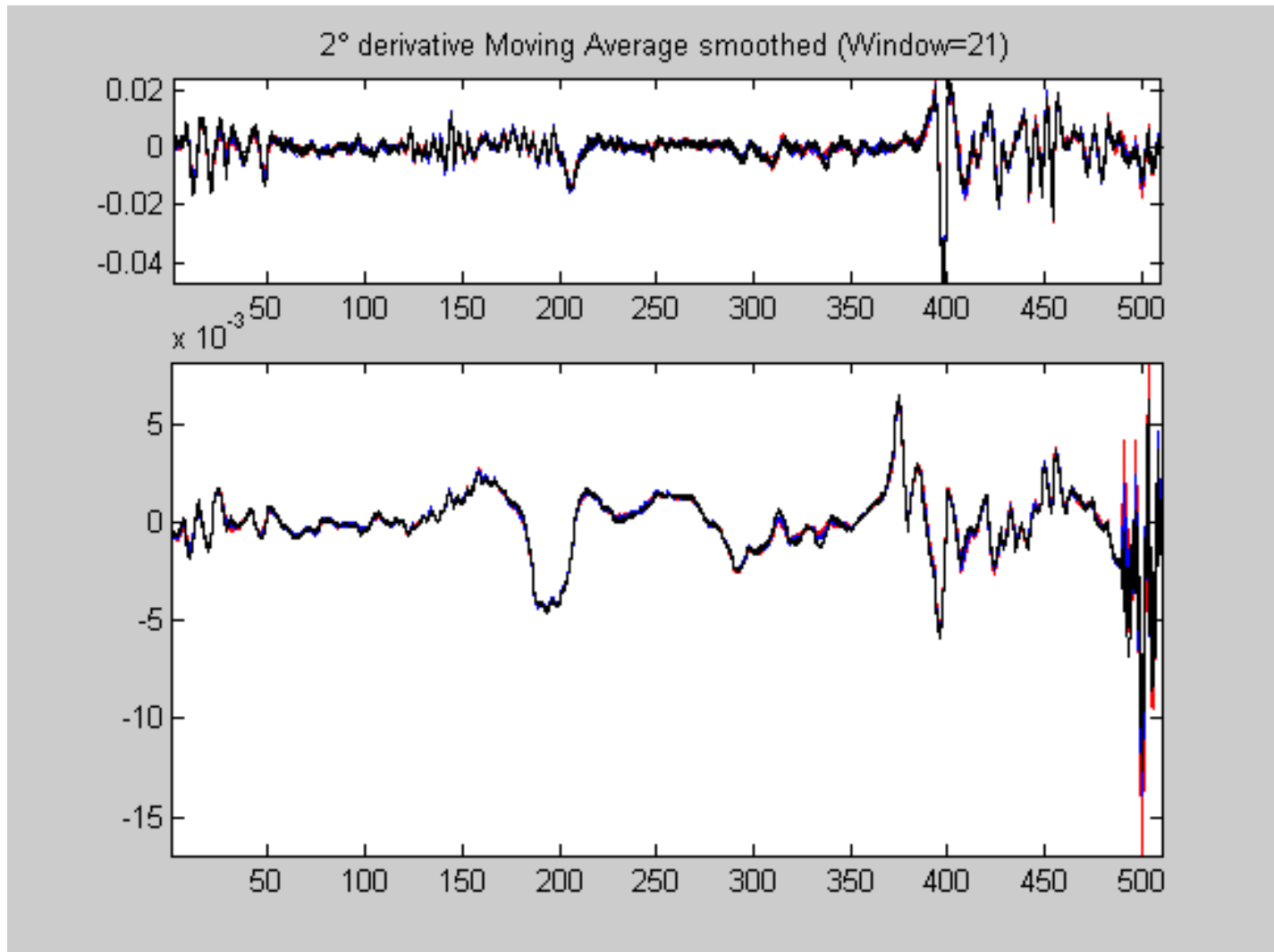Windows size : Noise decrease *vs.* loss of resolution

