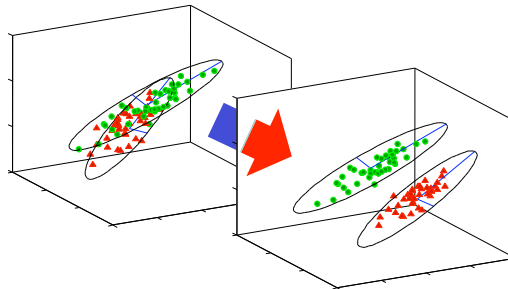


Advanced Chemometrics

Without Equations
(or hardly any)



©Copyright 2008-2017
Eigenvector Research, Inc.
No part of this material may be photocopied or
reproduced in any form without prior written
consent from Eigenvector Research, Inc.



Contact Information

Eigenvector Research, Inc.
196 Hyacinth Road
Manson, WA 98831
USA

Barry M. Wise
509-662-9213
bmw@eigenvector.com

Neal B. Gallagher
(509) 687-1039
nealg@eigenvector.com



Outline

- Introduction
- PCA Review
- PLS Regression Review
- Advanced Preprocessing
- Variable Selection
- Summary

3



Chemometrics - Use of Mathematics, Chemistry, Physics and Logic to Perform:

- **Experimental Design** - How to take measurements in such a way as to maximize the chances of obtaining the desired information at the least cost.
- **Data Analysis** - How to get as much information out of a set of measurements as possible and relate *measurements* made on a *chemical* system to the *state* of the system

4



Chemometrics Tools

- Simple exploratory analysis (e.g., PCA and PLS) are useful and help us understand the data.
 - The goal is to see trends and gain a better understanding about the measurements and system generating the data.
 - Can provide insight into how to preprocess the data
- Mathematical tools allow us to extract information from the signal (typically multivariate) that isn't always easy to see.

5



Some Resources

- Books
 - *Multivariate Calibration*, H. Martens and T. Næs, John Wiley & Sons Ltd. (1989) ISBN 0-471-90979-3
 - *Techniques and Applications of Hyperspectral Image Analysis*, Grahn, H. F.; Geladi, P., Eds. John Wiley & Sons: West Sussex, England (2007).
 - Smilde, A., Bro, R., and Geladi, P., "Multi-way Analysis with Applications in the Chemical Sciences", John Wiley & Sons, New York, NY (2004).
 - Magnus, J.R. and Neudecker, H., "Matrix Differential Calculus with Applications in Statistics and Economics, Revised Edition", John Wiley & Sons, New York, NY (1999).
- Journals
 - Journal of Chemometrics; Chemometrics and Intelligent Laboratory Systems; Analytical Chemistry; Analytica Chimica Acta; Applied Spectroscopy; Critical Reviews in Analytical Chemistry; Journal of Process Control; Computers in Chemical Engineering; Technometrics
- Special Journal Papers
 - Sanchez, E. and Kowalski, B.R., "Tensorial Calibration: II. Second Order Calibration", *J. Chemometrics*, **2**, 247–263 (1988).
 - Martens, H., Nielsen, J. P., Engelsen, S. B., "Light Scattering and Light Absorbance Separated by Extended Multiplicative Signal Correction. Application to Near-Infrared Transmission Analysis of Powder Mixtures", *Anal. Chem.*, **75**(3), 394–404 (2003).

6



Advanced Chemometrics

- Advanced concepts combine our understanding of the physics and chemistry of the system, and knowledge of how the mathematical tools work to provide better experimental designs and to ...
- maximize signal-to-noise → **signal-to-clutter**

7



Principal Components Analysis Review

- We'll come back to
"maximize signal-to-noise → **signal-to-clutter**"
- First let's review PCA and follow through an example
 - start software, load data, perform a PCA decomposition and define PCA terms

8



Outline

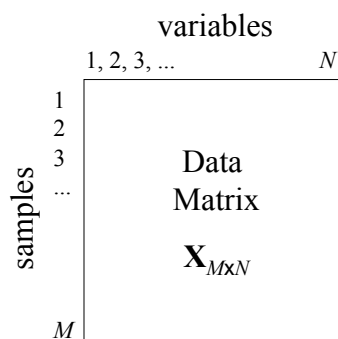
- Introduction
- PCA Review
 - Mean-centering and autoscaling
- PLS Regression Review
- Advanced Preprocessing
- Variable Selection
- Summary

9



Data Matrix X: Variables and Samples

- Examples of variables:
 - absorbance at each λ
 - ion current at each m/e
 - pressure, temperature, flow
 - chromatographic peak area
- Examples of samples:
 - samples taken to lab
 - data samples at time points
 - data from specific batches
 - etc....



10



PCA Decomposition

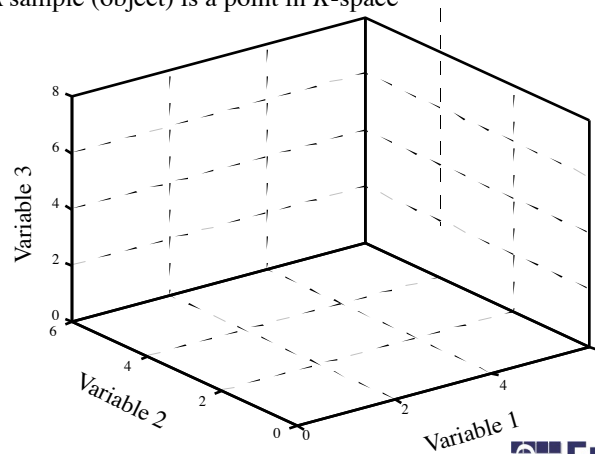
- PCA partitions a data matrix into
 - sample related information (scores) and
 - variable related information (loadings).
- Useful for multivariate exploratory data analysis.
- Scores and loadings are determined by maximizing capture of **variance**
 - information, sum-of-squares
 - show this graphically
 - Many methods in multivariate analysis are "factor based" – PCA factors are scores and loadings.

11



Principal of Projections

- K -space has K dimensions where each variable, or measurement on an object, is a coordinate axis
- A sample (object) is a point in K -space

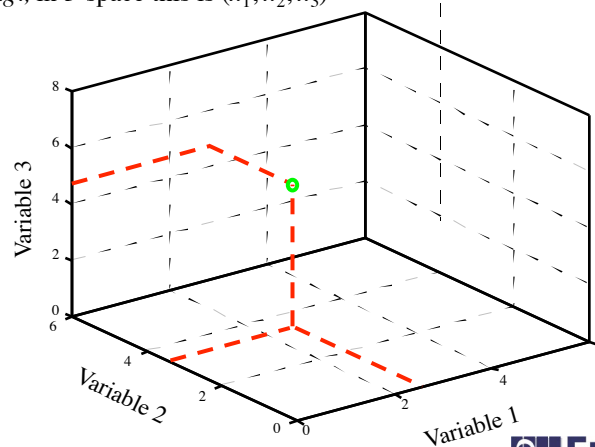


12



Projection in K-Space

- The projection of an object onto the K -space yields the coordinates of the object in that space
- e.g., in 3-space this is (x_1, x_2, x_3)

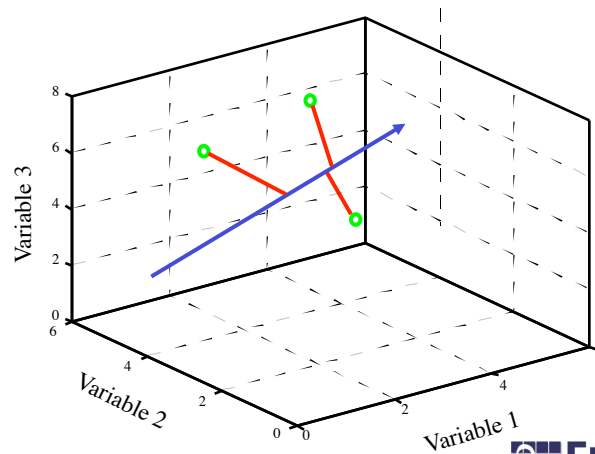


13



Projection onto a Vector

- Projection lines are perpendicular to the vector

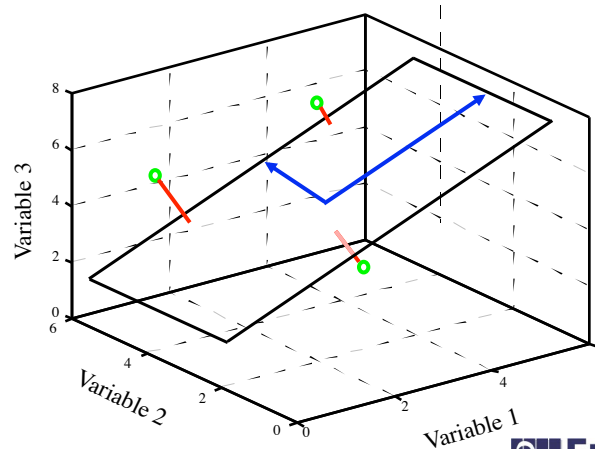


14



Projection onto a Plane

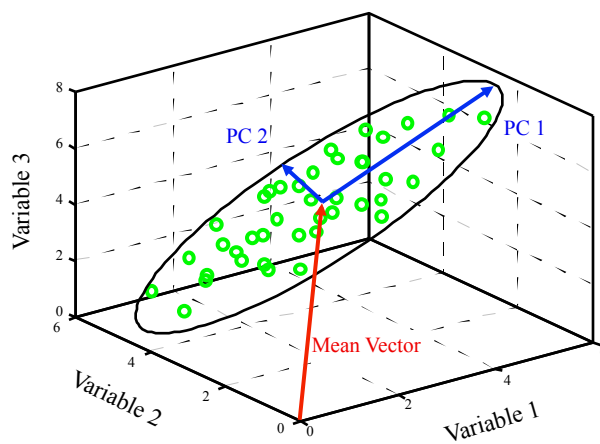
- Projection lines are perpendicular to the plane



15

 **EIGENVECTOR**
RESEARCH INCORPORATED

PCA

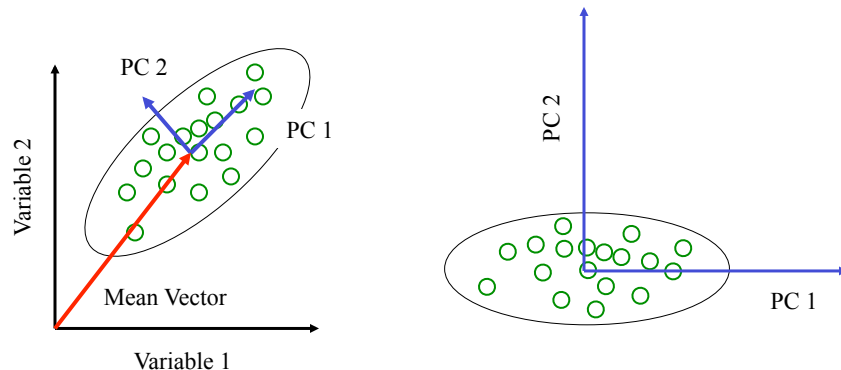


16

 **EIGENVECTOR**
RESEARCH INCORPORATED

PCA

- Geometry for 2 variables



17

 **EIGENVECTOR**
RESEARCH INCORPORATED

How Does PCA Find the PC's?

- The 1st principal component (PC) passes through the **origin** and the **maximum variance of the data**.
- The 2nd PC is **orthogonal** (perpendicular or independent) to PC1 and passes through the **second most variance**.
- The process can be continued until the **number of new PC's = number of old variables**.

18

 **EIGENVECTOR**
RESEARCH INCORPORATED

What Does PCA Give Me?

- Most of the **variance** (information) is concentrated in the first few PC's.
 - Some may be relevant to the problem of interest
- **Small random noise** is sifted into the later PC's
 - and may be thrown away - **data filtering**.
 - or used in a residuals analysis
- **Important Assumption:**
 - The signal/noise is > 1
 - *i.e.*, most of the variance is from sources other than random noise

19



What Does PCA Give Me?

- **Loadings:** Compositions of the new PC axes in terms of the old **variables**. **May be able to interpret the loadings in chemical terms, shows how variables are correlated.**
 - Loadings \Leftrightarrow Variables
- **Scores:** The position of the **samples** in the new PC coordinate system. **The closer samples are to each other in the first few PC space, the more they are alike.**
 - Scores \Leftrightarrow Samples

20



PCA

$$\begin{array}{c} \text{variables} \\ \boxed{\mathbf{X}} \end{array} = \begin{array}{c} \overline{\mathbf{p}_1} \\ \boxed{\mathbf{t}_1} \end{array} + \begin{array}{c} \overline{\mathbf{p}_2} \\ \boxed{\mathbf{t}_2} \end{array} + \dots + \begin{array}{c} \overline{\mathbf{p}_k} \\ \boxed{\mathbf{t}_k} \end{array} + \boxed{\mathbf{E}}$$

For \mathbf{X} with M samples and N variables, the PCA decomposition is:

$$\mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^T + \mathbf{t}_2 \mathbf{p}_2^T + \dots + \mathbf{t}_K \mathbf{p}_K^T + \dots + \mathbf{t}_R \mathbf{p}_R^T$$

$$\mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^T + \mathbf{t}_2 \mathbf{p}_2^T + \dots + \mathbf{t}_K \mathbf{p}_K^T + \mathbf{E} = \mathbf{T}_K \mathbf{P}_K^T + \mathbf{E}$$

$R \leq \min(M, N)$ is the mathematical **rank** of the data.

$K \ll R$ is the pseudo- or **chemical-rank** of the data.

The \mathbf{p}_i are eigenvectors of the covariance matrix of \mathbf{X} and λ_i are eigenvalues. Amount of variance captured by each $\mathbf{t}_i \mathbf{p}_i^T$ is proportional to λ_i .

21



Properties of PCA

- $\mathbf{t}_i, \mathbf{p}_i$ ordered by amount of *variance captured* $\propto \lambda_i$
- the chemical-rank K is the number of PCs that captures other than random noise
- \mathbf{t}_i or *scores* form an orthogonal set \mathbf{T}_k which describe relationship between *samples*
- \mathbf{p}_i or *loadings* form an orthonormal set \mathbf{P}_k which describe relationship between *variables*
- scores and loadings plots are interpreted in pairs
 - e.g. plot \mathbf{t}_i vs sample number and \mathbf{p}_i vs variable number
- it is useful to plot \mathbf{t}_{i+1} vs. \mathbf{t}_i and \mathbf{p}_{i+1} vs. \mathbf{p}_i

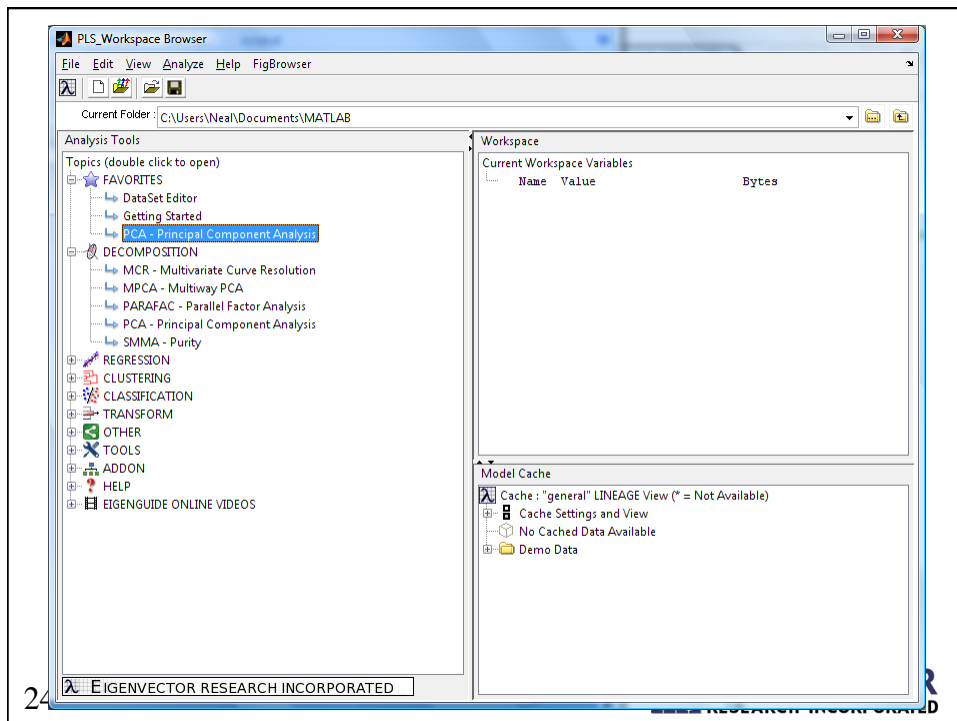
22



Example: Olive Oils

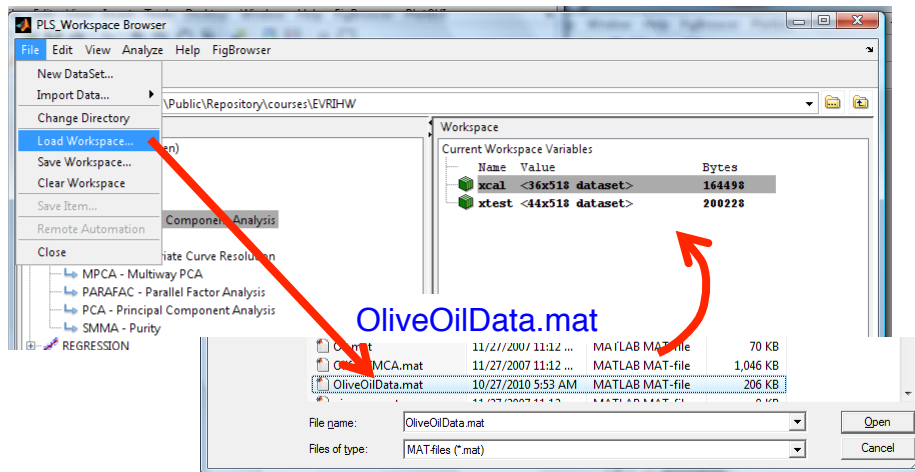
- Use FT-IR spectra and PCA for pattern recognition to distinguish authentic olive oil from counterfeit or adulterated olive oil.
- Shall see some special properties associated with **Spectral Data**.
 - Dahlberg, D.B., Lee, S.M, Wegner, S.J. and Vargo, J.A., "Classification of Vegetable Oils by FT-IR," *Appl. Spec.*, **51**(8), 1118-1124 (1997).
 - FT-IR spectra ($3600 - 600 \text{ cm}^{-1}$) using a fixed pathlength NaCl cell

23



24

Load OliveOilData.mat



25

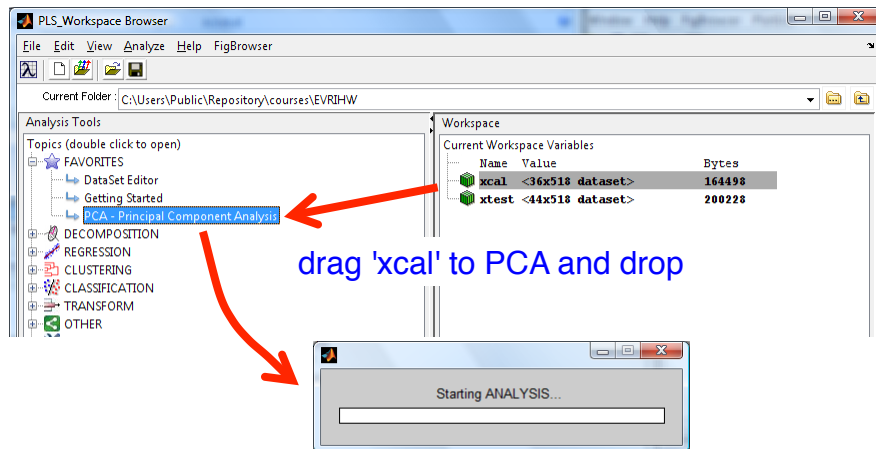


Olive Oil Samples

Learning set: **xcal** Start with this data set

Corn Oil	9 samples	(#1-9)
Olive Oil	15 samples	(#10-24)
Safflower Oil	8 samples	(#25-32)
Corn Margarine	4 samples	(#33-36)

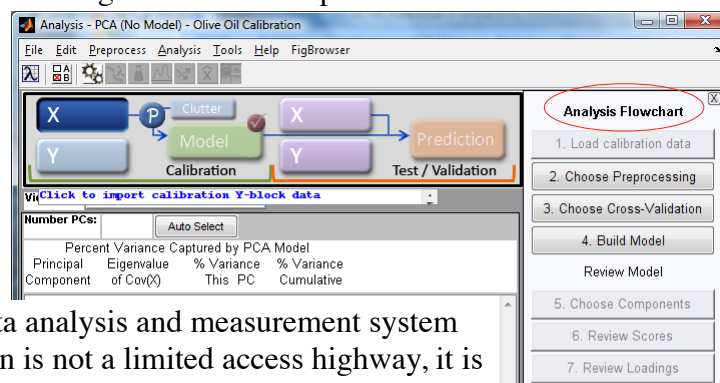
Start Analysis for PCA



27



Plot the data. (for this example let's include all the variables)
Use knowledge and logic during the analysis – this is *not* a black box.
Before modeling ask, "what will PCA give me for this data and this preprocessing? What is the expected **rank** of the data?"

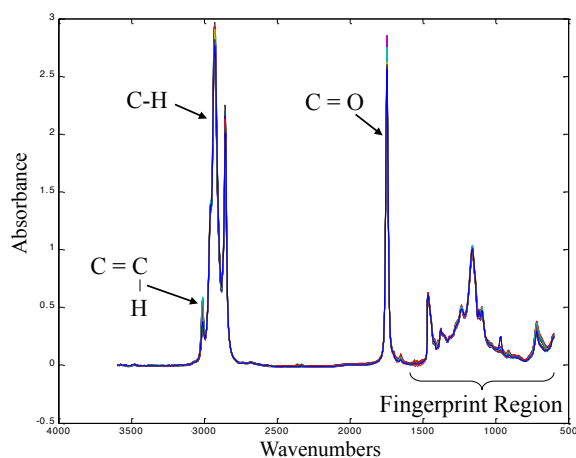


Data analysis and measurement system design is not a limited access highway, it is more like a worn path in the dirt.
Often what is learned at one step leads us back to the beginning.

28



Spectra of 36 Sample Learning Set



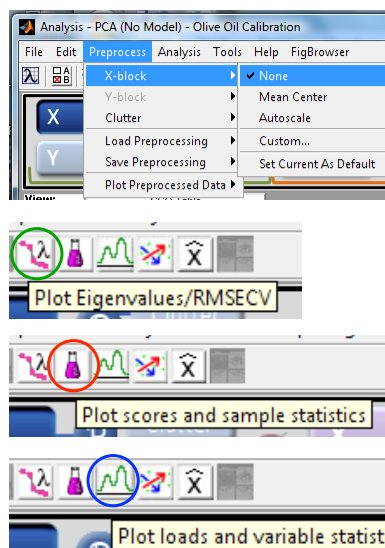
Notice how spectra look alike.