# Signal Smoothing

- Reduces effects of random noise

- Several algorithms :
  - Boxcar smoothing
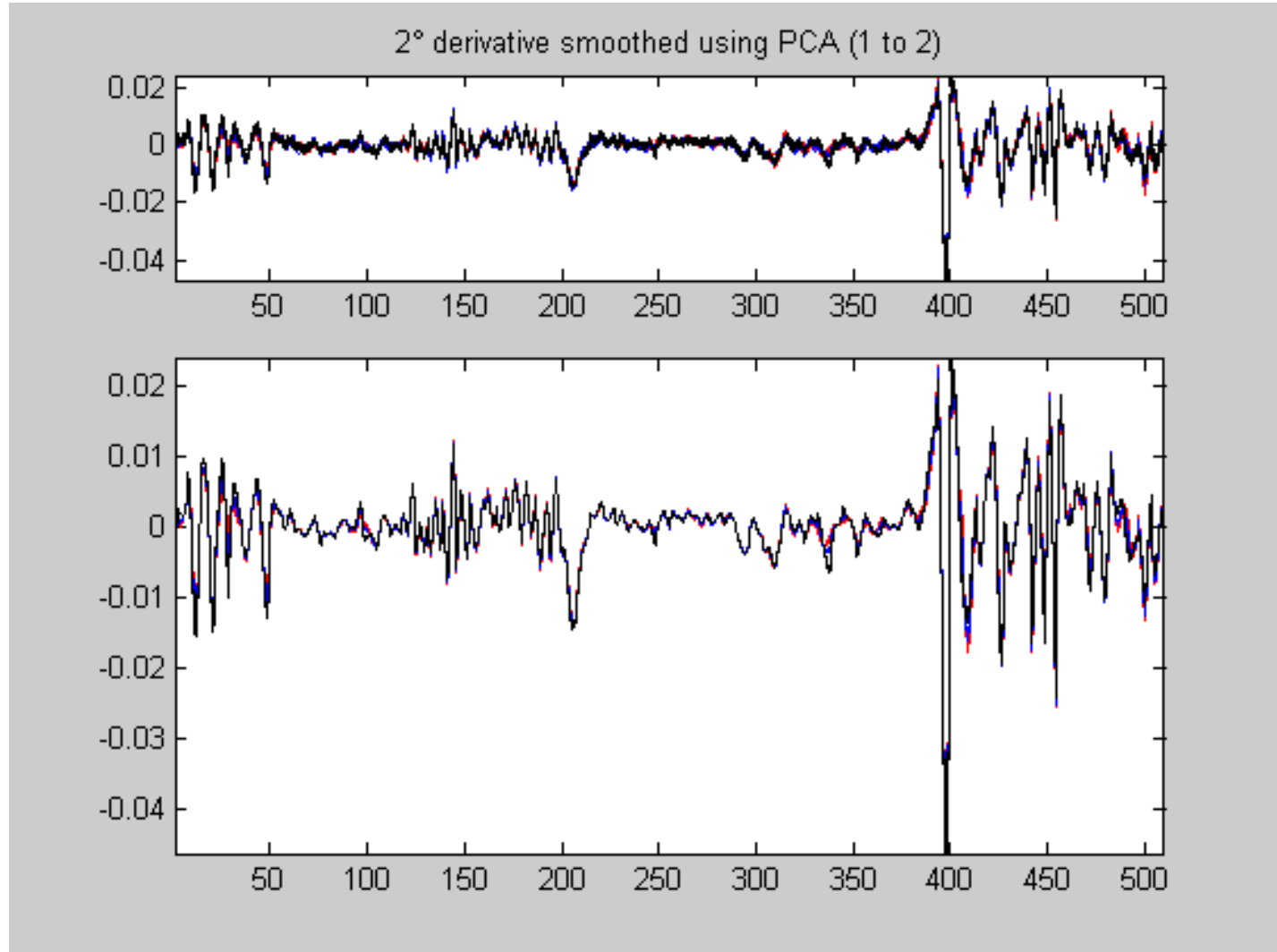  - Savitsky-Golay polynomial smoothing
  - PCA smoothing

# Savitsky-Golay smoothing

Windows size : Noise decrease *vs.* loss of resolution

# PCA smoothing



2° derivative smoothed using PCA (1 to 2)

# **Othogonalisation**

- Eliminate variability in signals not related to studied factor

- Eliminate that part of **X** which is orthogonal to ***y***
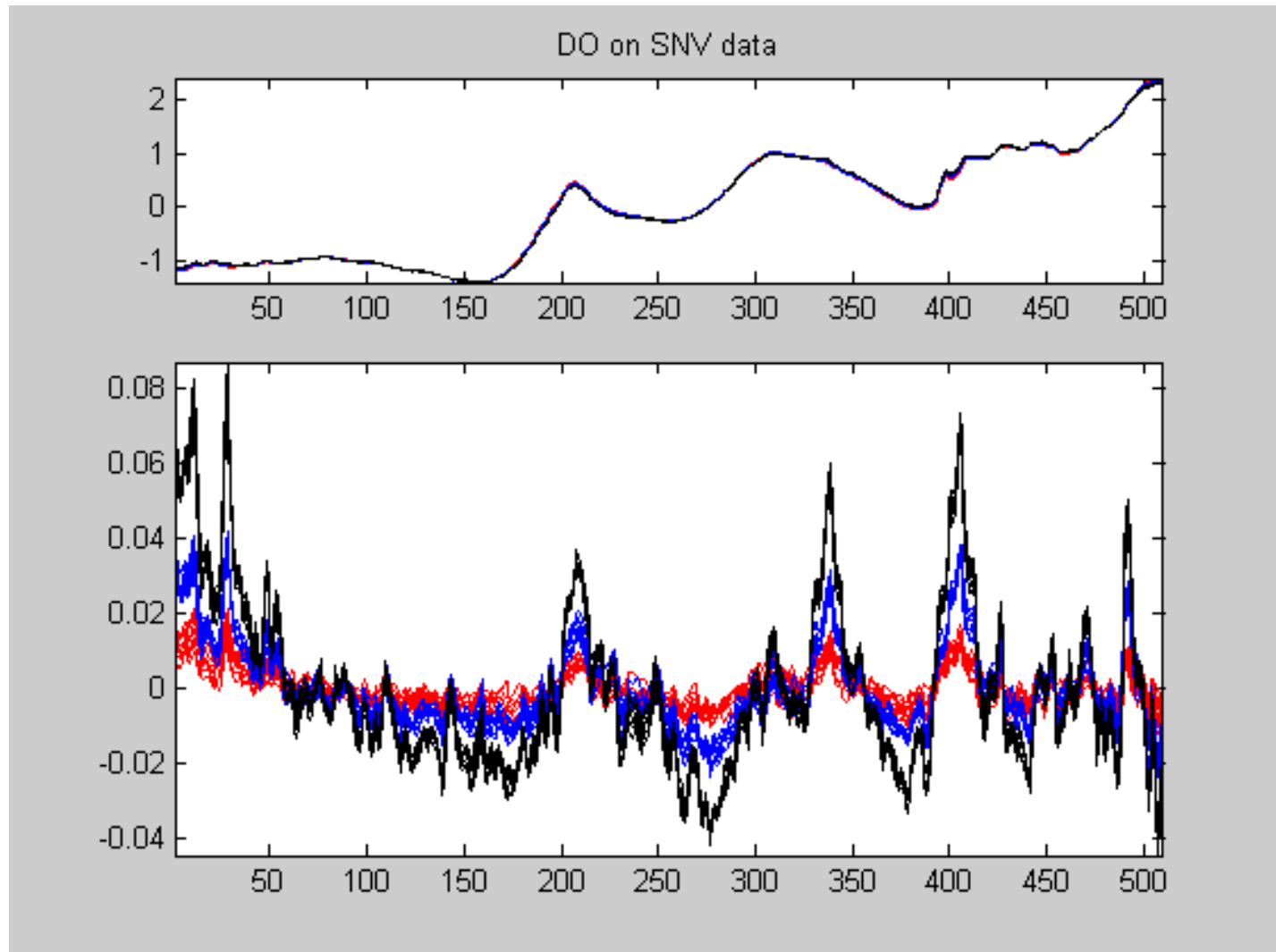  - Direct Othogonalisation
  - O-PLS
  - OSC
  - DOSC
  - …