29

EIGENVECTOR
RESEARCH INCORPORATED

Try PCA

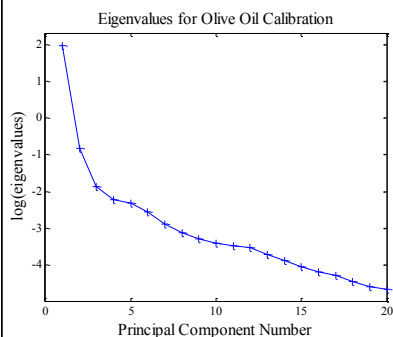
- Use no preprocessing
- plot the **eigenvalues**
 - choose number of PCs
- plot **scores** and **loadings**
 - interpret the results



30

EIGENVECTOR
RESEARCH INCORPORATED

Choosing Number of PCs

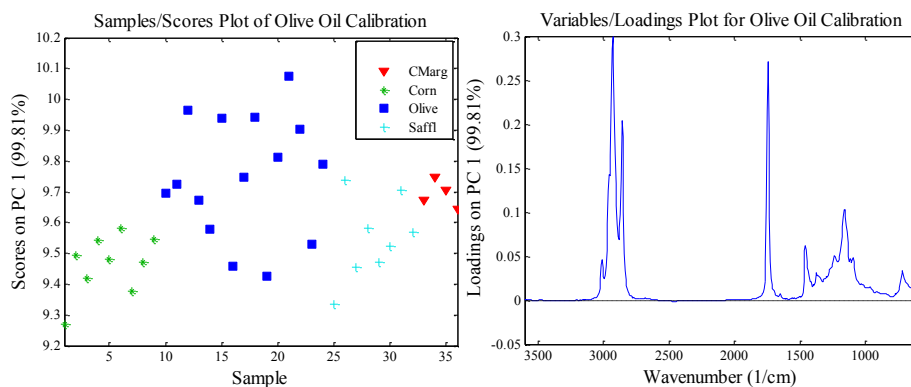


- It's not always easy
- In exploratory analysis it doesn't really matter
- Compare total % **variance of model** with **error of data**.
- **Eigenvalue Plot** - PCs before a break.
- **PRESS plot from cross-validation** - PCs at first minimum or near plateau.
- **Chemical intuition** to choose between conflicting results.
 - e.g., do the PCs discriminate?

31



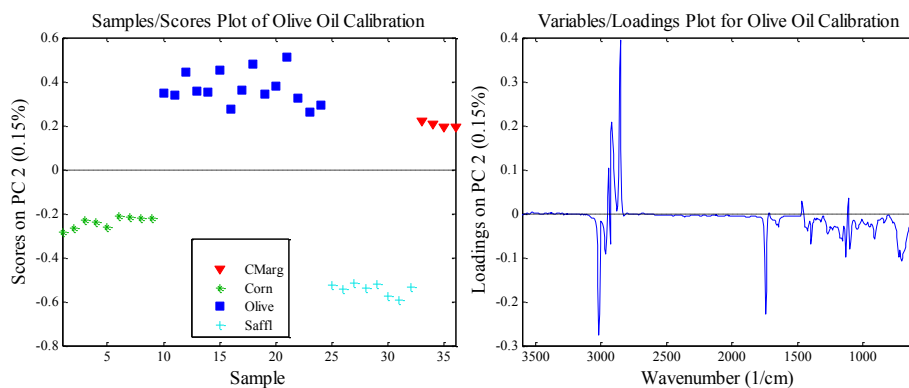
Scores and Loadings, PC 1



32



Scores and Loadings, PC 2



33



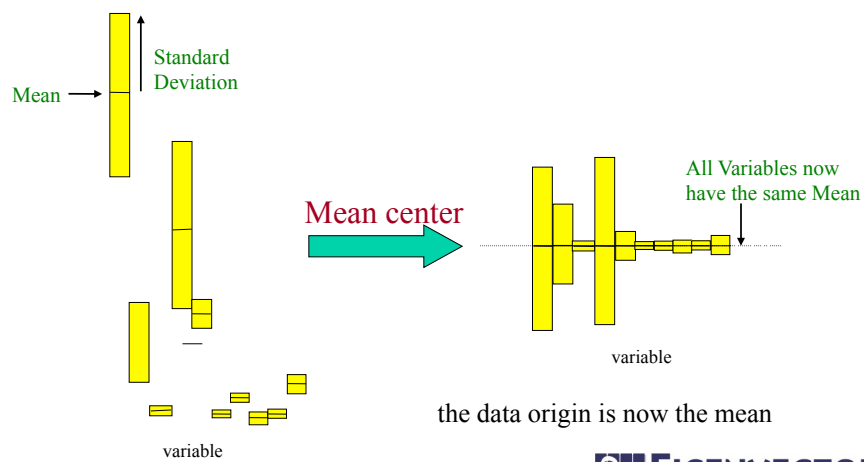
PCA Summary

- No preprocessing
 - PC 1 captured variance that was in the general direction of the mean
 - although it is not strictly the mean of the data
 - PC 2 discriminated the oils
 - some variables associated with differences between the oils were seen on PC 2
 - discrimination wasn't great, can we do better?
 - PCA is designed to capture sum-of-squares from the origin
 - that's why PC 1 was in direction of the mean!

34



Repeat the PCA with Mean-Centering

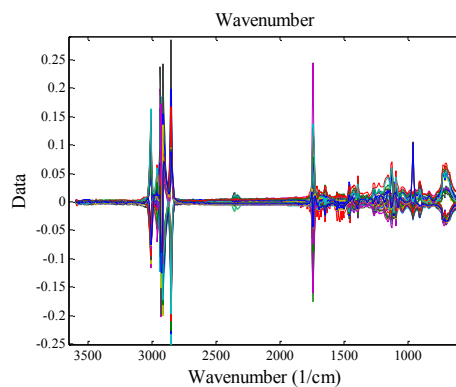
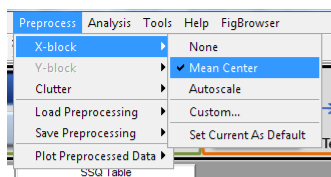


35

EIGENVECTOR
RESEARCH INCORPORATED

Mean-Centering

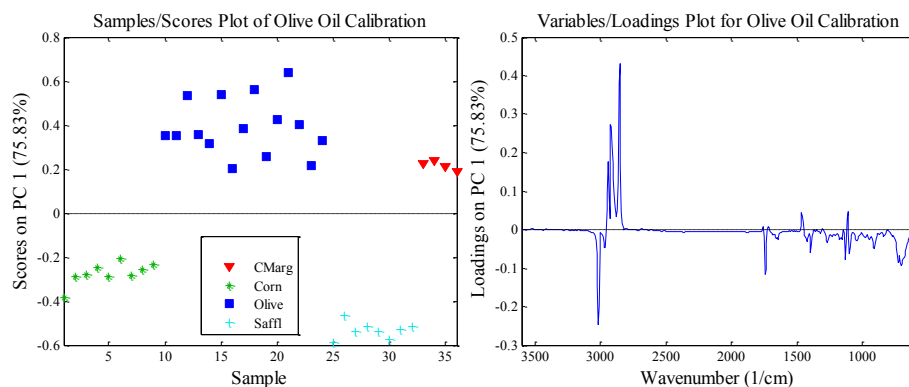
- PCA now captures sum-of-squares about the mean of the data



36

EIGENVECTOR
RESEARCH INCORPORATED

Scores and Loadings, PC 1



37

EIGENVECTOR
RESEARCH INCORPORATED

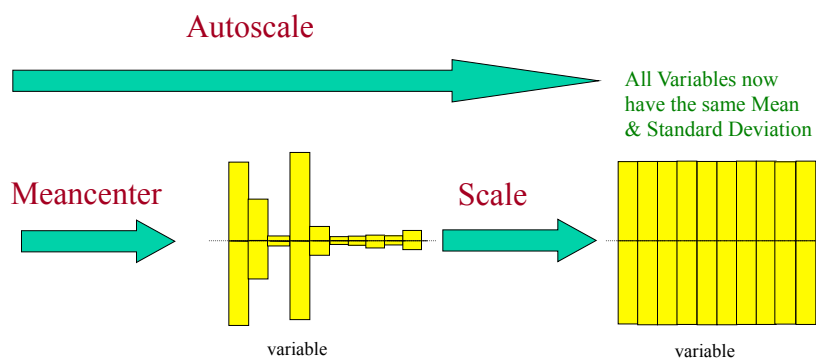
PCA Summary

- Mean-centering
 - removing the mean now focused PC 1 on variance about the mean and PC 1 discriminated the oils
 - we're bringing relevant variance closer to the top
 - median-centering can be used when there are expected to be outliers that might influence the mean
 - the outliers are easier to identify and then remove
 - additionally, we identified
 - regions with little or no signal
 - sloping baseline variability
 - can we do better, how about auto-scaling?

38

EIGENVECTOR
RESEARCH INCORPORATED

Autoscaling



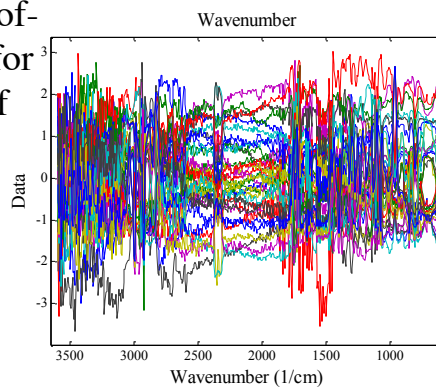
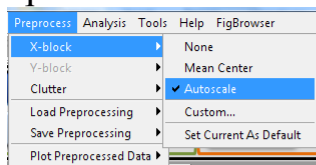
Remember: autoscaling includes mean centering

39



Auto-Scaling

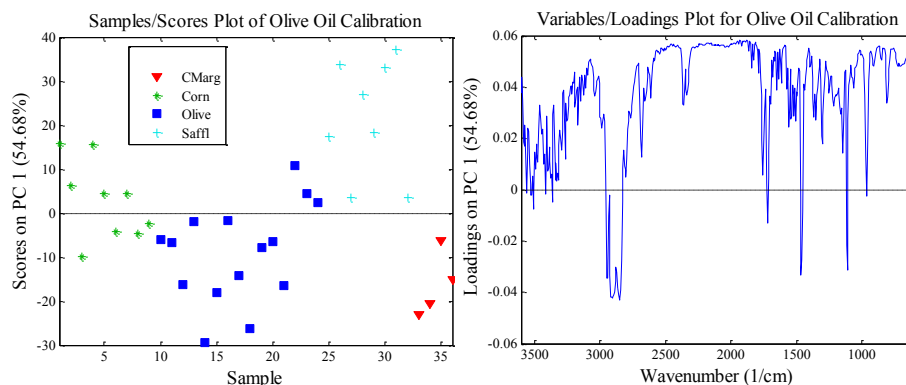
- PCA now captures sum-of-squares about the mean for the data with variables of equal variance



40



Scores and Loadings, PC 1



41



PCA Summary

- Autoscaling
 - PC 1 no longer discriminates
 - giving all the variables an equal weight "blew up" noisy variables that had small signal and subsequently added lot's of sum-of-squares
 - we're bringing irrelevant variance closer to the top
 - for many data sets, autoscaling is a good thing, but not often used in spectra
 - autoscaling ~assumes that each variable has a similar S/N
 - but clearly not the case over the entire spectral range
 - often used when variables are of different units
 - e.g., in engineering applications

42



Q Statistic in PCA

- Recall that the PCA model was truncated to keep only K PCs.
- What about \mathbf{E} ? \mathbf{E} is the lack of fit.
- The Q statistic is the sum-of-squares of each row of \mathbf{E} and is a measure of lack of fit of each sample.
 - It is a measure of the distance from the plane of the PCA model.

$$\begin{array}{c} \text{variables} \\ \hline \mathbf{X} \\ \hline \text{samples} \end{array} = \begin{array}{c} \overline{\mathbf{p}_1} \\ \hline t_1 \end{array} + \begin{array}{c} \overline{\mathbf{p}_2} \\ \hline t_2 \end{array} + \dots + \begin{array}{c} \overline{\mathbf{p}_k} \\ \hline t_k \end{array} + \mathbf{E}$$

43



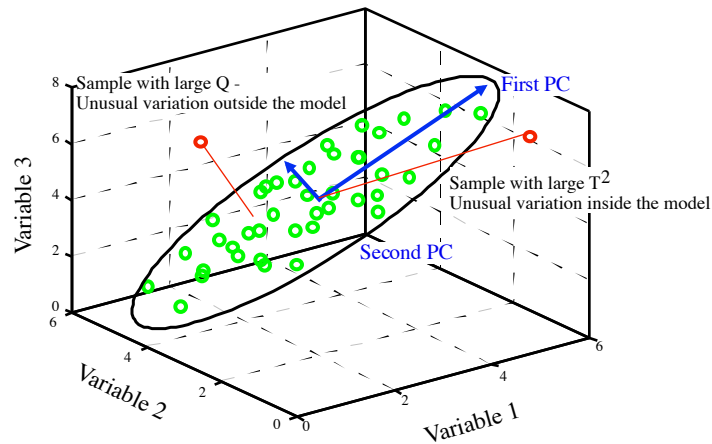
Hotelling's T^2

- Hotelling's T^2 statistic can be calculated from the PCA scores.
- T^2 accounts for the different amounts of variance in each direction to calculate a distance from the origin within the plane of the PCA model.

44



Geometry of Q and T^2



45

 **EIGENVECTOR**
RESEARCH INCORPORATED

Outline

- Introduction
- PCA Review
- PLS Regression Review
 - cross-validation
 - Savitsky-Golay
 - model validation
- Advanced Preprocessing
- Variable Selection
- Summary

46

 **EIGENVECTOR**
RESEARCH INCORPORATED

We Can't Always Measure What We Want*

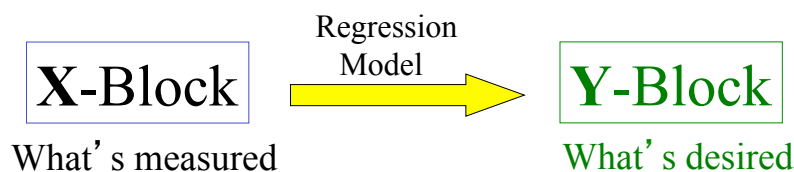
- Often measurements must be made on something else and the property of interest must be inferred from these measurements.
- This is the idea behind inferential sensing where variables are measured that are available in a timely manner to predict something that is more difficult (or more expensive) to obtain.

*"You Can't Always Get What You Want," Rolling Stones, Let it Bleed (1969)

47



Regression



PCA was used to explore the correlation structure within a single data block **X**.

Regression analysis identifies the dependency between two blocks of data **X** and **Y**.

Regression models are often used to obtain estimates (or predictions) for one block of data from the other.

48



Many Forms of Regression

- Classical Least Squares (CLS)
 - Generalized least squares
 - Extended least squares
- Multiple Linear Regression (MLR)
- Principal Components Regression (PCR)
- Partial Least Squares Regression (PLS)

$$\mathbf{X}\mathbf{b} = \mathbf{y} + \mathbf{e} \text{ (this is our focus)}$$

$$\mathbf{X}\mathbf{B} = \mathbf{Y} + \mathbf{E} \text{ (PLS-2: multivariate } \mathbf{Y})$$

49



PLS Description

- PCA decomposes \mathbf{X} into factors called PCs
- PLS decomposes \mathbf{X} (and \mathbf{Y}) into latent variables
- Selection of the number of LV's is ~more important in PLS than in PCA but it's also a bit easier
 - Cross-validation

50



Cross-Validation

- Divide data set into J sample subsets to leave out one at a time.
- For *each subset*:
 - build a PLS model using all samples in the *remaining* subsets (i.e., build J models) and using different numbers of LVs (1,2,...)
 - apply the model to predict the J^{th} subset samples
 - calculate PRESS (Predictive Residual Sum of Squares) for the subset samples and sum over all J subsets and LVs:

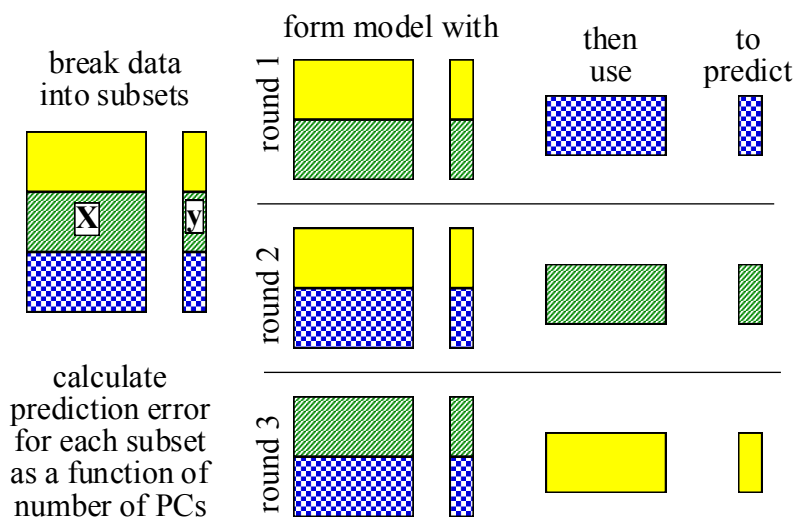
$$\mathbf{e}^2 = (\mathbf{y} - \mathbf{X}\mathbf{b})^2$$

- Look for minimum or “knee” in PRESS curve

51



Cross-validation Graphically



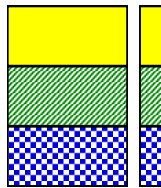
52



Formation of Test Sets



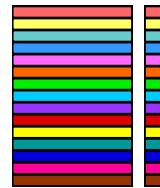
“Venetian blinds” - OK when data already in random order



contiguous blocks-best for time series



random selection-usually repeated several times



leave-one-out, used when not much data available

What else?

Custom selection, based on prior knowledge!

53



Cross-validation Considerations

- Cross-validation method selection criteria
 - *Number* of objects in dataset
 - *Order* of objects in dataset
 - *Objective* of cross-validation (specific type of error?)
 - Presence/absence of *replicates*
- “Traps” to avoid
 - “Repeat sample trap”
 - Repeat measurements in both model and test set
 - “External subset selection trap”
 - Test set “space” outside of model set “space”

54



Cross-Validation Rules of Thumb

- Divide data set into \sim square root of number of samples subsets.
- “Genuine Replicates” can be split between the Learning and Test Sets
 - “'genuine replicates' are repetitions which are subject to all the sources of error that affect runs made at different experimental conditions”*
 - If simple repeat measurements, keep them together, *i.e.* have all in either the Learning Set or Test Set.

*Box, Hunter, and Hunter, “Statistics for Experimenters”, Wiley (1978)

55



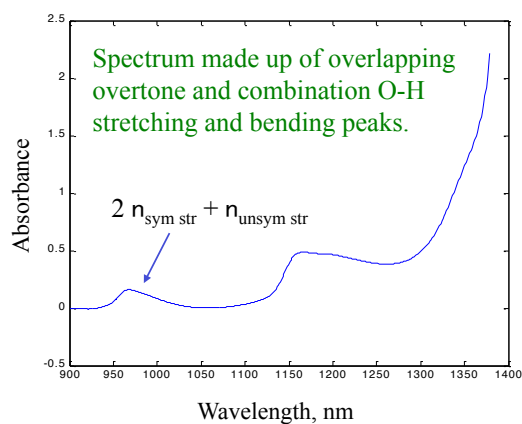
Example Application of PLS

- Estimate the concentration of NaOH in aqueous caustic brine solutions using SW-NIR
 - measured 12 solutions of NaCl and NaOH in water
 - peaks shift with changes in NaCl , NaOH and temperature, T
- Since T will vary in the application, T variation must be included in the Learning Set
 - although T need not be known to calibrate for NaOH, it must vary in the Learning Set for the model to be robust to changes in T

56



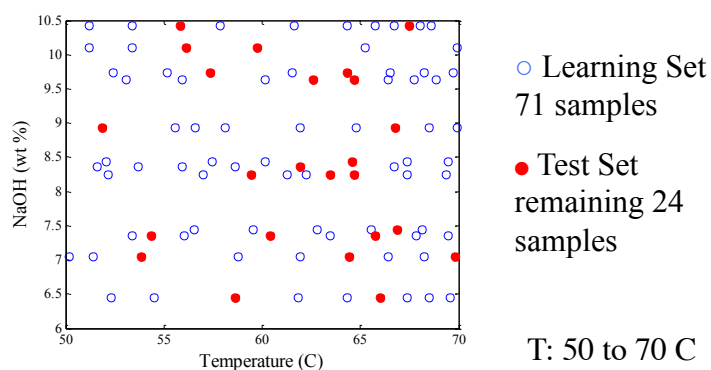
Typical SW-NIR Spectrum of Caustic Brine Solution



57



Learning and Test Sets Randomly Selected Samples



58



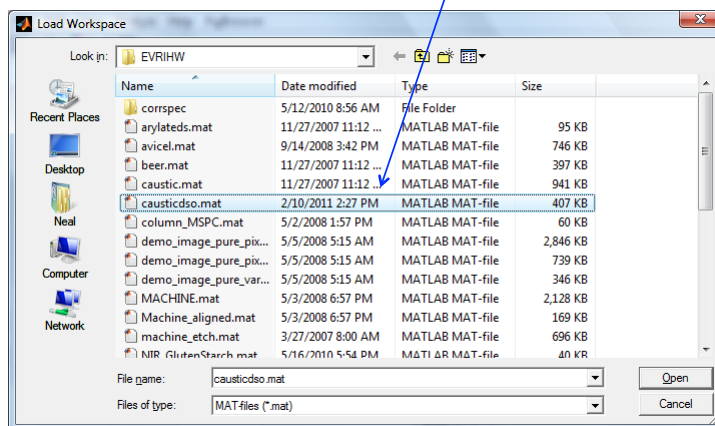
PLS Example

- How to preprocess?
 - just mean-center for now
- How many latent variables?
 - cross-validation using venetian blinds
 - split the data $\sqrt{71} \sim 8$ times
 - examine out to 20 LVs (expect true number < 20)

59



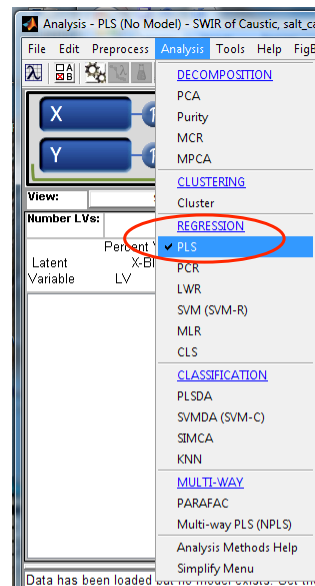
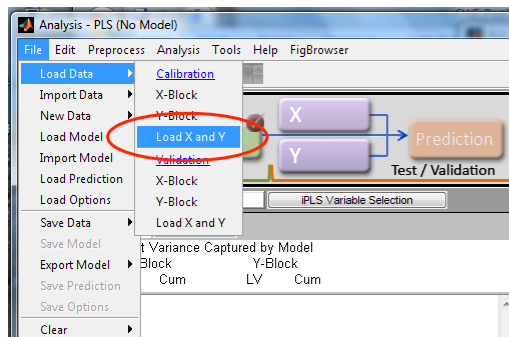
- Analysis: File: Clear: All
- Browser: File: Clear Workspace
- Browser: File: Load Workspace: [causticdso.mat](#)



60



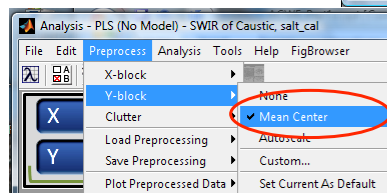
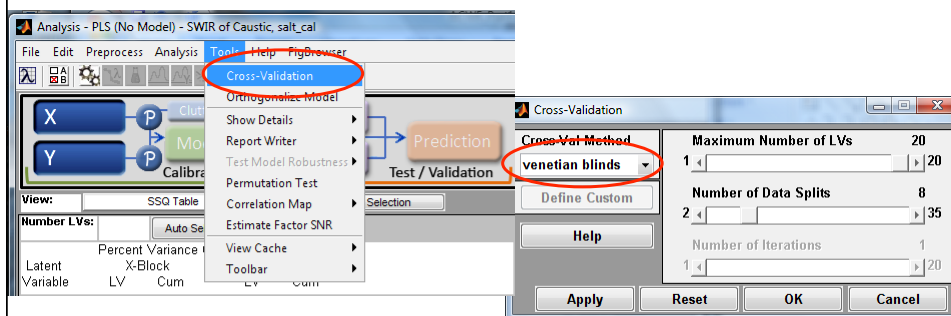
- Analysis: File: Load X and Y
 - load xcal and ycal
- Analysis: Analysis: PLS



61



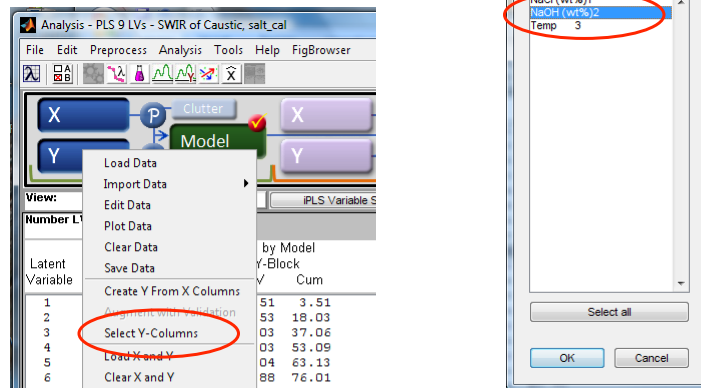
- Analysis: Tools: Cross-Validation
- Cross-Validation: venetian blinds: OK
- Analysis: Preprocess: X,Y-Block: Mean Center



62

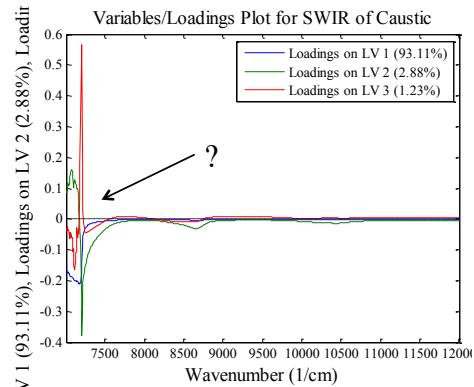
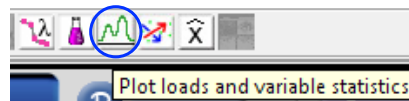
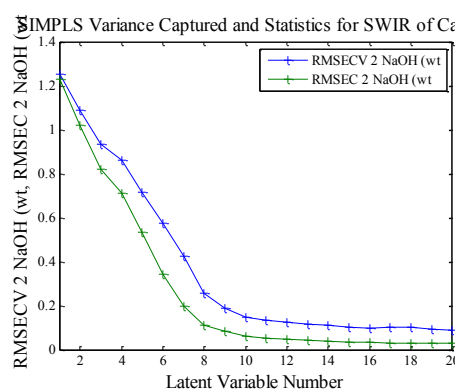
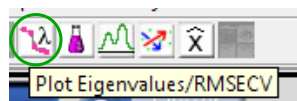


- Analysis: right-click Y: Select Y-Columns
- Select NaOH (wt%)2: OK
- Analysis: Model



63

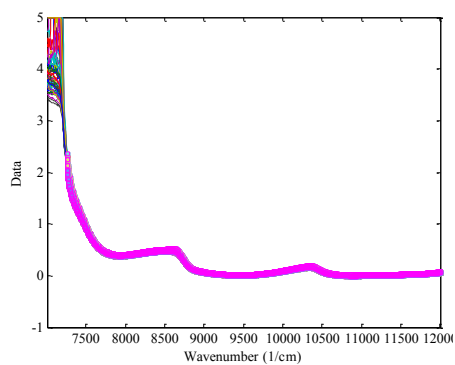
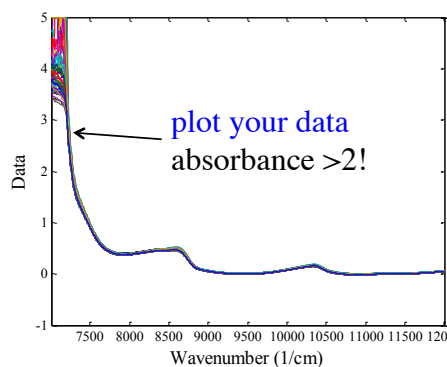
EIGENVECTOR
RESEARCH INCORPORATED



64

EIGENVECTOR
RESEARCH INCORPORATED

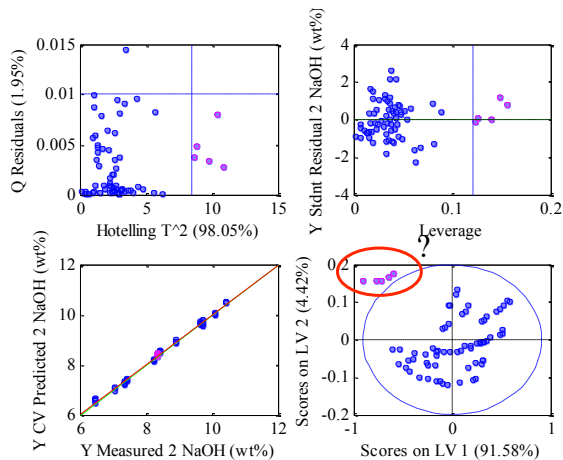
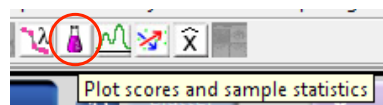
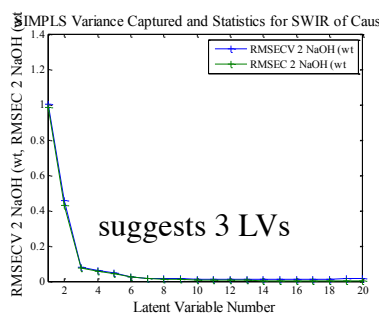
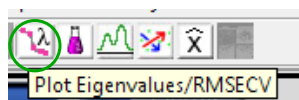
- Analysis: right-click X: Plot Data
- Plot Controls: Y: Data
- Plot Controls: Select {choose values < 2}
- Plot Controls: Edit: Include Only Selection



Variable selection starts by using what's known about the physics and chemistry of the measurement system.

65

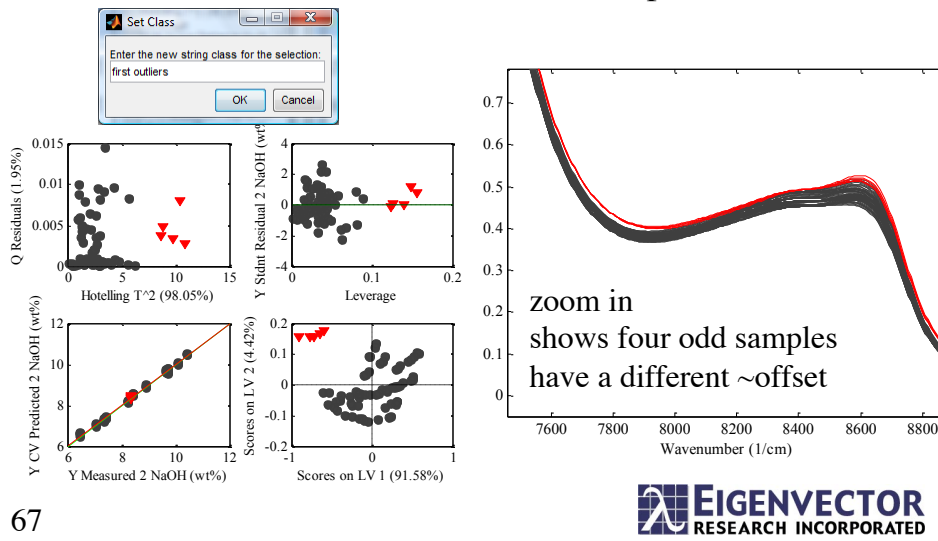
EIGENVECTOR
RESEARCH INCORPORATED



66

EIGENVECTOR
RESEARCH INCORPORATED

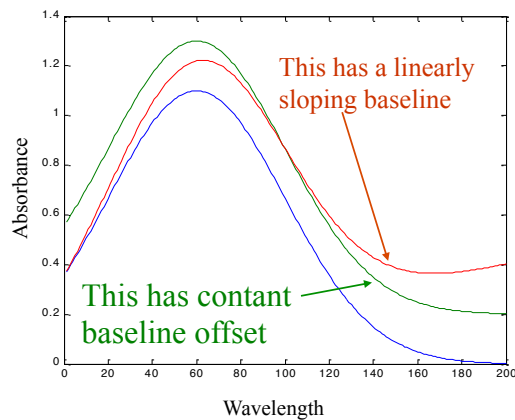
- Plot Controls: Select {select the four odd samples}
- Plot Control: Edit: Set Class of Selected: "1 outliers": OK
- Plot Controls: data {select several samples}



67

Baseline Problems Example

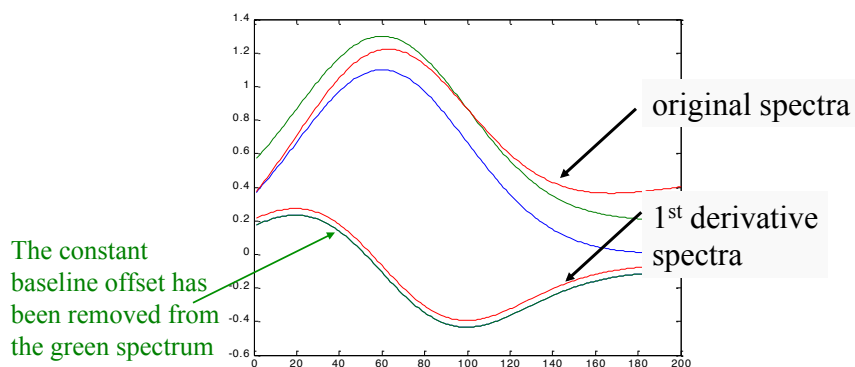
three identical spectra, except:



68

Using Derivative Spectra

Take 1st derivative of the Three Spectra

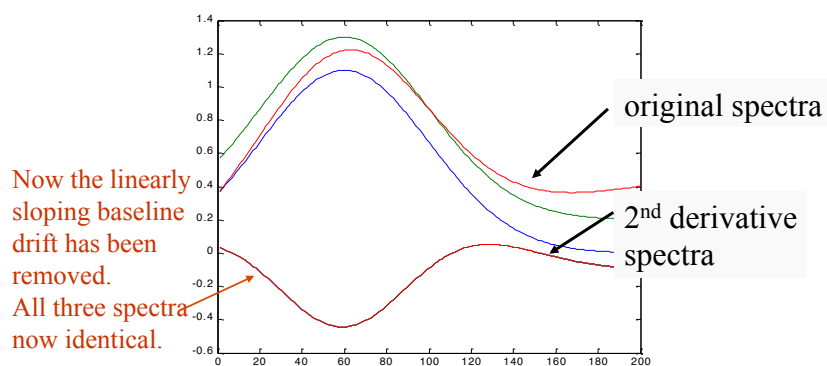


69

EIGENVECTOR
RESEARCH INCORPORATED

Using Derivative Spectra

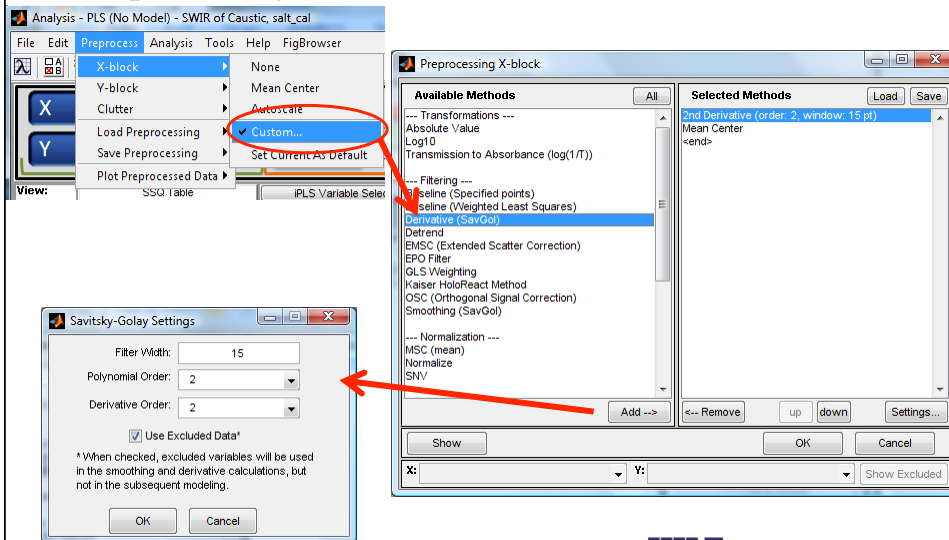
Take 2nd derivative of the Three Spectra



70

EIGENVECTOR
RESEARCH INCORPORATED

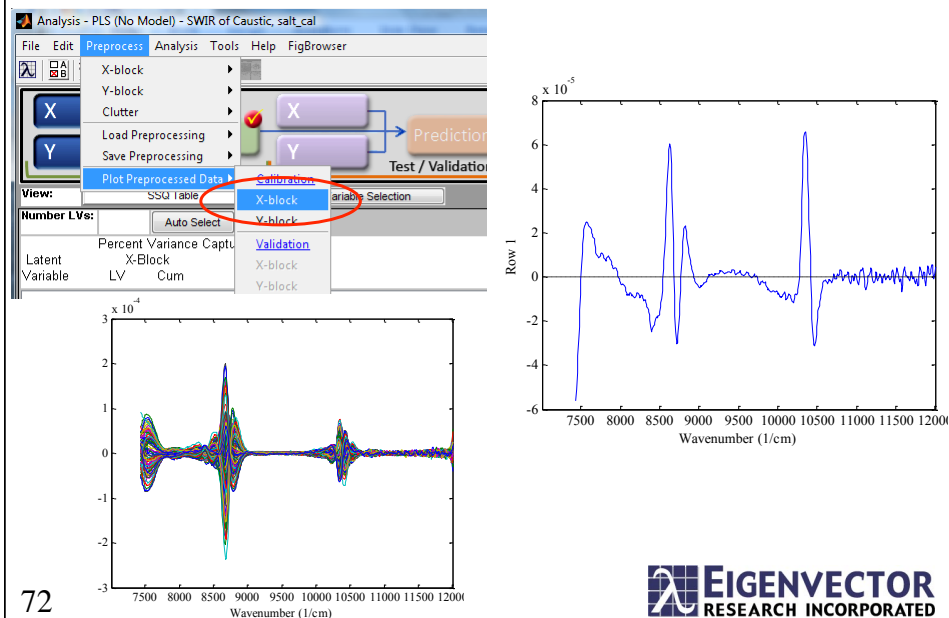
- Analysis: Preprocess: X-block: Custom...
- Preprocessing X-block: Derivative: Add



71

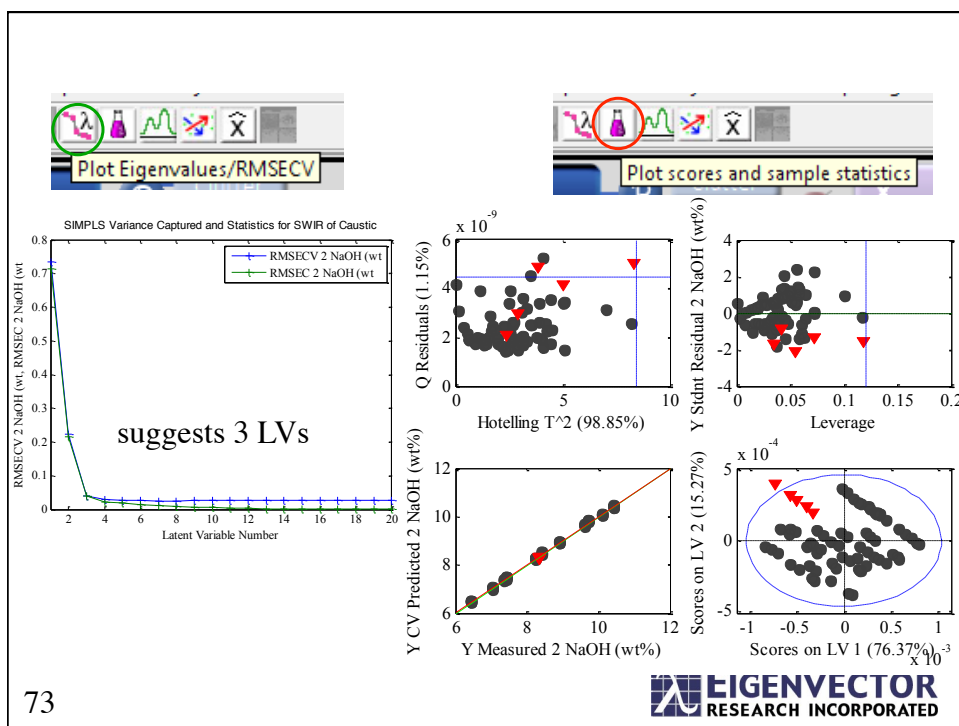
EIGENVECTOR
RESEARCH INCORPORATED

- Analysis: Preprocess: Plot Preprocessed Data: X-block



72

EIGENVECTOR
RESEARCH INCORPORATED

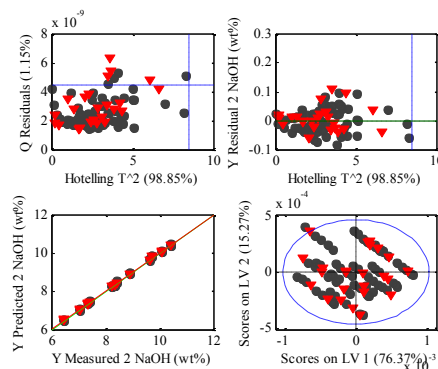


- Analysis: Load X, and Y-block test data / validation data and click Model

Linear regression model using
Partial Least Squares calculated with the SIMPLS algorithm
Developed 11-Feb-2011 09:42:59.84
Author: Neal@NEAL-VAIO_07
X-block: SWIR of Caustic, Test Set 24 by 593 Included: [1-24] [57-649]
Included (in axis units): [n/a] [7436.1-12002]
Preprocessing: 2nd Derivative (order: 2, window: 15 pt), Mean Center
Y-block: salt_cal 24 by 1
Included: [1-24] [2]
Preprocessing: Mean Center
Num. LVs: 3
Cross validation: venetian blinds w/ 8 splits
RMSEC: 0.0387602
RMSECV: 0.0410341
RMSEP: 0.0410703
Bias: 0
CV Bias: -0.00018375
Pred Bias: 0.00680956
R² Cal: 0.999081
R² CV: 0.99897
R² Pred: 0.998991

Percent Variance Captured by Regression Model

Comp	---X-Block---		---Y-Block---	
	This	Total	This	Total
1	76.37	76.37	68.95	68.95
2	15.27	91.64	28.20	97.14
3	7.22	98.85	2.76	99.91



PLS Example Summary

- Variable Selection
 - use *what you know* to remove irrelevant variables
 - plot your data
- Preprocessing
 - mean-centering was used to remove overall offsets
 - not mean-centering is a force fit through zero
 - Savitzky-Golay smoothing and derivatives were used to remove offsets and slopes in the spectra
- Fit and Prediction *are not* the same thing
 - model validation is very important and continues...

75



Before Applying Models to Real Unknowns

**Validate Them Thoroughly With a Well
Designed Test Set!!**

Models Do Not Last Forever

**Revalidate Them Often and Rebuild
Them If Necessary.**

76



Outline

- Introduction
- PCA Review
- PLS Regression Review
- Advanced Preprocessing
 - clutter
 - GLS, MSC, EMSC, SNV, normalization
- Variable Selection
- Summary

77



What is Clutter?

- Signal is defined as the measurement associated with the target of interest.
 - e.g., it is the part of the FTIR spectrum corresponding to discriminating the olive oils, or
 - the relationship between temperatures in a distillation column and the tray compositions
- Clutter is everything else in the measurement
 - interferences
 - noise

78



Measured Signal

- Measured signal includes target
- and Clutter (X-, Y-block, ...)



$$\left[\mathbf{X}_{\text{target}} + \mathbf{X}_{\text{clutter}} \pm \delta \mathbf{X} \right] \mathbf{b} = \left[\mathbf{y}_{\text{target}} + \mathbf{y}_{\text{clutter}} \pm \delta \mathbf{y} \right] + \mathbf{e}$$

- Use physics to create a linear relationship
 - non-linearity w/in X-block adds factors (digs deeper into noise)
 - non-linearity between X- and Y-blocks adds error

79



Sources of Clutter

- Instrument physics
 - offset and gain changes, drift, hardware changes, smile, wavelength registration, temperature, humidity, operator ...
- Sample / sampling
 - interferences chemical and physical
 - presence of other analytes
 - pathlength changes, particulate and size distribution changes, ...

80



Preprocessing Objective

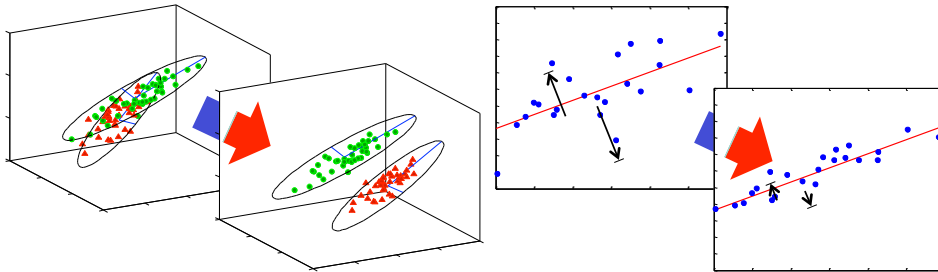
- Typical analysis methods of interest are based on maximizing capture of sum-of-squares or minimizing least-squares.
- The objective of preprocessing is to minimize variance due to clutter so that the analysis can focus on signal of interest
 - Clutter: sensor noise and the confounding effects of interferences
 - Radar Clutter Definition: (DOD, NATO) Unwanted signals, echoes, or images on the face of the display tube, which interfere with observation of desired signals.