# Direct Othogonalisation

- Calculate space orthogonal to $y = y_o$
- Project **X** onto $y_o = X_o$
- Do PCA on $X_o = T_o$ and $P_o$
- Use $T_o$ and $P_o$ to calculate interesting part of $X_o = X_o'$
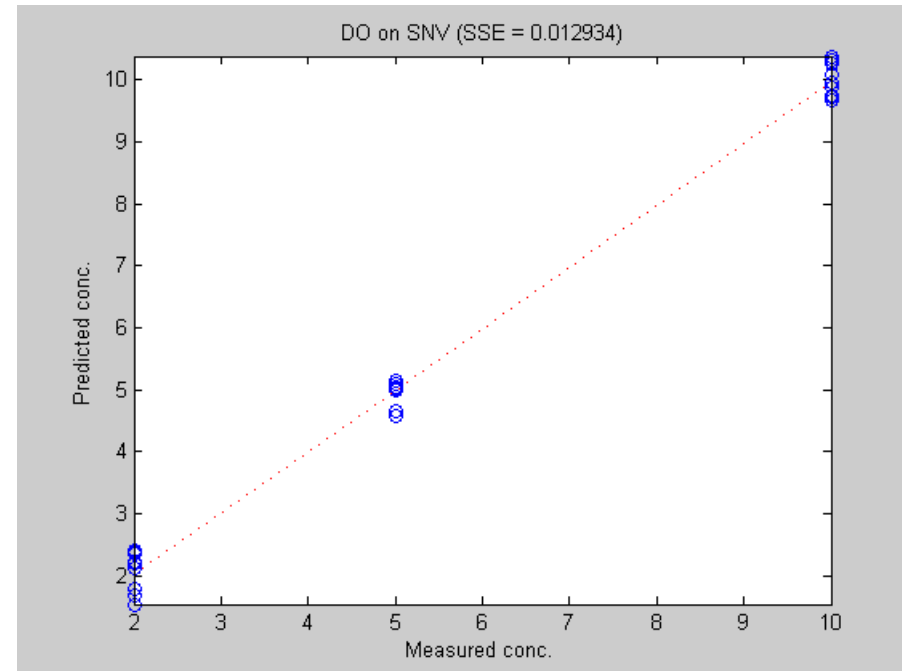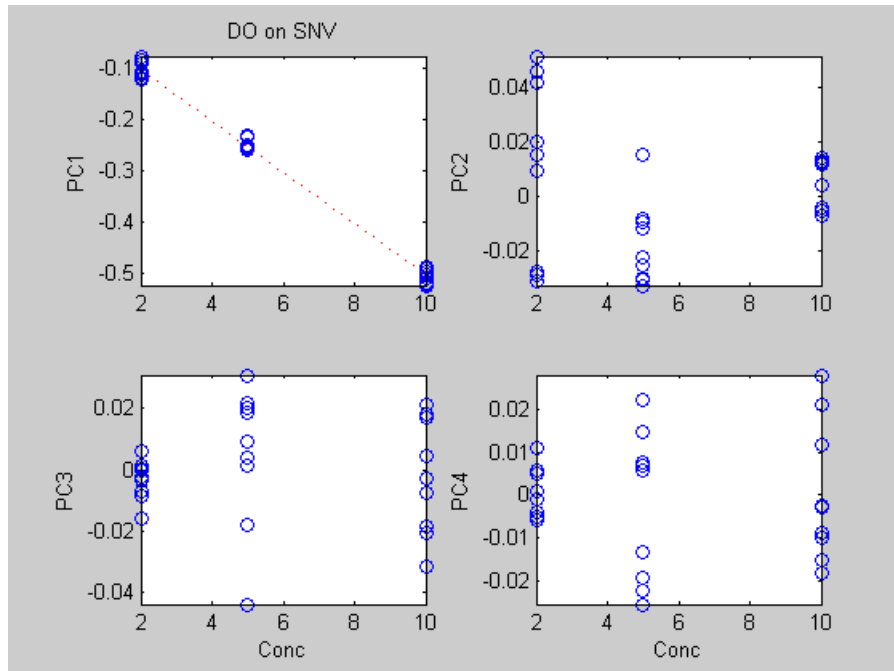
- $X_{DO} = X - X_o'$