81



Advance Preprocessing

- Introduce concepts and methodologies to maximize signal-to-clutter for use in PCA and PLS
 - maximize between-class distance / within-class distance
 - minimize the prediction error



82



Advanced Chemometrics

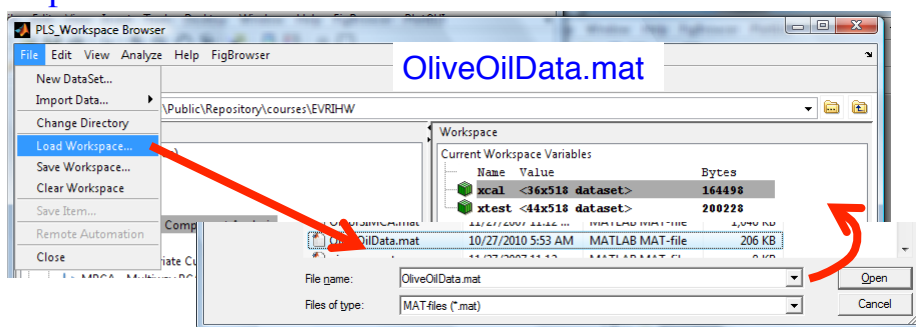
- Advanced concepts combine our understanding of the physics and chemistry of the system, and knowledge of how the mathematical tools work to provide better experimental designs and to ...
- maximize signal-to-noise → **signal-to-clutter**
- Data analysis and preprocessing should *not* be treated as a black box

83



Reload OliveOilData.mat

- Analysis: File: clear all
- Browser: Clear Workspace
- Analysis, PCA, mean-centering, cross-val none, and plot scores

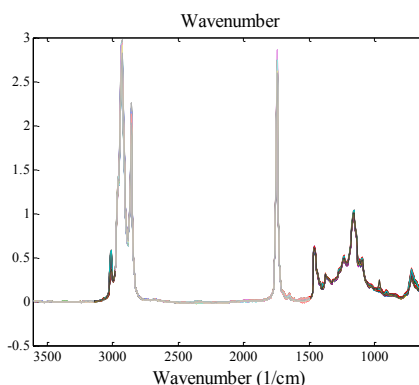
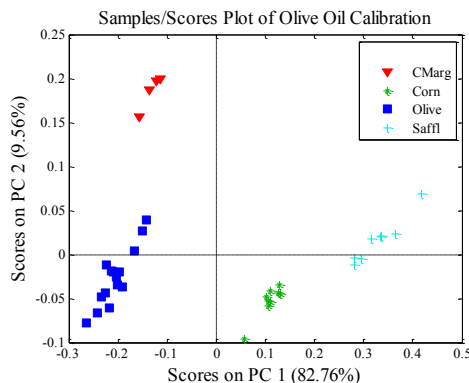


84



This scores plot shows better discrimination than our previous PCA with the same preprocessing, why?

Plot your data
Plot Controls: View: Excluded Data
Variable selection removed
'irrelevant variables.'



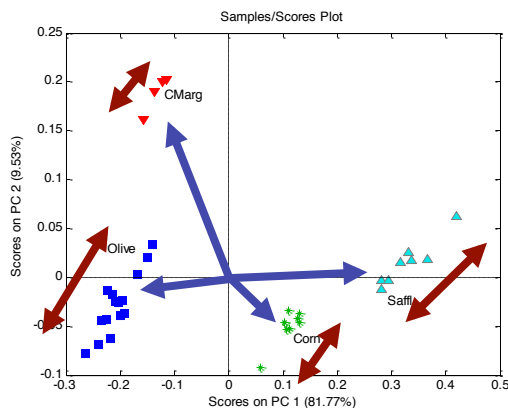
... but can we do better?

85



PCA Results

- PCA shows that the four classes in the calibration data set are separate from each other (high **between class variance**) but ...
- have significant **within class variance**

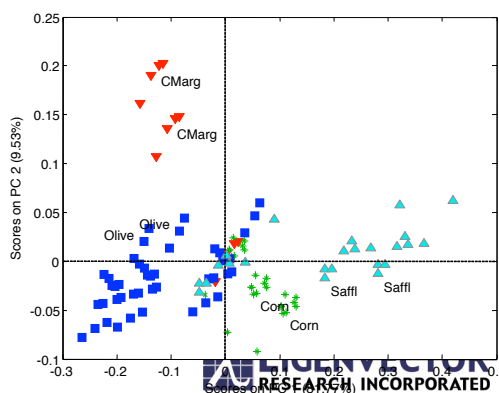


86

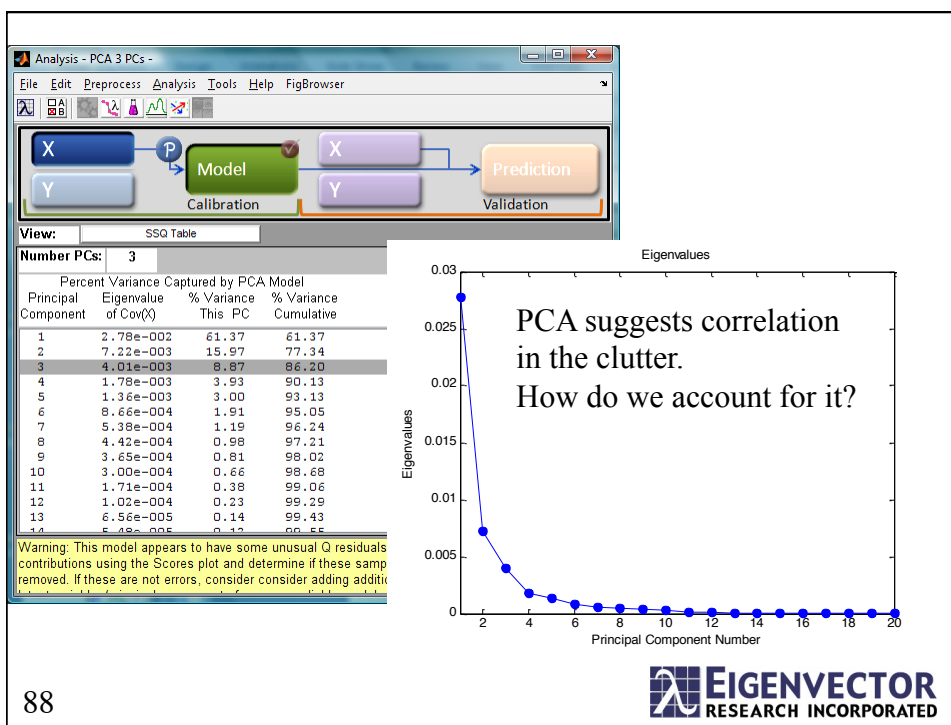


Replicates

- Ideally, replicates would lie on top of each other.
- Variance within each class is clutter variance.
 - Is it random noise? Is the clutter correlated?
- Center each class to its own mean and do PCA on the result.



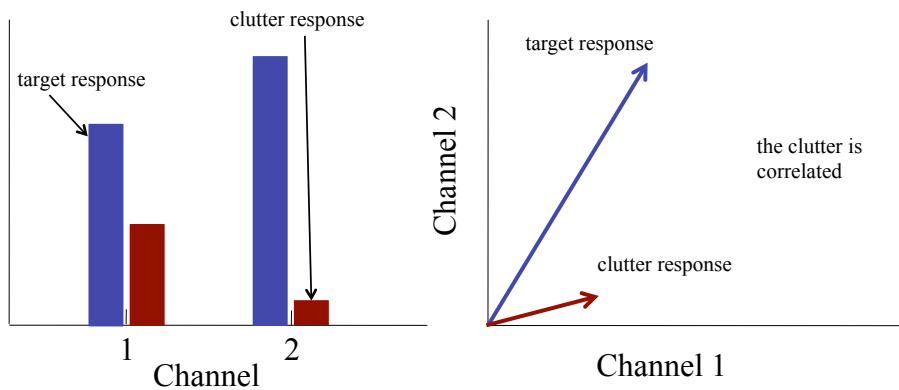
87



88

How does the clutter affect the measurements?

- Imagine a 2-channel spectrometer



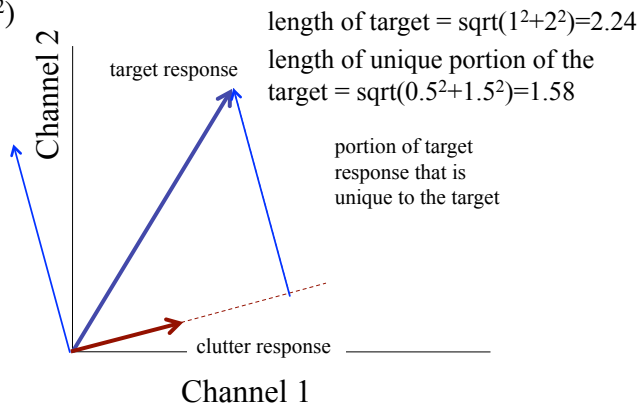
89

EIGENVECTOR
RESEARCH INCORPORATED

How does the clutter affect the measurements?

- characterize the signal as the length of the vector

$$\sqrt{x_1^2 + x_2^2}$$



90

EIGENVECTOR
RESEARCH INCORPORATED

Why is clutter bad?

- The signal-to-clutter is ~proportional to the length of the unique portion of the target's response.
 - in absence of clutter it was 2.24
 - in the presence of clutter it was 1.58
- In regression, clutter-to-signal is related to the estimation error.
 - higher clutter-to-signal → higher estimation error
 - in the presence of clutter the estimation error is 2.24/1.58 times the error when clutter is absent

91



Effect of Clutter

- The effect of clutter is to remove target signal
 - for olive oil example the target signal is the differences between the classes
- Instrument related clutter can be minimized by
 - good instrument design that accounts for the environment (noise+interferences) in which measurements are to be made
 - instrument standardization
 - remove drifts in offsets and gains that adds to the clutter
- Can't always be eliminated → what to do?

92



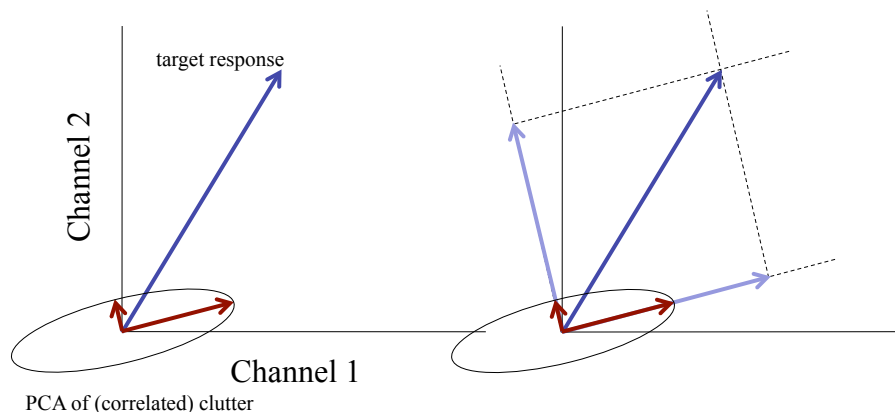
Accounting for Clutter

- One method used to account for clutter is a weighting scheme
 - similar to that used in **generalized least squares (GLS)**
- Autoscaling scales each variable to unit variance
- GLS weighting scales each clutter direction (as determined using PCA) to unit variance
 - directions of high clutter are deweighted
 - directions of low clutter are given more opportunity to allow signal through

93



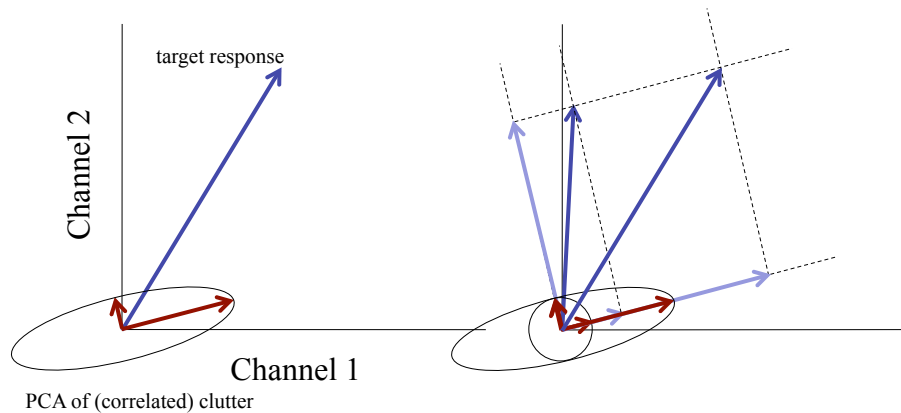
Target Projected onto Clutter Directions



94



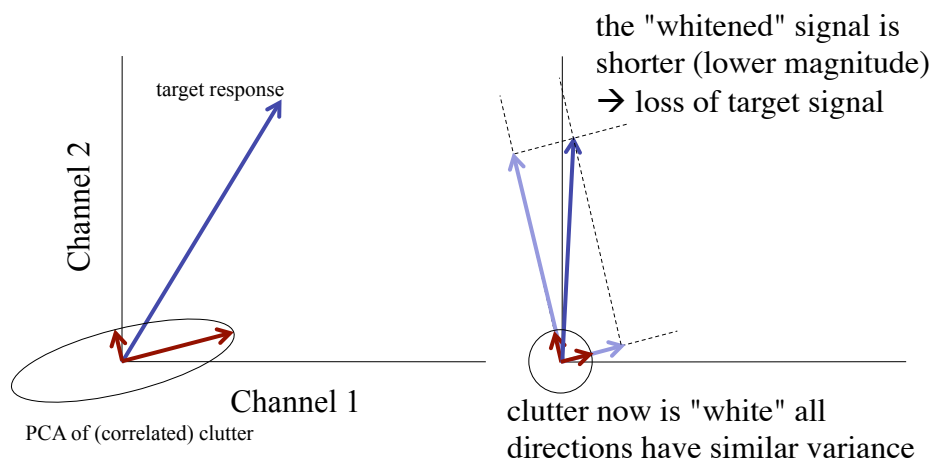
Scale Target by Clutter



95

 **EIGENVECTOR**
RESEARCH INCORPORATED

Whitened Signal

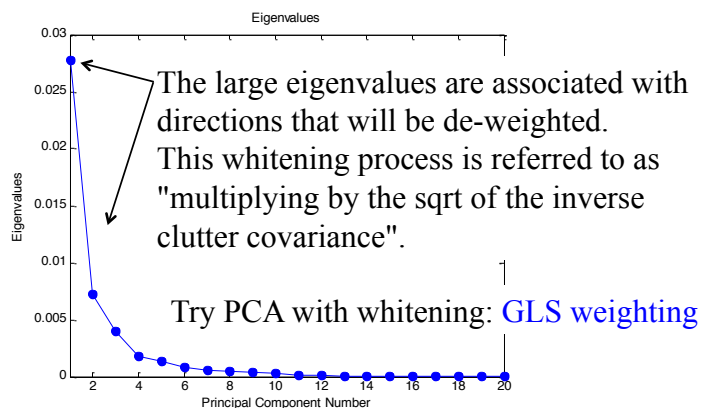


96

 **EIGENVECTOR**
RESEARCH INCORPORATED

Olive Oil Clutter

Eigenvalue distribution of the within class variance.



Based on concepts outlined by Aitken, A., "On Least Squares and Linear Combinations of Observations", *Proceedings of the Royal Society of Edinburgh*, 1935, **55**, 42-48, and used in Maximum Noise Fractions Green AA, Berman M, Switzer P, Craig MD (1988) *IEEE Trans Geosci Remote Sens* 26:65-74

97



Click Preprocessing Shortcut
GLS Weighting: Add: Details:
x-block classes
Mean Center: Add: OK

Analysis - PCA (No Model) - Olive Oil Calibration

View: **Preprocessing X-block**

Number PCs: 1

Available Methods:

- Derivative (SavGol)
- Detrend
- BMSC (Extended Scatter Correction)
- EPO Filter
- GLS Weighting**
- Kaiser-Holm-React Method
- OSC (Orthogonal Signal Correction)
- Smoothing (SavGol)
- Normalization ---
- MSC (mean)
- Normalize
- SNV
- Scaling and Centering ---
- Autoscale
- Group Scale
- Log Decay Scaling
- Mean Center
- Median Center
- Multiscale
- Multiscale

Selected Methods:

- GLS Weighting (classes alpha 0.02)
- Mean Center

Declutter Settings

Clutter Source:

- ☐ automatic
- ☐ y-block gradient
- ☒ x-block classes
- ☐ external data

Algorithm:

- ☒ GLSW
- ☐ EPO
- ☐ EMM / ELS
- ☐ None (disable filter)

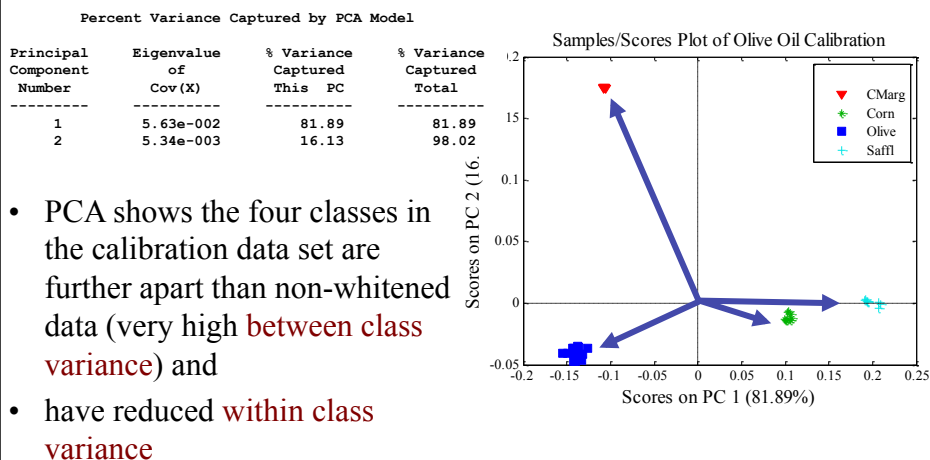
Declutter Threshold: 0.02

Number of PCs: 1

98



PCA of Whitenes Spectra



- PCA shows the four classes in the calibration data set are further apart than non-whitened data (very high **between class variance**) and
- have reduced **within class variance**

99

EIGENVECTOR
RESEARCH INCORPORATED

Olive Oil Samples

Learning set: xcal Start with this data set

Corn Oil	9 samples	(#1-9)
Olive Oil	15 samples	(#10-24)
Safflower Oil	8 samples	(#25-32)
Corn Margarine	4 samples	(#33-36)

Test set: xtest NEW DATA

Corn Oil	9 samples	(#1-9)
Olive Oil	15 samples	(#10-24)
Safflower Oil	8 samples	(#25-32)
Corn Margarine	4 samples	(#33-36)
Corn Oil in Olive Oil	5 samples	(#37-41)
5, 10, 20, 30 & 40%		
Almond Oil	1 sample	(#42)
Peanut Oil	1 sample	(#43)
Sesame Oil	1 sample	(#44)

Analysis - PCA 4 PCs - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

Click Load X Validation Shortcut and load **xtest**

Import

Import from file type:

- Experimental Data (CSV, CSV, CSV, CSV)
- Grams Thermo Galactic File (SPC)
- Hamilton Sundstrand ASF File (ASF, AIF, BKH)
- Hamilton Sundstrand PIONIR File (PDF)
- Horiba JY File (NOS, NGC)
- Image (Workspace/MAT file)
- Image standard (JPG, TIFF, GIF, BMP, PNG)
- JCAMIP (general) (DX, DX)
- Lispix Raw Formatted Image (RAW)
- OPOTEK ENVI Image Format (HDR)
- OPOTEK multi-layer TIFF files (TIFF)
- Physical Electronics RAW Image (RAW)
- Stellarnet ABS File (ABS)
- XML Paste from Clipboard
- XML file (XML)
- XY... Delimited Text Files (TXT, XY)
- Other...

OK Cancel

Load

Create DataSet from

Look In: >> Base Workspace <<

Items

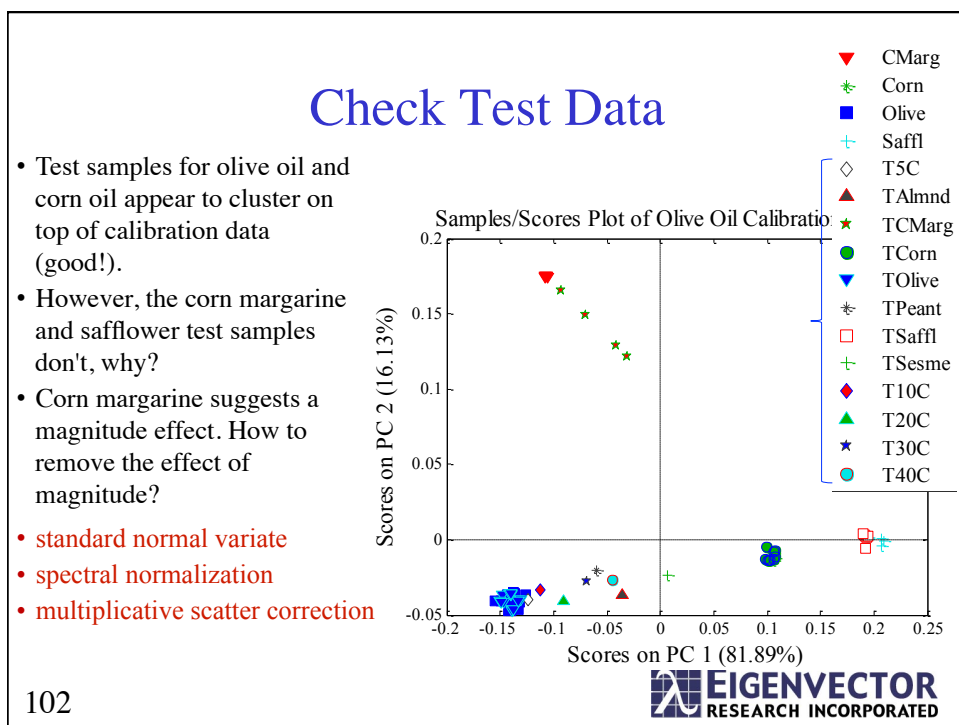
Item	Size	Type
xcal	36x518	(dataset)
xtest	44x518	(dataset)

Item: xtest

From File Load Cancel

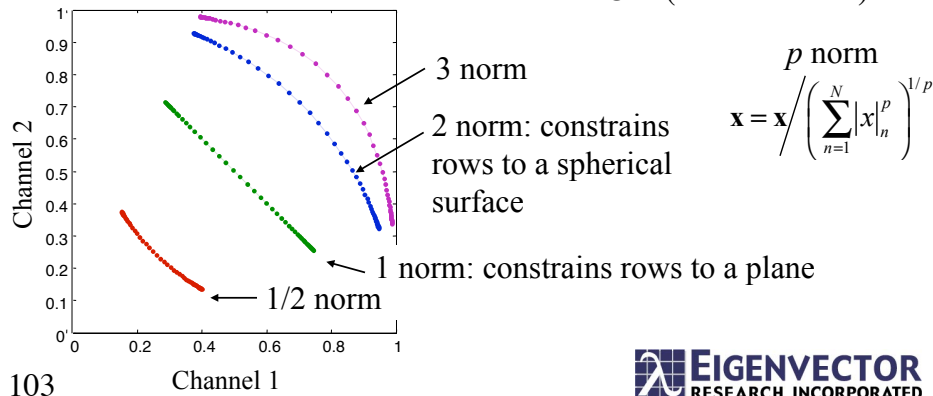
101

EIGENVECTOR
RESEARCH INCORPORATED



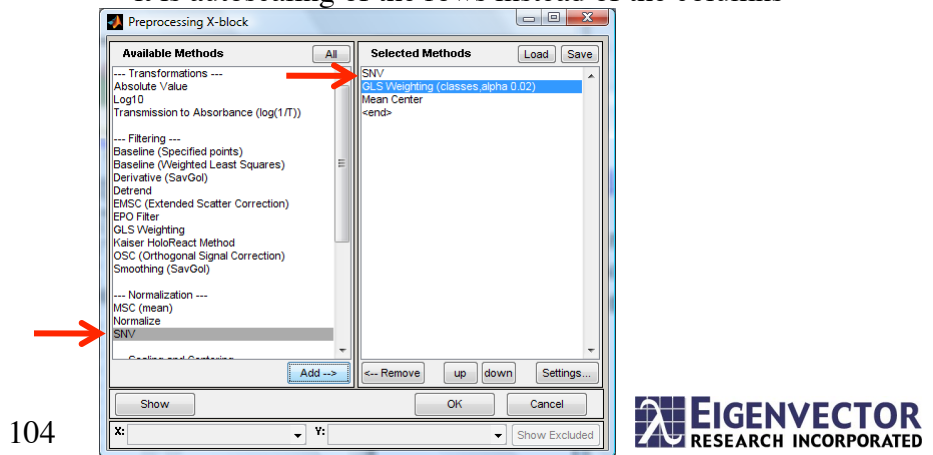
Normalization

- Normalize each row / spectrum
 - 1-norm: normalize to unit AREA (area = 1)
 - 2-norm: normalize to unit LENGTH (vector length = 1)
 - inf-norm: normalize to unit MAXIMUM (max value = 1)

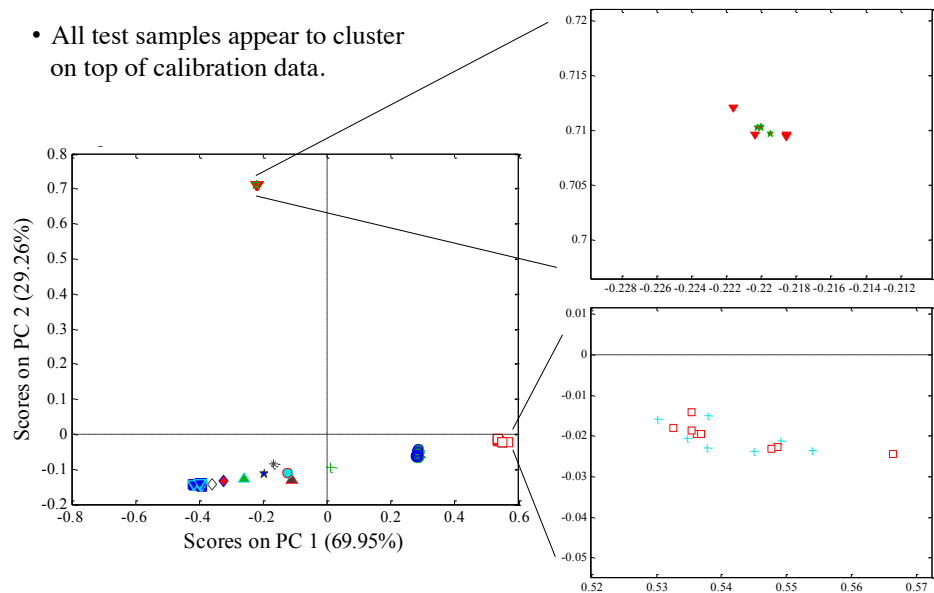


Standard Normal Variate

- Mean-centers each row / spectrum and scales it by its standard deviation
 - it is autoscaling of the rows instead of the columns



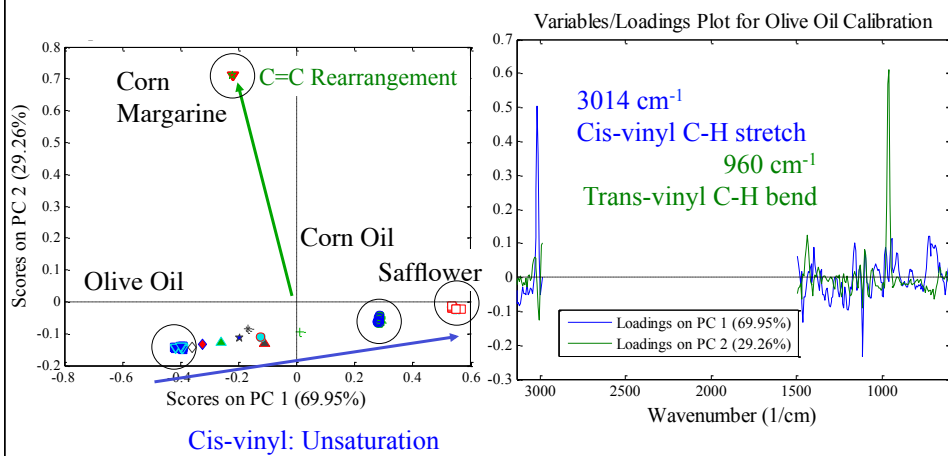
- All test samples appear to cluster on top of calibration data.



105

EIGENVECTOR
RESEARCH INCORPORATED

Interpret Scores and Loadings



106

EIGENVECTOR
RESEARCH INCORPORATED

GLS Weighting

- GLS Weighting of the spectral data accounted for some of the clutter observed in the spectra, but didn't account for magnitude changes.
- SNV was used to account for magnitude changes.
- The result was
 - clusters were further apart and tighter
 - the ratio of between-class to within-class variance was increased making discrimination easier
 - clusters were so tight and far apart that confidence bounds defining each class could be wider

107



Multiplicative Effect in Spectra

- Two spectra are identical except one is a multiple of the other
 - Changing sample pathlength, *e.g.* changing light scattering with particle size.
 - Changing sample density, *e.g.* changing temperature of sample.
 - Changing gain of instrument.
- Plotting a measured spectrum versus a reference spectrum (usually the mean) looks linear

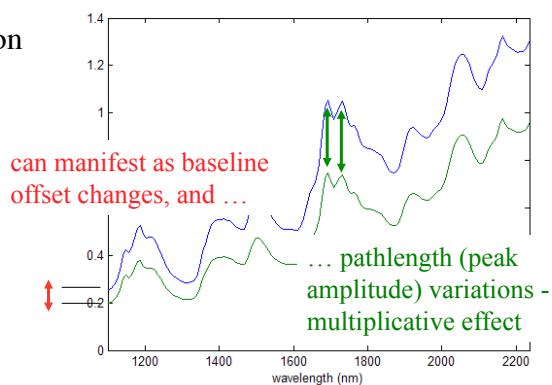
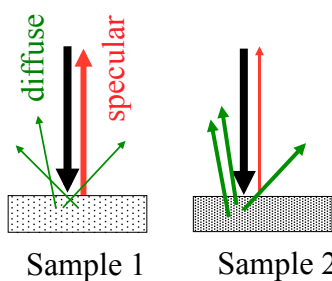
108



Scattering Effects in Reflectance

Caused by variations in:

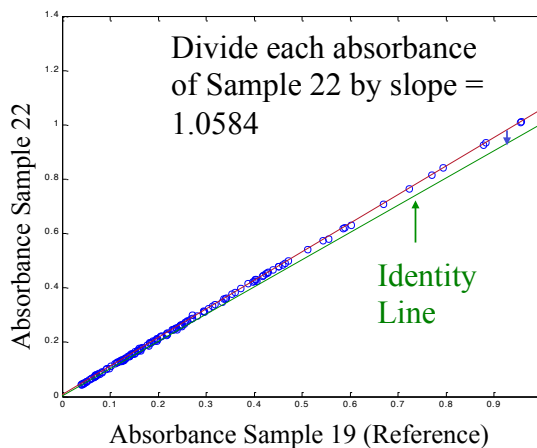
- Particle mean & distribution
- Sample opacity
- Sampling packing density
- Sample placement



109

EIGENVECTOR
RESEARCH INCORPORATED

MSC Multiplicative Signal (Scatter) Correction

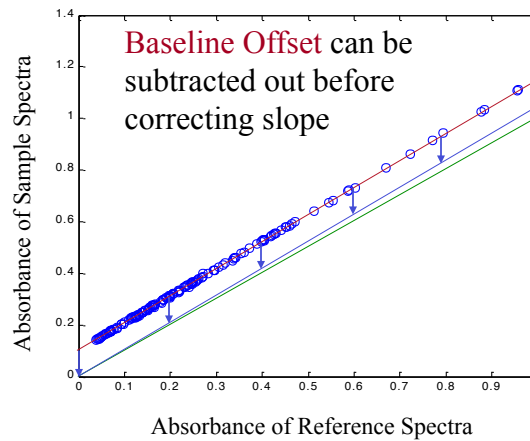


Geladi P, MacDougall D, Martens H., *Appl. Spectrosc.*, **39**(3), 491-500 (1985)

110

EIGENVECTOR
RESEARCH INCORPORATED

If there is also an Offset



111



What to use as a Reference Spectrum?

- Anything we want that looks like the spectra in the Learning Set.
- Usually choose **Mean Spectrum** of the Learning Set.
 - The same spectra subtracted when mean centering.

112



Extended MSC (EMSC)

- Like MSC, EMSC is used to account for offset and gain (multiplicative effects). Also:
 - clutter by using an **extended mixture model**
 - using interference spectra or PCA loadings of clutter data
 - instrument artifacts like slope and smile
 - can allow desired target spectra to 'pass the filter'
- The **extended mixture model** is a classical least squares-like model that is used to explicitly account for clutter using extended least squares.

113



EMSC

Provide spectra of:

- Known **target analytes S**
- Polynomial baselines **P**
- Known **interferences Q**
 - e.g., loadings from a PCA model of clutter
 - the coefficients for each linear effect are estimated using least-squares (indicated by "hat")

$$\mathbf{s}_{2,measured} = \mathbf{s}_{ref}c_{ref} + \mathbf{S}\mathbf{c}_S + \mathbf{P}\mathbf{c}_P + \mathbf{Q}\mathbf{c}_Q \quad \mathbf{P} = \begin{bmatrix} \mathbf{L} & \mathbf{v}^2 & \mathbf{v} & \mathbf{1} \end{bmatrix}$$

$$\mathbf{s}_{2,corrected} = (\mathbf{s}_2 - \mathbf{P}\hat{\mathbf{c}}_P - \mathbf{Q}\hat{\mathbf{c}}_Q) / \hat{c}_{ref} \quad \mathbf{Q} = \text{loadings}$$

114



Analysis - PCA 2 PCs - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

View: SSO Table [Click to import validation X-block data](#)

Number PCs: 2

Percent Variance Captured by PCA Model

Principal Component	Eigenvalue	% Variance of Cov(X)	This PC	% Variance Cumulative
1	4.73e-001	46.31	46.31	46.31
2	3.59e-001	35.12	81.44	81.44
3	6.39e-002	6.26	87.70	87.70
4	4.13e-002	4.05	91.74	91.74
5	2.44e-002	2.39	94.14	94.14
6	2.18e-002	2.13	96.27	96.27
7	1.37e-002	1.34	97.61	97.61
8	7.03e-003	0.69	98.29	98.29
9	5.49e-003	0.54	98.83	98.83
10	2.60e-003	0.25	99.09	99.09
11	2.21e-003	0.22	99.30	99.30
12	1.27e-003	0.12	99.43	99.43
13	1.16e-003	0.11	99.54	99.54
14	7.05e-004	0.08	99.62	99.62

A model has been calibrated from the data. Review the model using the toolbar button(s) save the model (File) preprocess

```
>> z = xcal;
>> for i1=1:4
    z.data(find(xcal.class{1}==i1), :) = mncn(z.data(find(xcal.class{1}==i1), :));
end
>> z.description = char(z.description, 'Each class center to its own mean. ');
>> p = zeros(2, 518);
>> p(:, z.include{2}) = pcam.loads{2}';
```

PCA of clutter.

PCA of calibration data with classes centered to class mean.

Keep 2 PCs to model the clutter.

Save model to pcam

115



Analysis - PCA (No Model) - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

View: SSO Table [Click to import validation X-block data](#)

Number PCs: 2

Percent Variance Captured by PCA Model

Principal Component	Eigenvalue	% Variance of Cov(X)	This PC	% Variance Cumulative
1	4.73e-001	46.31	46.31	46.31
2	3.59e-001	35.12	81.44	81.44
3	6.39e-002	6.26	87.70	87.70
4	4.13e-002	4.05	91.74	91.74
5	2.44e-002	2.39	94.14	94.14
6	2.18e-002	2.13	96.27	96.27
7	1.37e-002	1.34	97.61	97.61
8	7.03e-003	0.69	98.29	98.29
9	5.49e-003	0.54	98.83	98.83
10	2.60e-003	0.25	99.09	99.09
11	2.21e-003	0.22	99.30	99.30
12	1.27e-003	0.12	99.43	99.43
13	1.16e-003	0.11	99.54	99.54
14	7.05e-004	0.08	99.62	99.62

Data has been loaded. Preprocess and T and edited with th

Preprocessing X-block

Available Methods

- Transformations ---
 - Absolute Value
 - Log10
- Filtering ---
 - Baseline (Specified points)
 - Baseline (Weighted Least Squares)
 - Derivative (SavGol)
 - Detrend
 - EMSC (Extended Scatter Correction)
 - EPO Filter
 - GLS Weighting
 - OSC (Orthogonal Signal Correction)
 - Smoothing (SavGol)
- Normalization ---
 - MSC (mean)
 - Normalize
 - SNV
- Scaling and Centering ---
 - Autoscale
 - Group Scale
 - Log Decay Scaling
 - Mean Center
 - Median Center
 - Multway Center
 - Multway Scale

Selected Methods

- EMSC (Extended Scatter Correction)
- Mean Center
- <end>

Add --> <-- Remove up down Settings...

Show OK Cancel

X: Y: Show Excluded

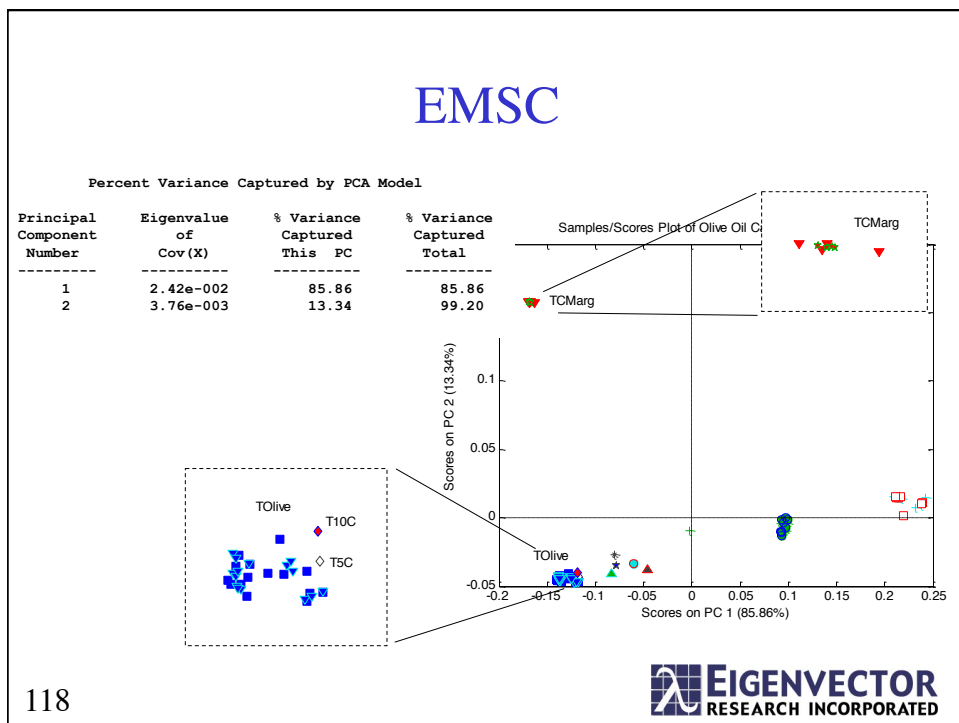
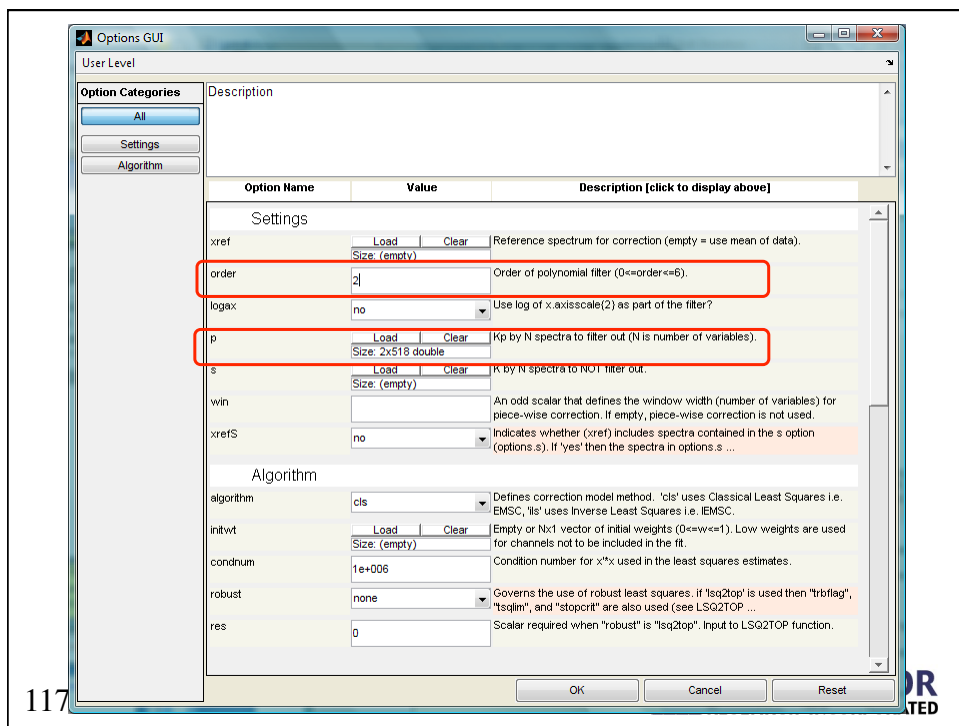
Click Preprocessing Shortcut

EMSC (Extended Multiplicative Scatter Correction)

Settings...

116





EMSC Summary

- EMSC attempts to account for clutter explicitly
 - e.g., model clutter with basis vectors (e.g., PCA loads)
 - analyst takes control of the model
 - requires good use of measurements: clutter and target spectra
 - use what you know!
 - interpretable
 - analyst control is more daunting than simple SavGol and MSC, but
 - results are much more interpretable than 2nd derivative spectra
 - Martens H, Stark E., *J. Pharm. and Biomedical Analysis*, **9**, 625–635 (1991).
 - Helland IS, Naes T, Isaksson T., *Chemom. Intell. Lab. Syst.*, **29**, 233–241 (1995).
 - Martens H, Nielsen JP, Engelsen SB., *Anal. Chem.*, **75**(3), 394–404 (2003).
 - Gallagher NB, Blake TA, Gassman PL., *J. Chemometr.*, **19**(5-7), 271–281 (2005).

119



Analysis - PCA 2 PCs - Olive Oil Calibration (2)

File Edit Preprocess Analysis Tools Help FigBrowse

X-block
Y-block

Load Preprocessing
Save Preprocessing
Plot Preprocessed Data

Calibration X-block
Validation Y-block

Prediction

View: SSQ Table

Number PCs: 2

Percent Variance Captured by PCA Model

Principal Component	Eigenvalue of Cov(X)	% Variance This PC	% Variance Cumulative
1	2.42e-002	85.86	85.86
2	3.76e-003	13.34	99.20
3	9.36e-005	0.33	99.53
4	5.83e-005	0.21	99.74
5	1.73e-005	0.06	99.80
6	1.17e-005	0.04	99.84
7	1.00e-005	0.04	99.88
8	8.27e-006	0.03	99.91
9	4.42e-006	0.02	99.92
10	3.90e-006	0.01	99.94
11	3.30e-006	0.01	99.95
12	2.23e-006	0.01	99.96
13	1.77e-006	0.01	99.96
14	1.48e-006	0.01	99.97

Warning: This model appears to have some unusual Q contributions using the Scores plot and determine if the removed. If these are not errors, consider consider adding them back.

Click Preprocessing Menu

Preprocess:Plot Preprocessed Data: Calibration: X-block

Plot:Rows

View:Classes:Oil

Plot Controls

File Edit View Plot FigBrowse

Fig 3: Calibration

X: Wavenumber

Y: Row 29
Row 30
Row 31
Row 32
Row 33
Row 34
Row 35
Row 36

Z: none

Color By...

Plot auto-update

Select Tool

Table Ctrl+T
Numbers Ctrl+U
Labels Ctrl+L

Classes Ctrl+Z

Excluded Data Ctrl+E

Dec clutter Labels Ctrl+K

Label Angle Ctrl+A

Axis Lines

Diagonal 1:1 line

Log Scales

Auto Y-Scale Ctrl+F

Subplots

Duplicate Figure Ctrl+D

Spawn Static View

Dock Controls

Settings...

View: SSQ Table

Number PCs: 2

Percent Variance Captured by PCA Model

Principal Component	Eigenvalue of Cov(X)	% Variance This PC	% Variance Cumulative
4	5.83e-005	0.21	99.74
5	1.73e-005	0.06	99.80
6	1.17e-005	0.04	99.84
7	1.00e-005	0.04	99.88
8	8.27e-006	0.03	99.91
9	4.42e-006	0.02	99.92
10	3.90e-006	0.01	99.94
11	3.30e-006	0.01	99.95
12	2.23e-006	0.01	99.96
13	1.77e-006	0.01	99.96
14	1.48e-006	0.01	99.97

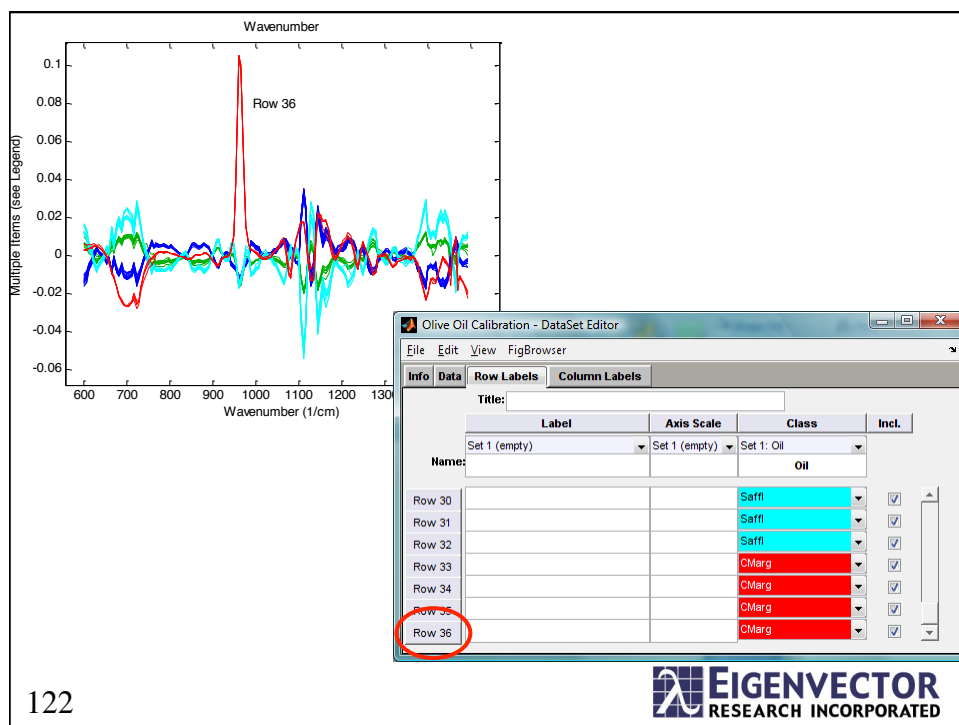
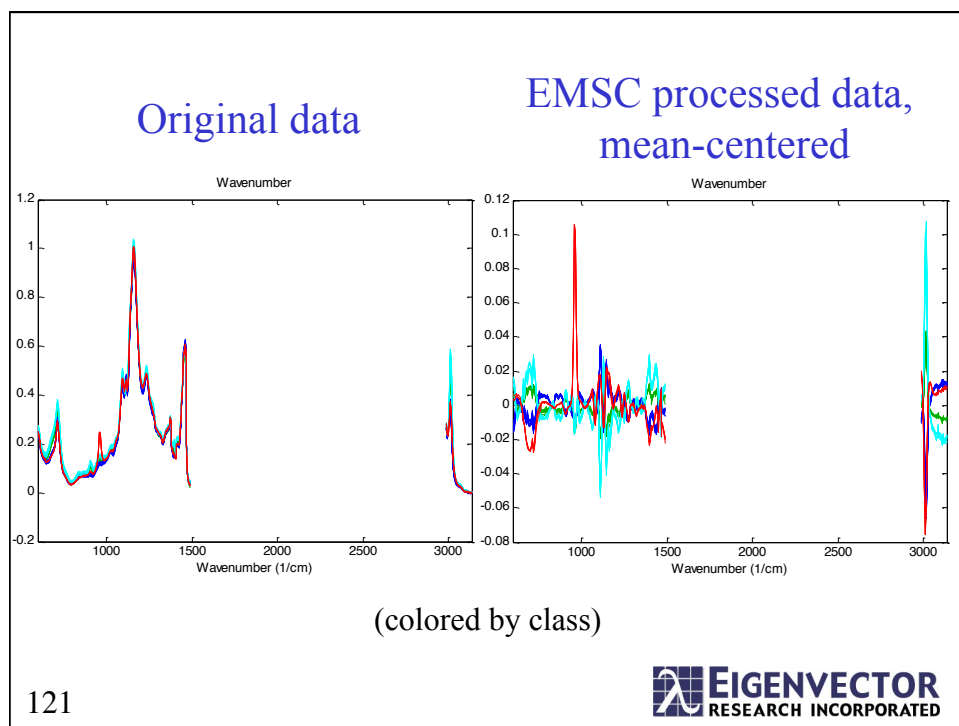
Warning: This model appears to have some unusual Q contributions using the Scores plot and determine if the removed. If these are not errors, consider consider adding them back.