# Direct Othogonalisation



DO on SNV data

# PCA & PLS on DO-corrected data

- It is easier to separate the three concentration levels
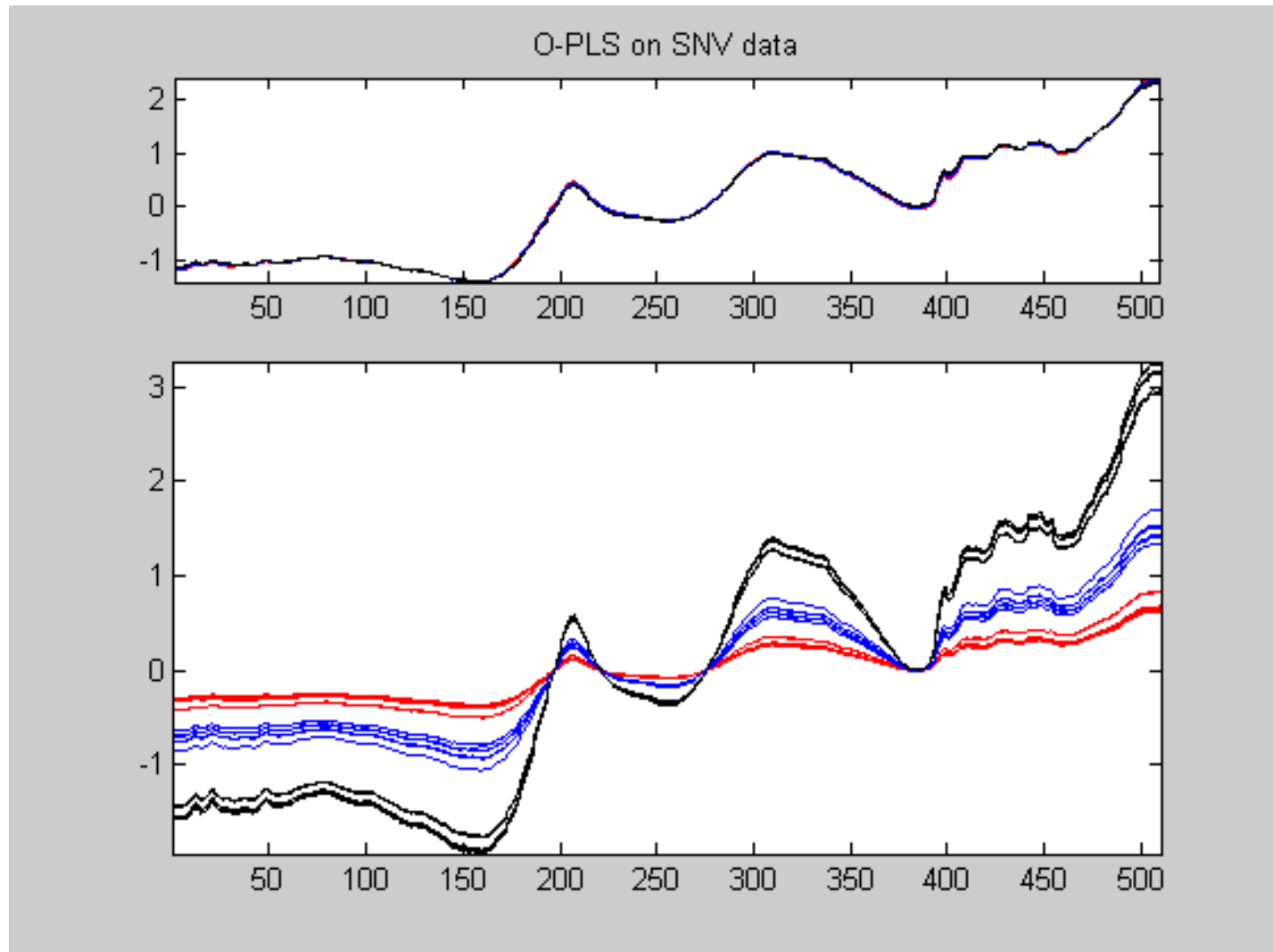- Need to determine optimal number of PCs for DO

**Othogonalisation**

# Orthogonal-PLS

- Do PLS between **X** and $y$
- Calculate $w$, $t$, $p$
- Project $w$ orthogonal to $p = \mathbf{w}_o$
- Use $\mathbf{w}_o$ to calculate orthogonal part of $t$ and $p = \mathbf{t}_o$, $p_o$
- Use $t_o$ and $p_o$ to calculate orthogonal part of **X** = $\mathbf{X}_o$