Outline Class Groups

Oil

EIGENVECTOR RESEARCH INCORPORATED

120



Outline

- Introduction
- PCA Review
- PLS Regression Review
- Advanced Preprocessing
- Variable Selection
 - why do it?
 - use what you know!
 - iPLS
- Summary

123



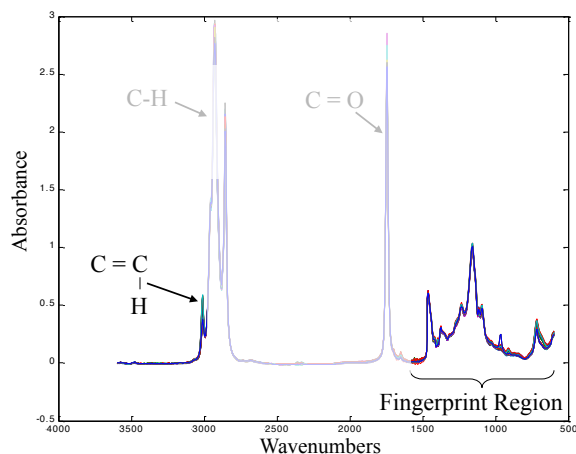
Why Variable Selection?

- Improvement of the model
 - Remove irrelevant, unreliable or noisy variables (clutter)
 - Improve predictions
 - Improve statistical properties
- Interpretation
 - Obtain a model that is easier to understand
- Costs
 - Use fewer measurements to replace expensive or time-consuming one
- Development of fast instruments/routines for on-line control
 - Find wavelength ranges for a filter-based instrument

124



Already performed variable selection
based on *a posteriori* knowledge...



125



Variable Selection Methods

- ***a priori***
 - Choose measurements
- ***a posteriori***
 - Use chemical/physical insight
- **Model based**
 - Look at loadings
- **"Random based"**
 - Genetic algorithms
 - Simulated annealing
- **"Spectral"**
 - i-PLS
 - fullsearch
- **Classical**
 - Forward, backward selection
 - Best subset selection
 - Significance tests
 - Significance based on Jack-knife
 - GOLPE
- **Other**
 - Pure variables
 - Principal variables
 - Iterative weighting with regression vector
 - ...

126

(see the Variable Selection Course at EigenU)



Variable Selection Methods

- How to choose which method?!?
- Different methods work in different situations
- Interval-PLS is a good “example” method to understand the considerations of variable selection. Simple to implement and use.
- Can be used on the Olive Oil data set, but first need to define PLS-DA

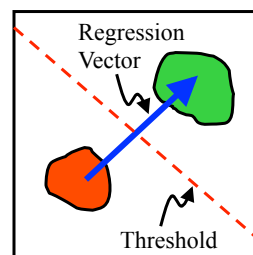
127



Partial Least Squares-Discriminate Analysis (PLS-DA)

- Use logicals (0,1) in Y-block to indicate if sample belongs to a class or not.
- Develop PLS model to predict class block
- Thresholds must be set between 0 and 1 to indicate if new samples are a member of each class...

Can use Bayes theorem to set threshold and include prior probability of each class

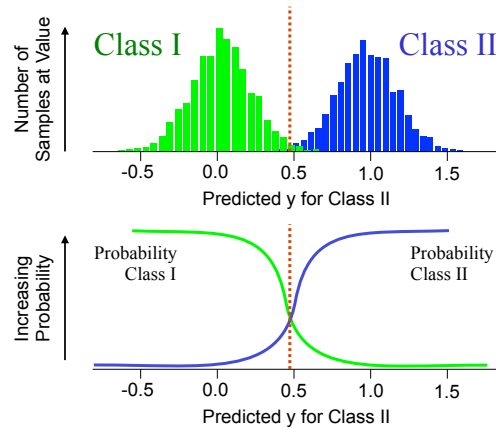


128



Thresholds in PLS-DA

Observed distribution of predictions can be handled in a straight-forward Bayesian way



129



PLS-DA for Olive Oil Data

- PLS-DA tends to capture variance which is useful in separating classes and ignoring variance within a class.
 - goal: maximize inter-class variance while minimizing intra-class variance
- For Olive Oils it seems reasonable to discriminate Corn Margarine from all the others first.
 - Other classes can be separated in turn
 - Two classes: Corn margarine and Everything else
 - this was evident based on the previous exploratory analysis

130



Analysis - PLSDA 1 LVs - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

Load Data Import Data **Edit Data** Plot Data Clear Data Save Data

Model Calibration Validation

Olive Oil Calibration - DataSet Editor

File Edit View FigBrowser

Info Data Row Labels **Column Labels**

Title: Wavenumber

Label	Axis Scale	Class	Incl.
Set 1 (empty)	Set 1: 0		
Wave	Wave		
Col 1	3600		
Col 2	3594		
Col 3	3588		
Col 4	3583		
Col 5	3577		
Col 6	3571		

Copy
Paste
Edit Incl.
Clear/Reset
Bulk Include Change
Use as Class
Sort By Selected (Ascend)
Sort By Selected (Descend)
Load Incl.
Extract Incl.

Re-Include all variables
Click "X" >> Edit Data
Select Column Labels tab
Right Click Incl. >> Bulk Include
Change
"Select all"
"OK"

Included Columns [axis scale : class]

1 [3600:]
2 [3594:]
3 [3588:]
4 [3583:]
5 [3577:]
6 [3571:]
7 [3565:]
8 [3559:]
9 [3554:]
10 [3548:]
11 [3542:]
12 [3536:]
13 [3530:]
14 [3525:]
15 [3519:]
16 [3513:]
17 [3507:]
18 [3501:]
19 [3496:]
20 [3490:]
21 [3484:]

Select all
OK Cancel

131

EIGENVECTOR
RESEARCH INCORPORATED

Analysis - PLSDA (No Model) - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

DECOMPOSITION
PCA
Purity
MCR
MPCA

CLUSTERING
Cluster

REGRESSION
PLS
PCR
LWR
SVM (SVM-R)
MLR
CLS

CLASSIFICATION
PLSDA
SVM (SVM-C)
SIMCA
KNN

MULTI-WAY
PARAFAC
Multi-way PLS (NPLS)

Analysis Methods Help
Simplify Menu

Analysis - PCA (No Model) - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

Model Calibration Validation

Preprocessing X-block

Available Methods

- Transformations ---
 - Absolute Value
 - Log10
 - Transmission to Absorbance (log(1/T))
- Filtering ---
 - Baseline (Specified points)
 - Baseline (Weighted Least Squares)
 - Derivative (SavGol)
 - Detrend
 - EMSC (Extended Scatter Correction)
 - EPO Filter
 - OLS Weighting
 - Kaiser HoldReact Method
 - OSC (Orthogonal Signal Correction)
 - Smoothing (SavGol)
- Normalization ---
 - MSC (mean)
 - Normalize
 - SNV

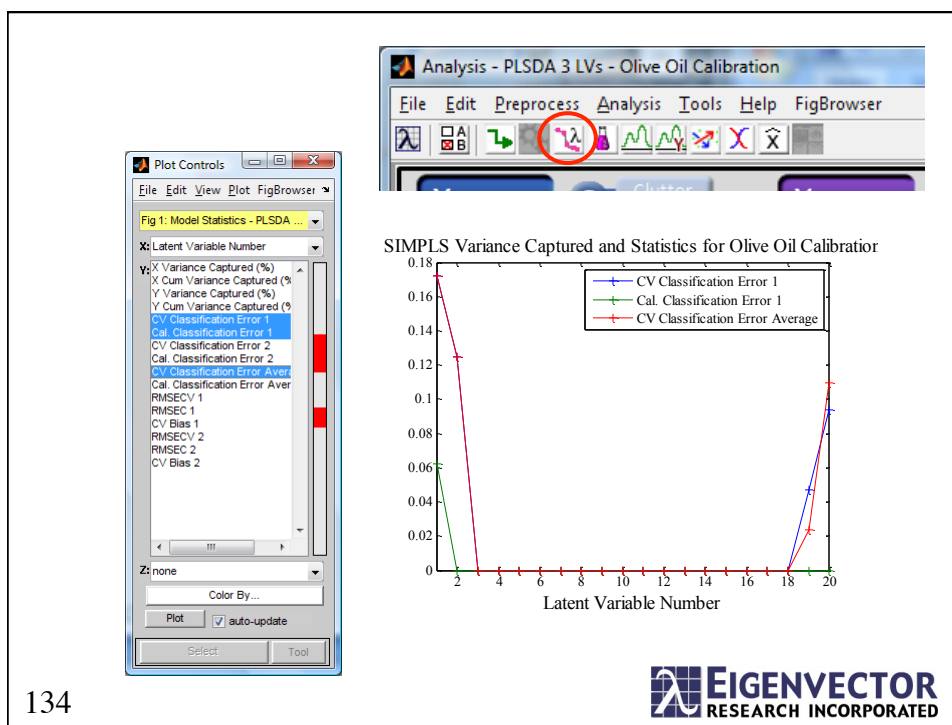
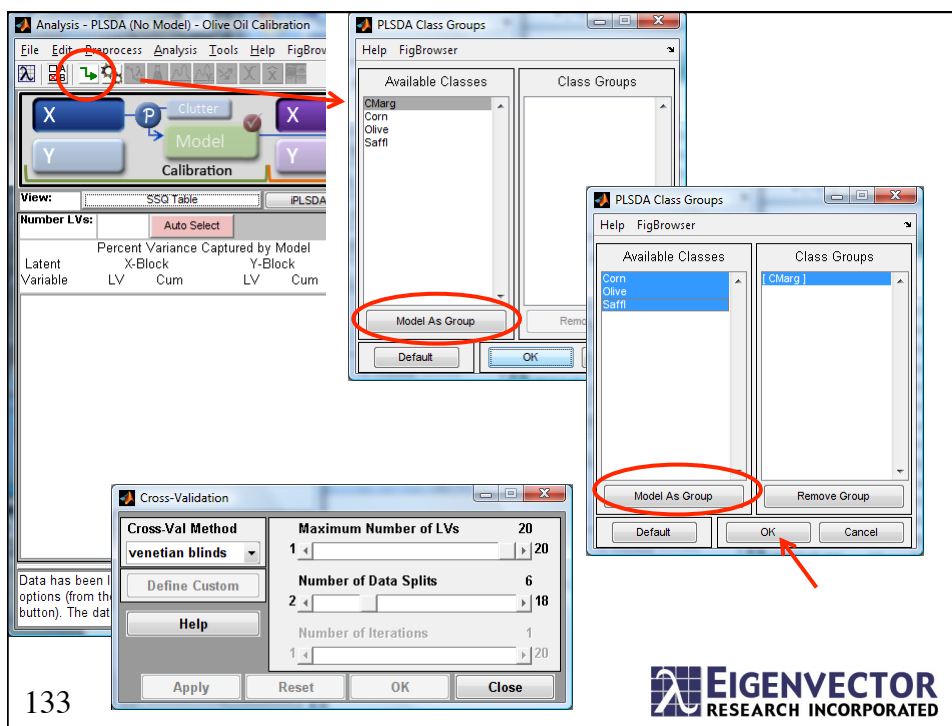
Selected Methods

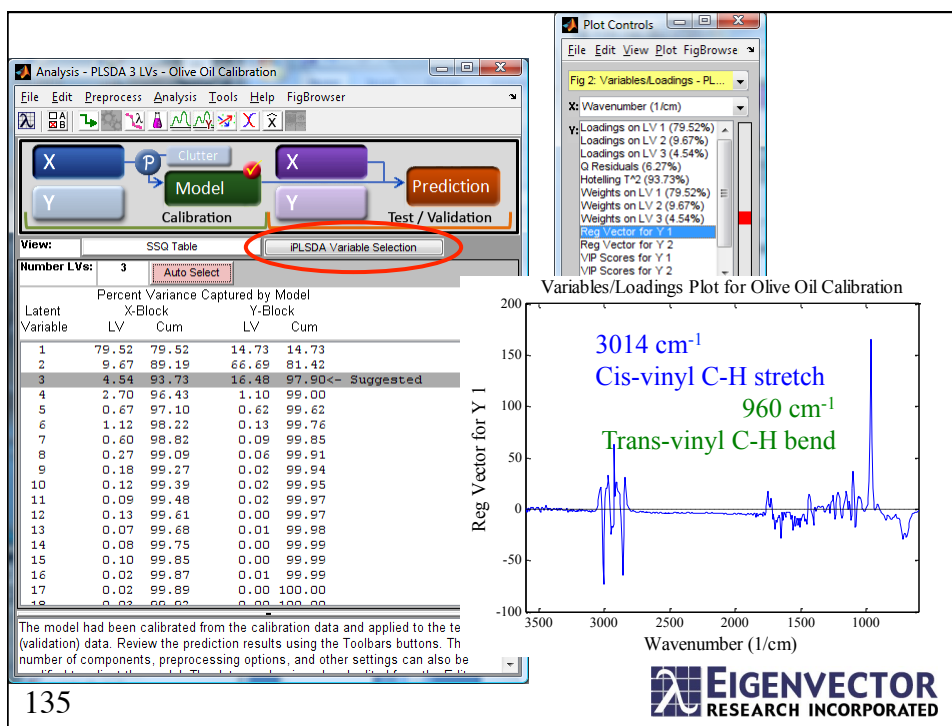
- Normalize (1-Norm, Area = 1)
- Mean Center

OK Cancel

132

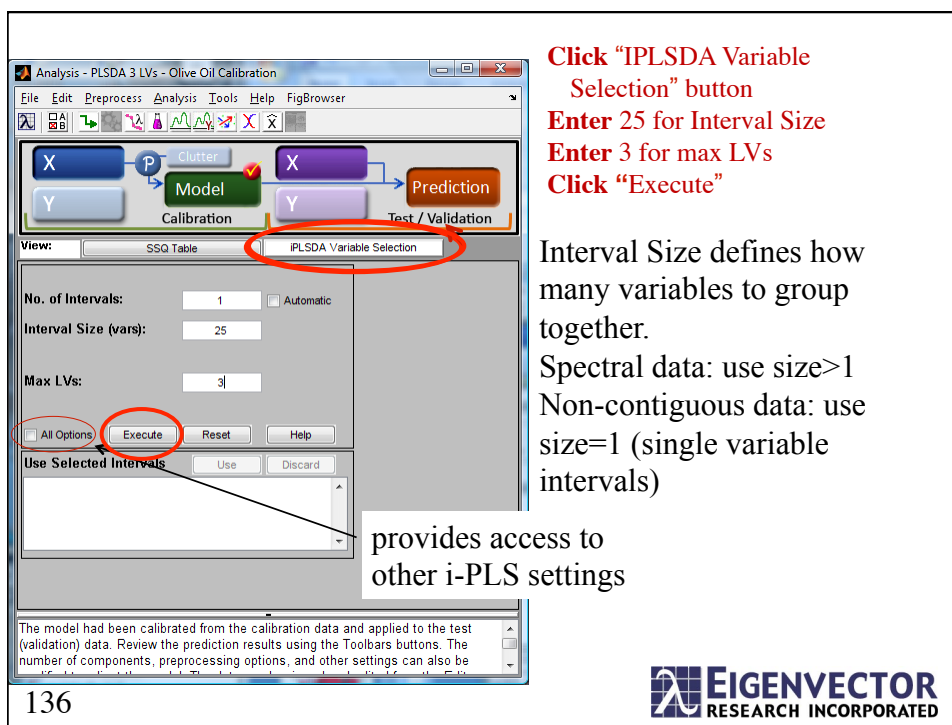
EIGENVECTOR
RESEARCH INCORPORATED





135

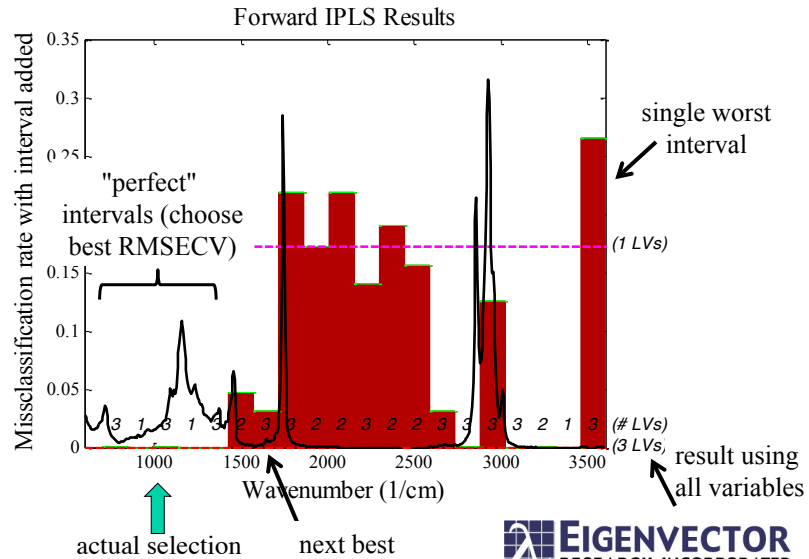
EIGENVECTOR
RESEARCH INCORPORATED



136

EIGENVECTOR
RESEARCH INCORPORATED

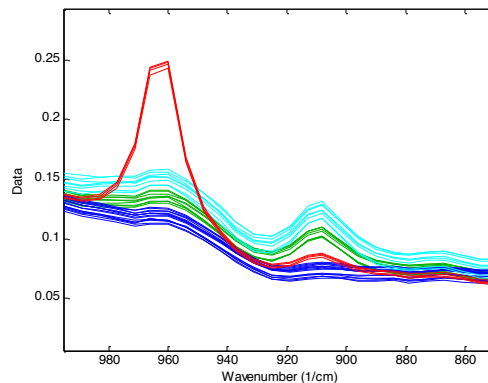
First Interval Result



137

Why the Low-Frequency Intervals? Information Content

- Best intervals contain useful signal from all four classes.
- Original data shown; Preprocessing only makes this separation better!

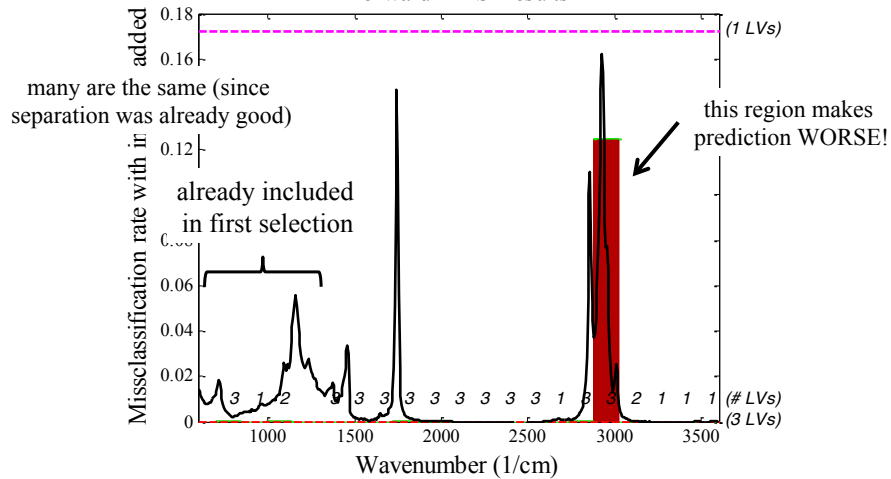


138

EIGENVECTOR RESEARCH INCORPORATED

Repeat "Execute" and "Add To Previous"

Forward IPLS Results



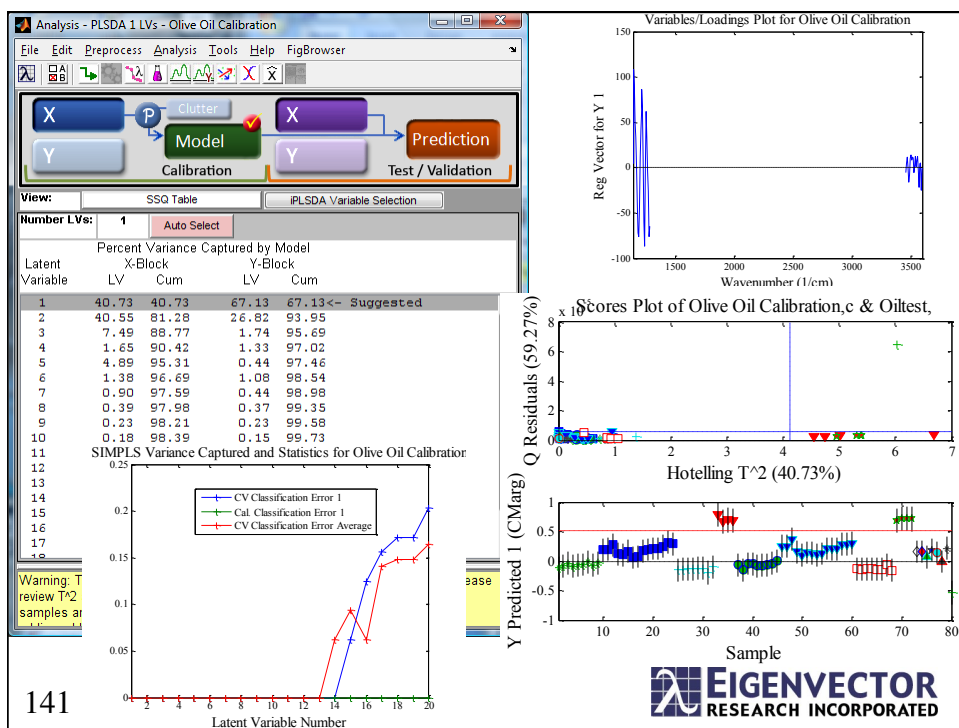
139

EIGENVECTOR
RESEARCH INCORPORATED

Click "Use" button
Click "OK"
Click "SSQ Table" button
Click "Model"

EIGENVECTOR
RESEARCH INCORPORATED

140

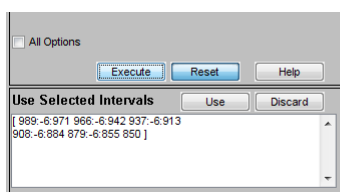


Number of Intervals

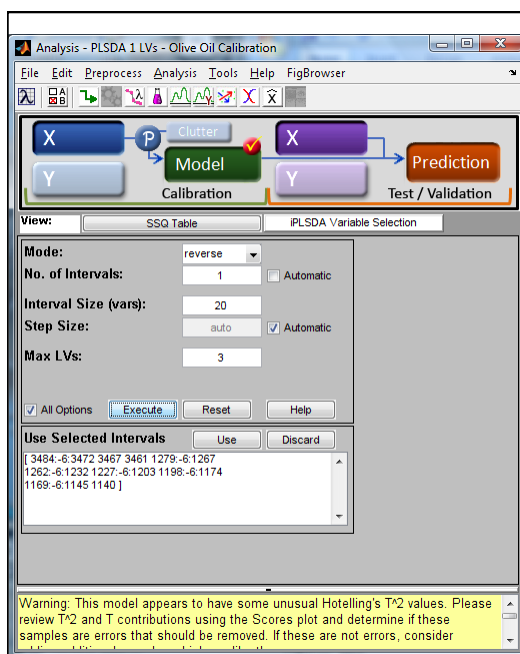
- Can choose a pre-set number of intervals to find
- Can also use “Automatic” to continue selecting intervals until RMSECV/misclassification does not improve
- This is **not** the same as exhaustive combinatorial search (fullsearch). It is sequential (choose one, “lock” it in, choose a second, “lock” it in...)
- For very complex data, may not give actual “best” windows, but probably not a bad one.