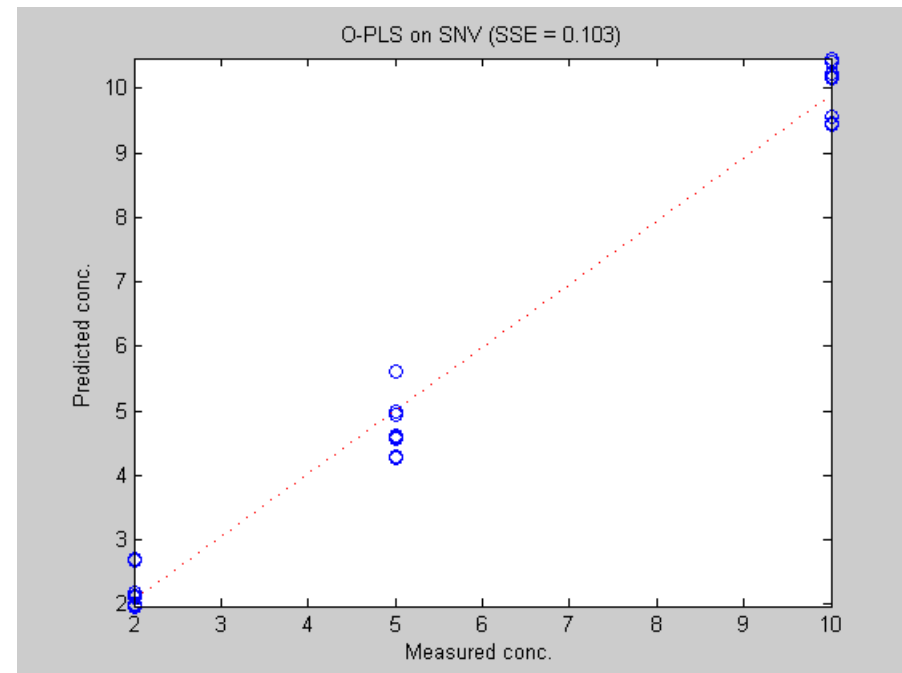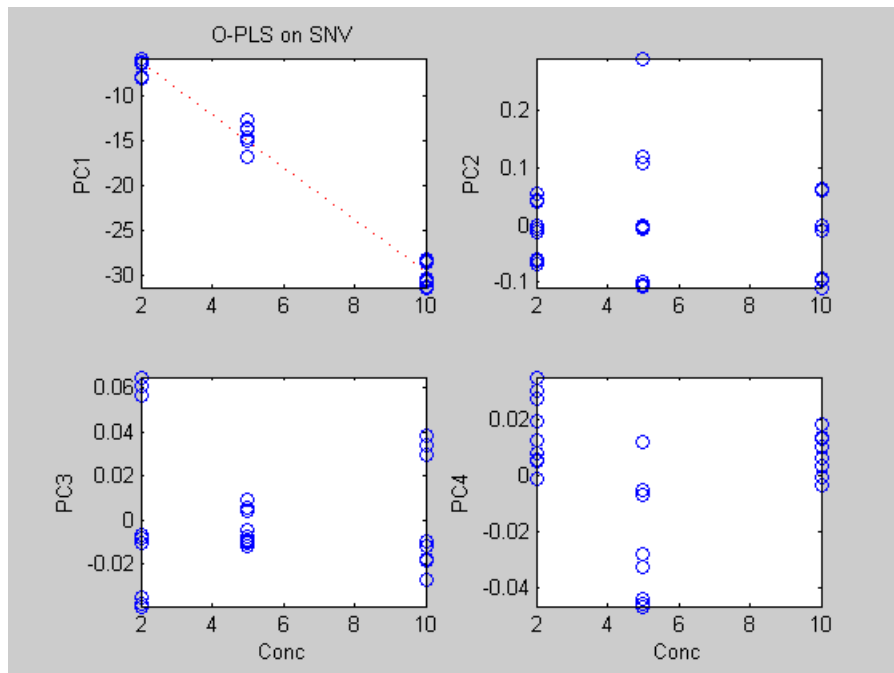- $\mathbf{X}_{O\text{-}PLS} = \mathbf{X} - \mathbf{X}_o'$

# O-PLS



O-PLS on SNV data

# PCA & PLS on O-PLS data

- It is easier to separate the three concentration levels
- Need to determine optimal number of LVs for O-PLS and for PLS !
- No real improvement in the model, just in its interpretability

# Can pretreatment of spectra improve regression models ?

| Preprocessing | SSE |
|---|---|
| None | 289.7 |
| SNV | 0.103 |
| SNV-2   Deriv. | 0.005 |
| SNV-Centering | 0.103 |
| SNV-Center/Std | 0.053 |
| SNV-DO | 0.013 |
| SNV-OPLS | 0.103 |

# Conclusions

- Pretreatments can eliminate interferences

- Pretreatments can facilitate extraction of information

- The optimal pretreatment depends on the data

# Reference

- M. Zeaiter, D. N. Rutledge
  Chapter 2 : "Preprocessing"
  Section : "Linear Regression Modeling"
  in "*Comprehensive Chemometrics*", Elsevier 2009