What is the Result?

- For PLSDA, lower RMSECV should indicate better class separation in predicted Y values
- Selecting additional intervals gives little improvement in RMSECV (on this data)
- Use ONLY first selected interval and build new model...



143



144



Outline

- Introduction
- PCA Review
- PLS Regression Review
- Advanced Preprocessing
- Variable Selection
- Summary

145



Summary

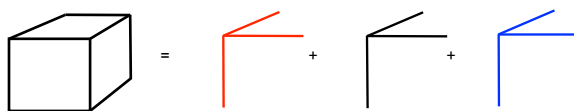
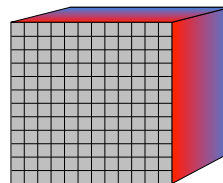
- Data analysis requires knowledge of
 - the system, physics, chemistry and math → ~~black box~~
- Advanced Preprocessing
 - uses knowledge of the clutter (GLS, ELS, etc.)
- Variable Selection
 - choose variables that are most predictive

146



If time ...

- Introduce
 - multivariate image analysis (hyperspectral image analysis)
 - multiway analysis



147



Outline

- Introduction
- PCA Review
- PLS Regression Review
- Advanced Preprocessing
- Variable Selection
- Introduction to Multivariate Image Analysis and Multi-way Analysis (if time)
- Summary

148



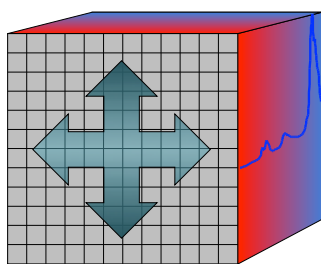
Multivariate Images

A data array of *dimension three* (or more) where the first two dimensions are *spatial* and the last dimension(s) is a function of another variable.

149



Multivariate Images



Spatial Information
between pixels

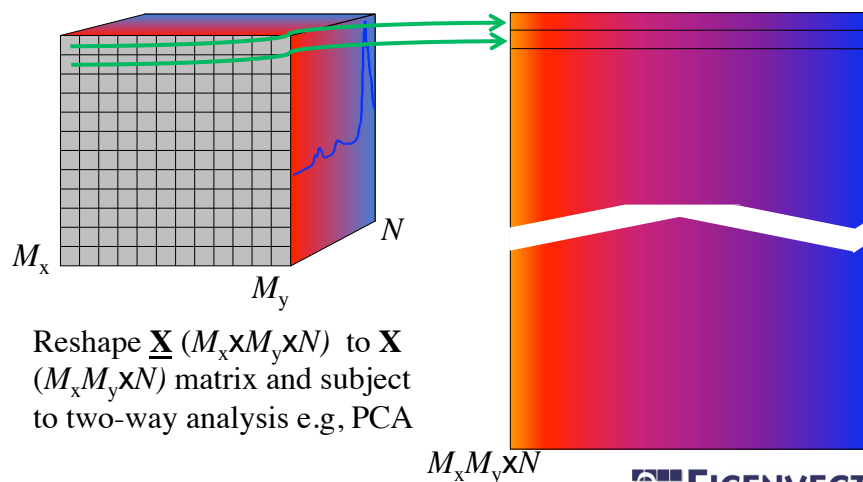
Spectral Information
between channels
(chemical information)

**Spatial distribution of
chemical analytes, physical
features, and other
properties**

150



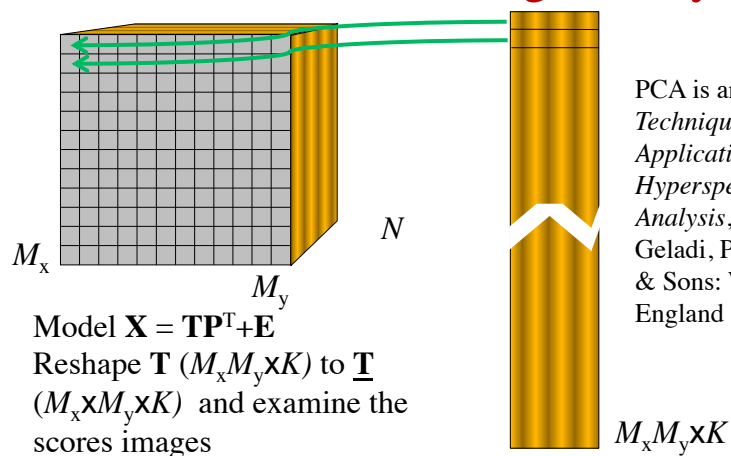
Reshaping Images for Analysis



151

 **EIGENVECTOR**
RESEARCH INCORPORATED

PCA for Multivariate Image Analysis



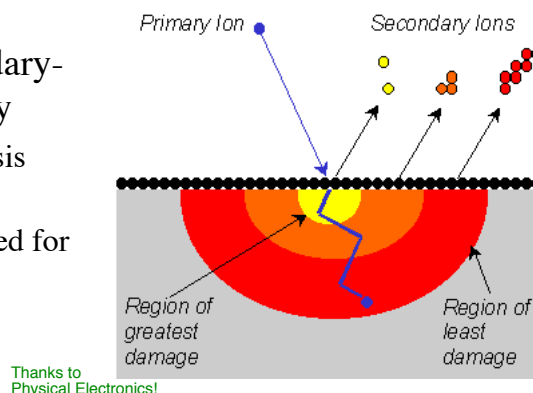
PCA is an example of MIA.
Techniques and Applications of Hyperspectral Image Analysis, Grahn, H. F.; Geladi, P., Eds. John Wiley & Sons: West Sussex, England (2007)

152

 **EIGENVECTOR**
RESEARCH INCORPORATED

Example: TOF-SIMS

- Time-of-Flight Secondary-Ion-Mass Spectrometry
 - common surface analysis technique
 - mass spectrum generated for each pixel



153



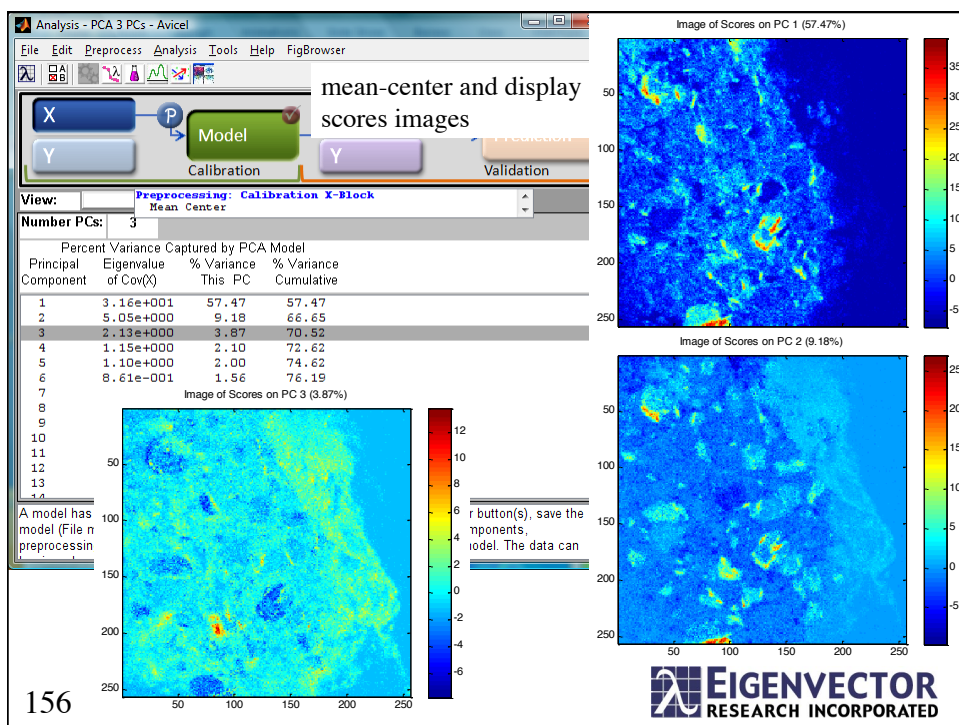
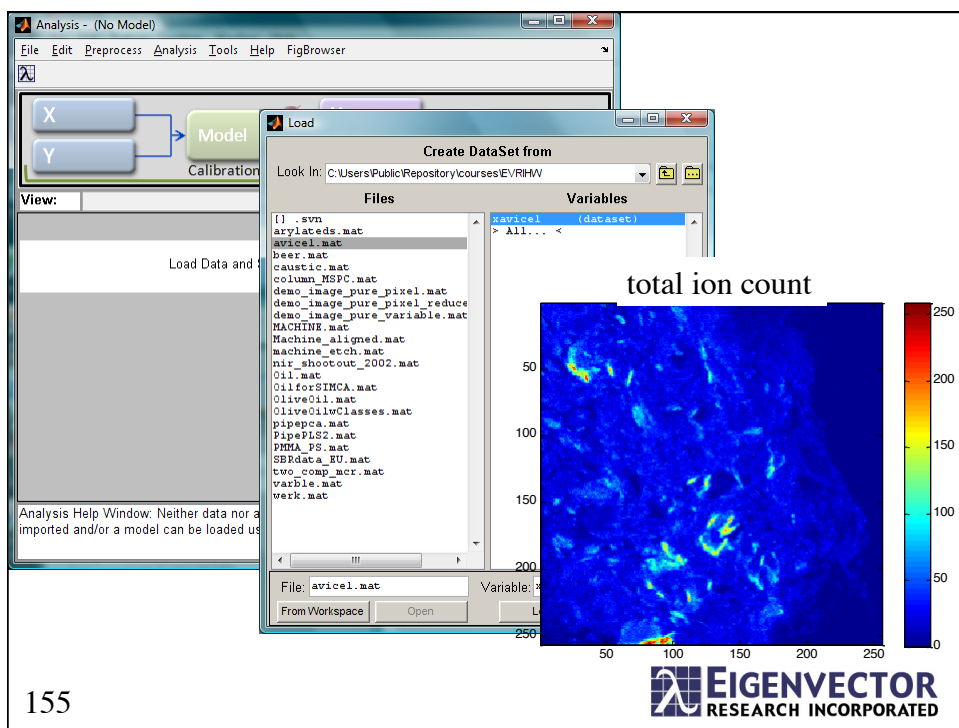
TOF-SIMS of Time Release Drug Delivery System

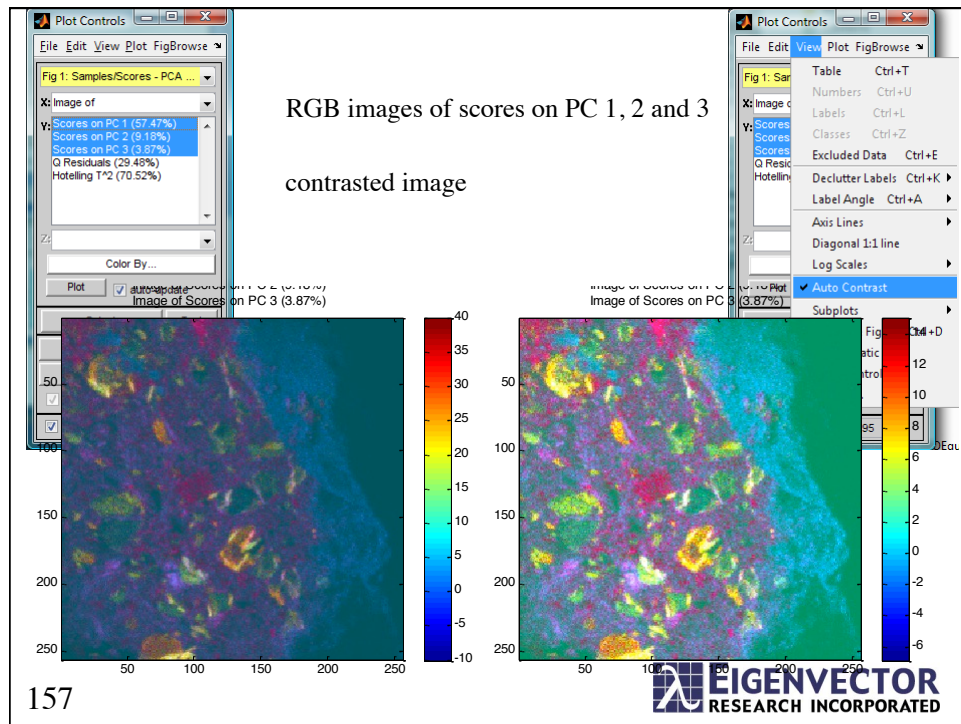
- Multi-layer drug beads serve as controlled release system
- TOF-SIMS of cross section of bead
- Evaluate the integrity of the layers and distribution of ingredients

A.M. Belu, M.C. Davies, J.M. Newton and N. Patel, "TOF-SIMS Characterization and Imaging of Controlled-Release Drug Delivery Systems," *Anal. Chem.*, **72**(22), 5625–5638 (2000).
 Gallagher, N.B., Shaver, J.M., Martin, E.B., Morris, J., Wise, B.M. and Windig, W., "Curve resolution for images with applications to TOF-SIMS and Raman", *Chemometr. Intell. Lab.*, **73**(1), 105–117 (2003).

154







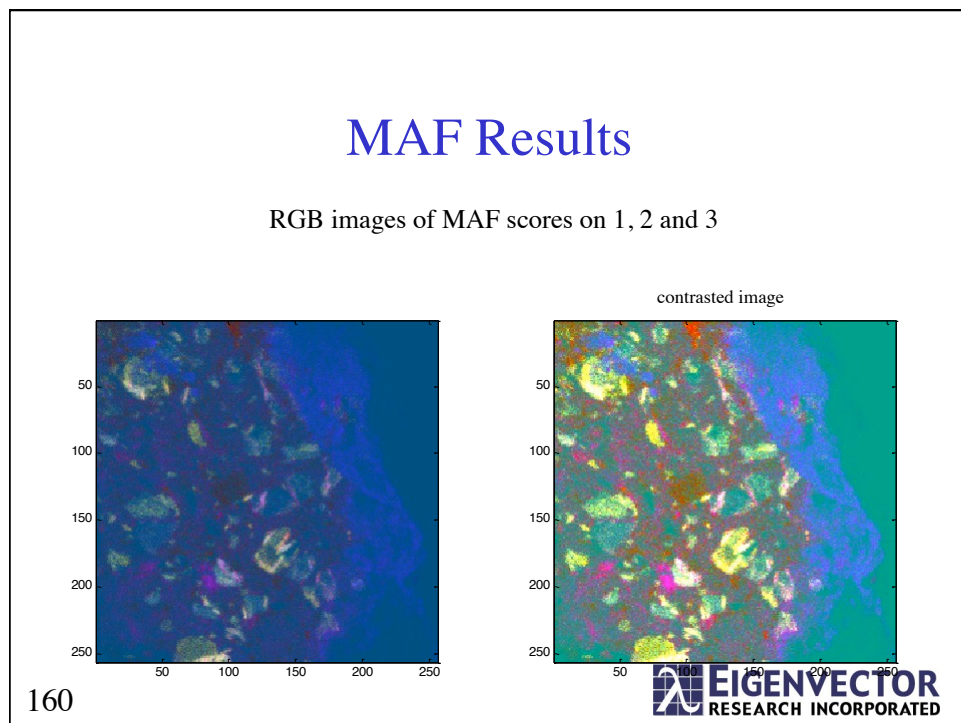
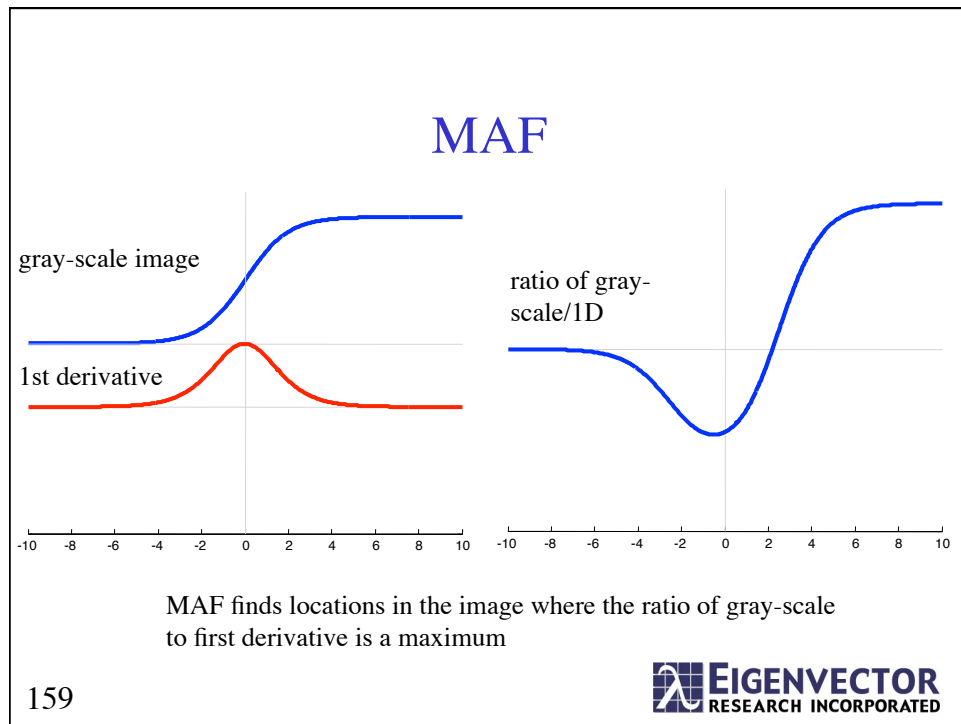
Minimum Noise Factors (MNF)

- MNF attempts find directions in the data that maximize the signal-to-clutter.

$$\max_{\mathbf{v}_i \neq 0} \left(\frac{\mathbf{v}_i^T \Sigma_X \mathbf{v}_i}{\mathbf{v}_i^T \Sigma_C \mathbf{v}_i} \right) \quad \text{the objective function}$$

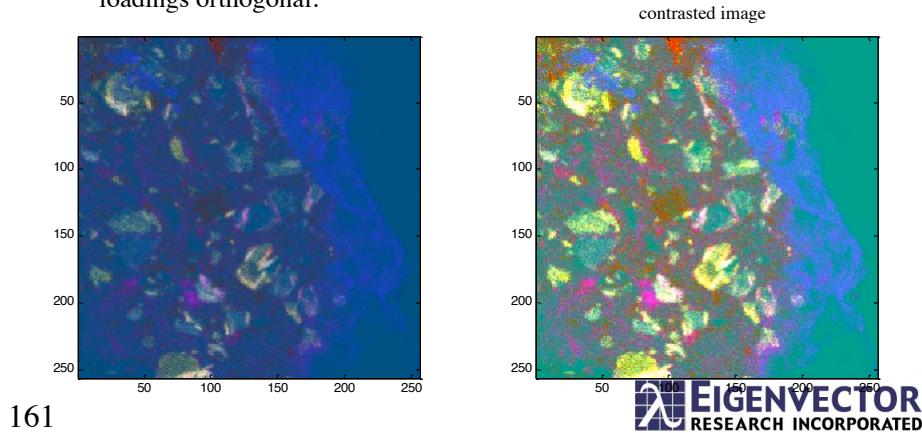
- Result is a PCA-like eigenvector problem
- In **maximum autocorrelation factors (MAF)** clutter is the first difference image (difference between near-by pixels)

158



PCA w/ GLS Weighting for ~MAF

RGB images of PCA w/ GLS weighting scores on 1, 2 and 3. Similar to MAF results. Objective function ~similar, but PCA scores and loadings orthogonal.

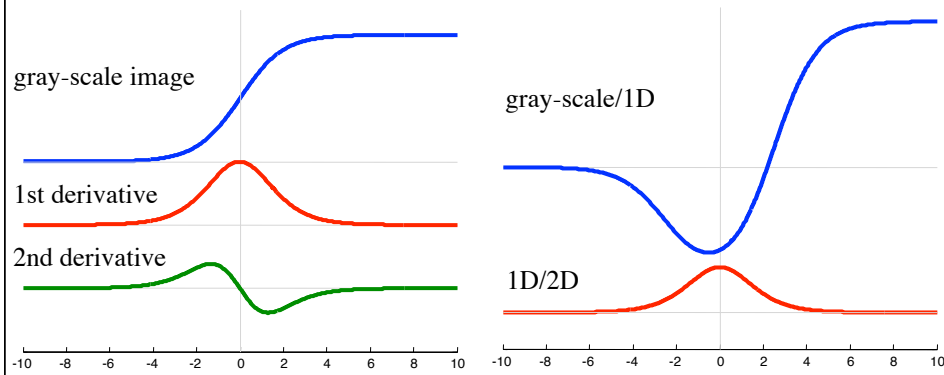


Maximum Difference Factors (MDF)

- In MDF the signal covariance corresponds to the first derivative across the spatial dimensions.
 - in MAF the first difference is the clutter
- The clutter corresponds to the second derivative across the spatial dimensions.
- Gives a multivariate analysis estimate of edges in an image.
 - analogous method available for GLS weighting w/ PCA

162

MDF



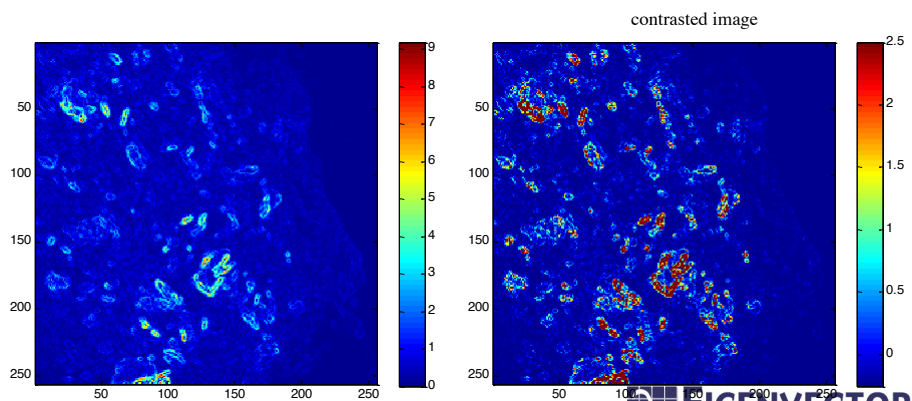
MDF finds locations in the image where the ratio of first to second derivative is a maximum

163



MDF Results

Scaled image of MDF scores on $\frac{1}{\sqrt{dx^2 + dy^2}}$

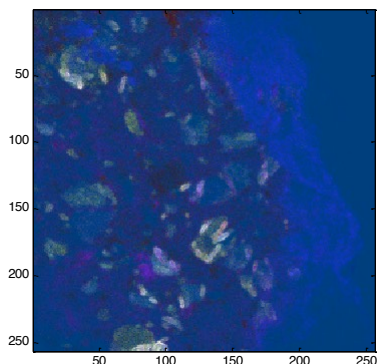


164

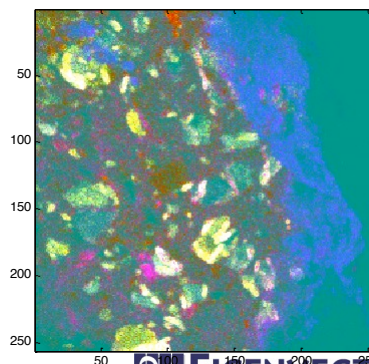


MAF+MDF Results

RGB image of MAF scores 1, 2 and 3 + MDF scores on 1
 $\sqrt{dx^2 + dy^2}$



contrasted image



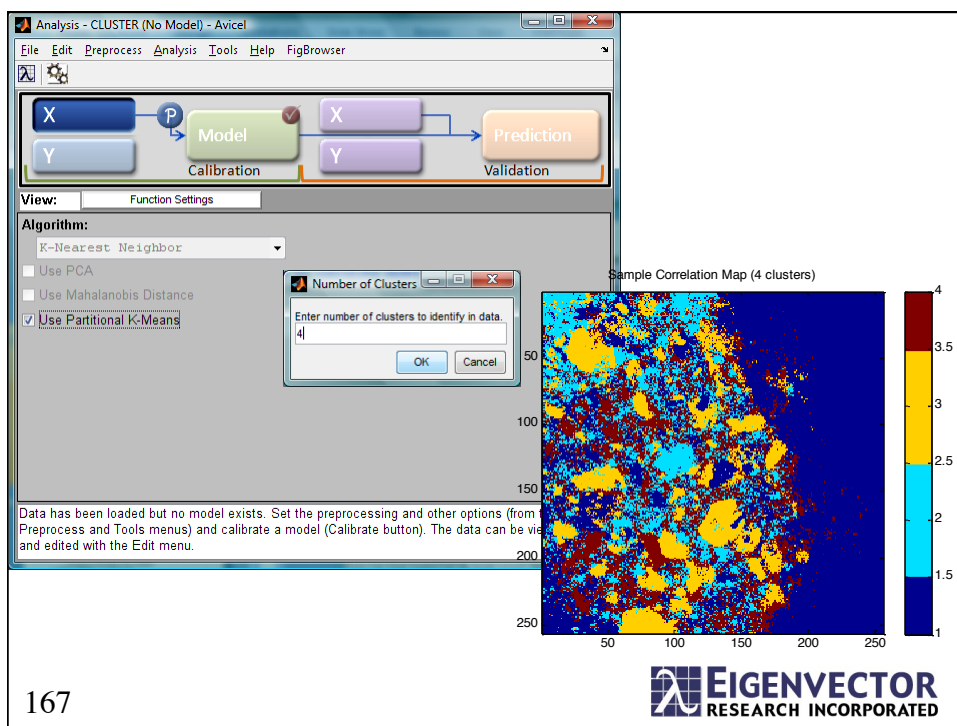
165



cluster analysis

166





MIA ...

- Much more to MIA
 - linked scores plots and density plots
 - interactive exploration of the image(s)
 - image SIMCA and PLS-DA
 - classification
 - curve resolution
 - chemical identification and mapping
 - image statistical process control (ISPC) for multivariate statistical process control (MSPC)
 - ...

168

Outline

- Introduction
- Advanced Preprocessing
 - Clutter and characterizing clutter
 - Generalized least squares weighting
 - Extended multiplicative scatter correction
 - Interval PLS (iPLS)
 - Model Robustness
- Multivariate image analysis
- Multi-way Analysis
- Summary

©Copyright 2008-2012
Eigenvector Research, Inc.
No part of this material may be photocopied or
reproduced in any form without prior written
consent from Eigenvector Research, Inc.



Why is Clutter Bad?

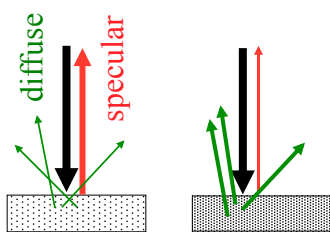
- What is clutter and how does clutter effect the measured signal?
- Use FT-IR spectra and pattern recognition to distinguish authentic olive oil from counterfeit or adulterated olive oil.



Sources of Clutter: Scattering Effects in Reflectance

Caused by variations in:

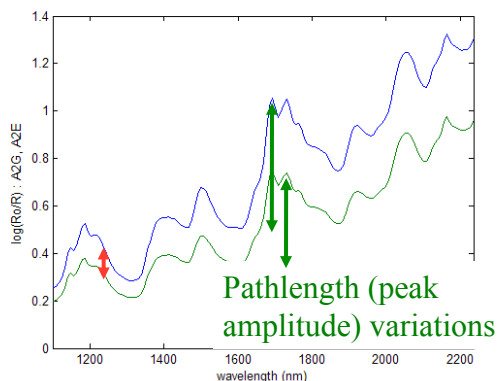
- Particle size (mean & distribution)
- Sample opacity
- Sampling packing density
- Sample placement



171

Sample 1

Sample 2



Baseline offset
changes

EIGENVECTOR
RESEARCH INCORPORATED

Olive Oil Samples Learning / Calibration Set:

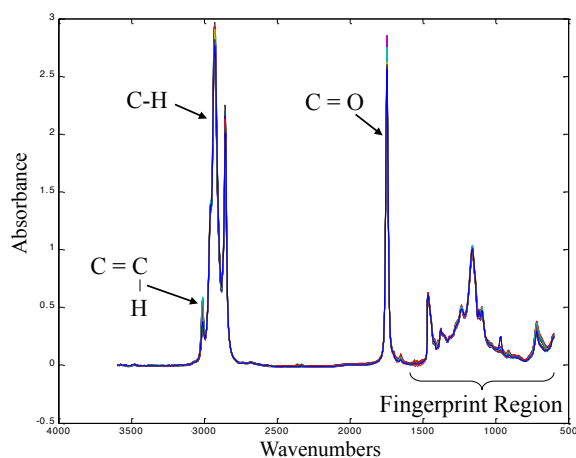
Corn Oil	9 samples	(#1-9)
Olive Oil	15 samples	(#10-24)
Safflower Oil	8 samples	(#25-32)
Corn Margarine	4 samples	(#33-36)

Took FT-IR spectra ($3600 - 600 \text{ cm}^{-1}$) of these oils using a fixed pathlength NaCl cell.

172

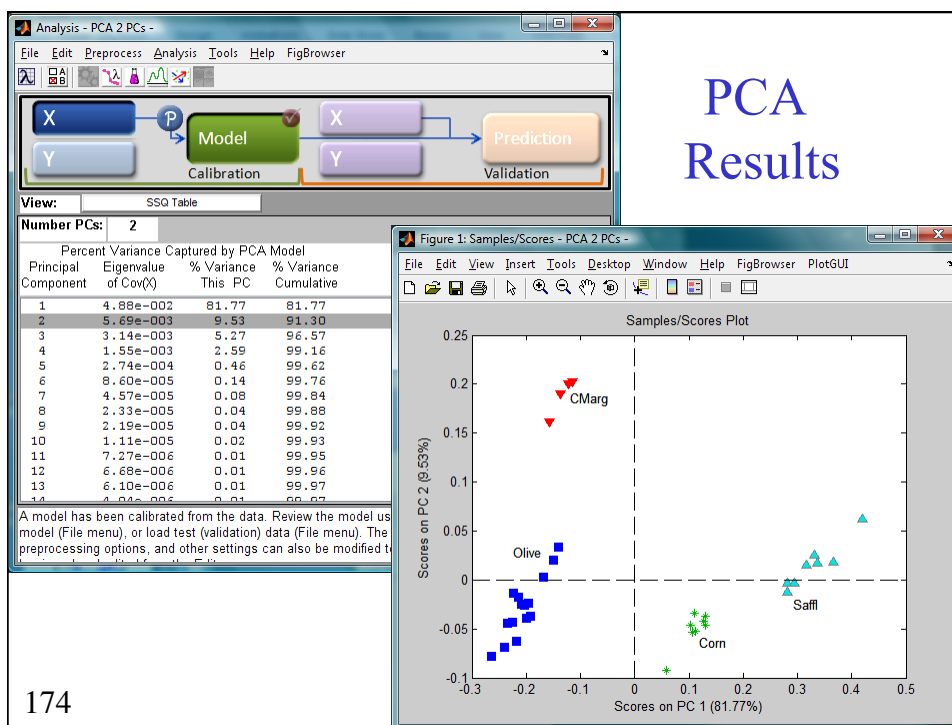
EIGENVECTOR
RESEARCH INCORPORATED

FTIR Spectra of 36 Sample Learning Set



173

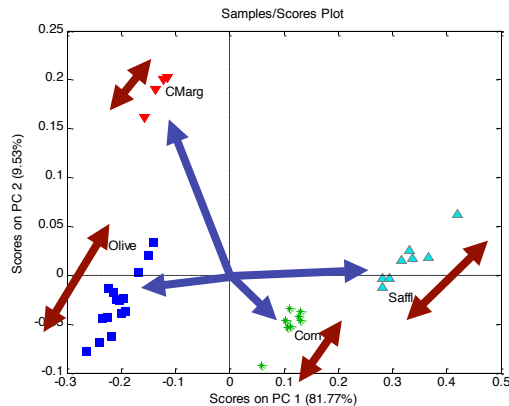
EIGENVECTOR
RESEARCH INCORPORATED



174

PCA Results

- PCA shows that the four classes in the calibration data set are separate from each other (high **between class variance**) but ...
- have significant **within class variance**

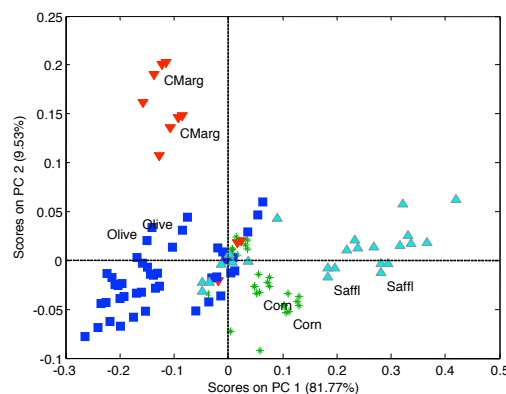


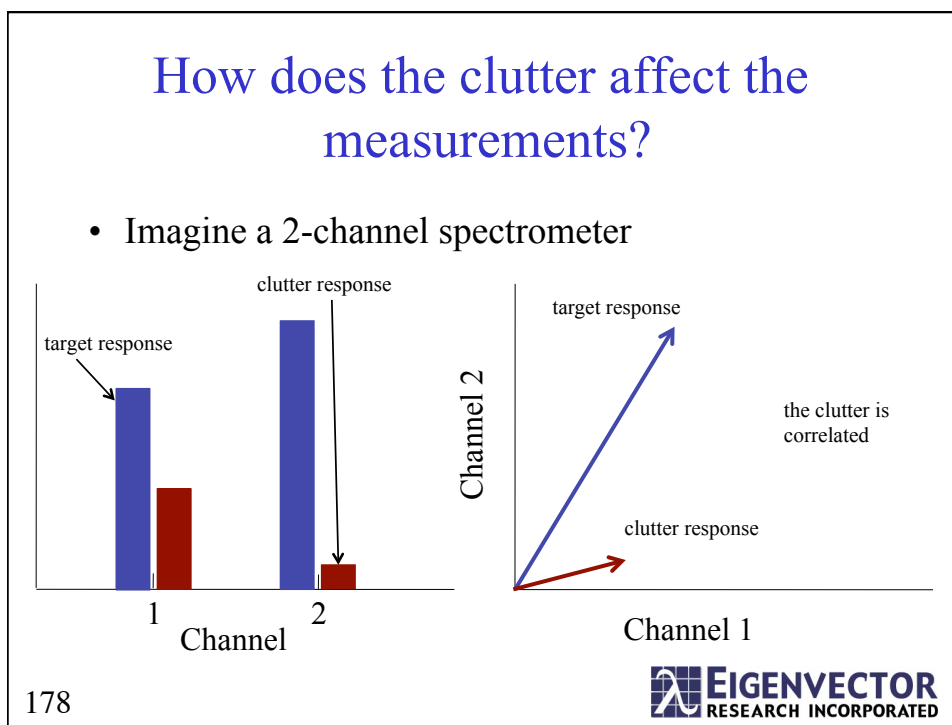
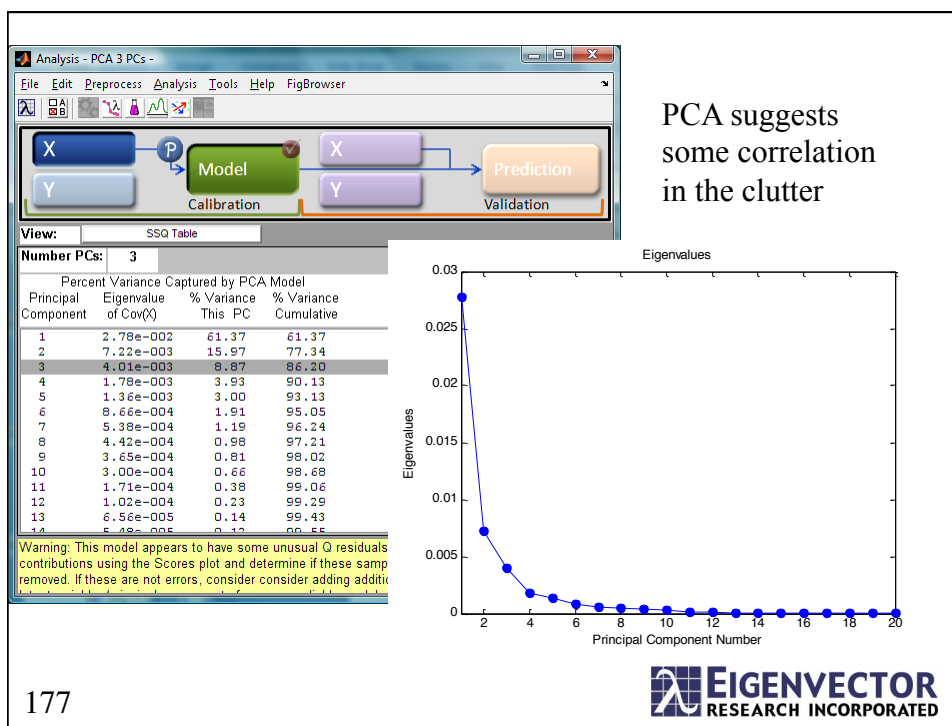
EIGENVECTOR
RESEARCH INCORPORATED

175

Replicates

- Ideally, replicates would lie on top of each other.
- Variance within each class is clutter variance.
 - Is it random noise? Is the clutter correlated?
- Center each class to it's own mean and do PCA on the result.

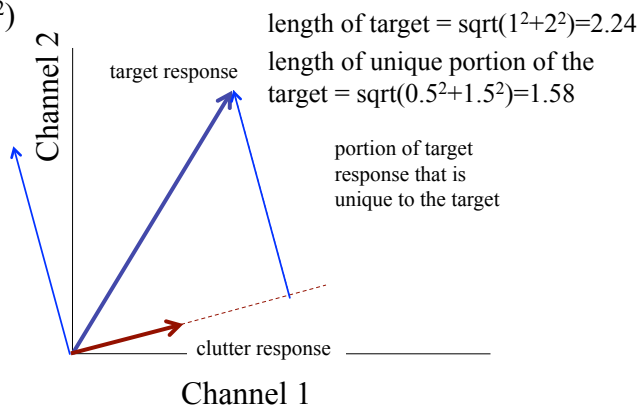




How does the clutter affect the measurements?

- characterize the signal as the length of the vector

$$\sqrt{x_1^2 + x_2^2}$$



179



Why is clutter bad?

- The signal-to-clutter is ~proportional to the length of the unique portion of the target's response.
 - in absence of clutter it was 2.24
 - in the presence of clutter it was 1.58
- In regression, clutter-to-signal is related to the estimation error.
 - higher clutter-to-signal → higher estimation error
 - in the presence of clutter the estimation error is 2.24/1.58 times the error when clutter is absent

180



Effect of Clutter

- The effect of clutter is to remove target signal
 - for olive oil example the target signal is the differences between the classes
- Instrument related clutter can be minimized by
 - good instrument design that accounts for the environment (noise+interferences) in which measurements are to be made
 - instrument standardization
 - remove drifts in offsets and gains that adds to the clutter
- Can't always be eliminated → what to do?

181



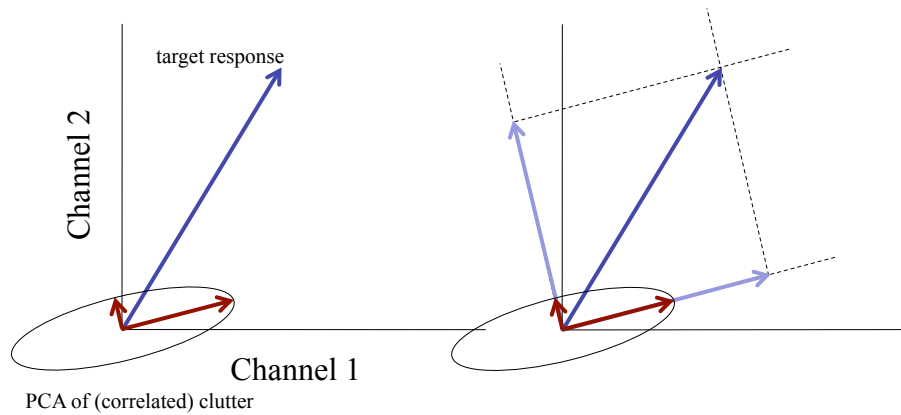
Accounting for Clutter

- One method used to account for clutter is a weighting scheme
 - similar to that used in **generalized least squares (GLS)**
- Autoscaling scales each variable to unit variance
- GLS weighting scales each clutter direction (as determined using PCA) to unit variance
 - directions of high clutter are deweighted
 - directions of low clutter are given more opportunity to allow signal through

182



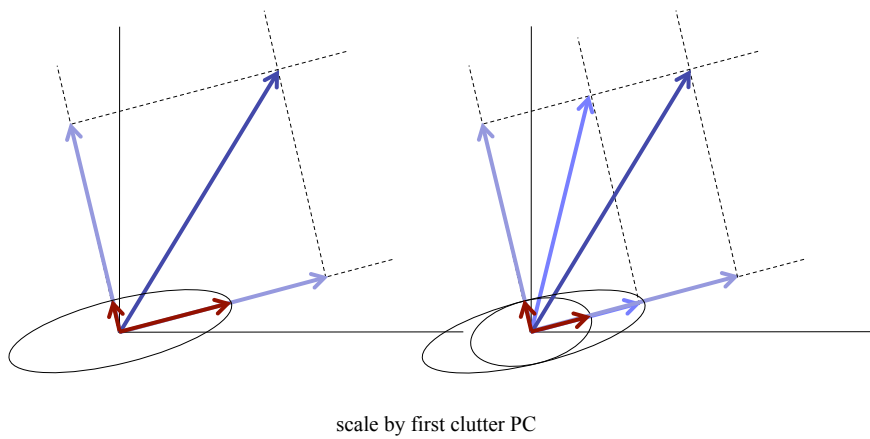
Target Projected onto Clutter



183

 **EIGENVECTOR**
RESEARCH INCORPORATED

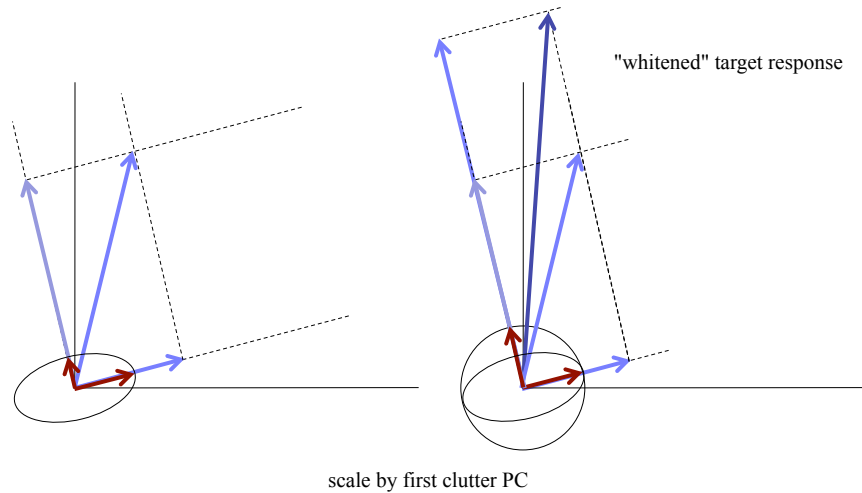
Scale Target by Clutter



184

 **EIGENVECTOR**
RESEARCH INCORPORATED

Scale Target by Clutter

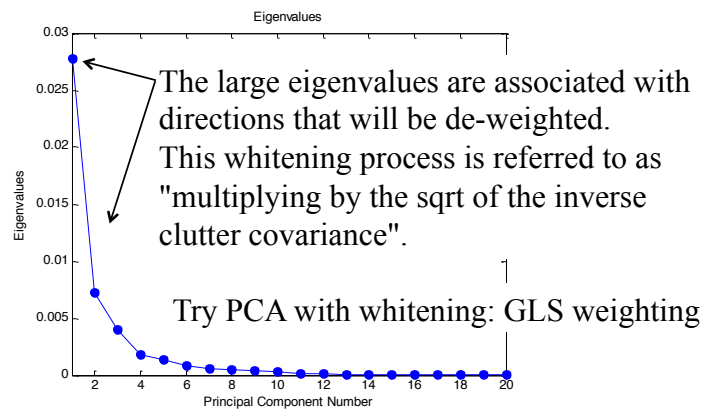


185



Olive Oil Clutter

Eigenvalue distribution of the within class variance.



186



Click Load X Calibration Shortcut
Workspace/Mat file
OliveOilwClasses.mat
xcal
Click Analysis:PCA

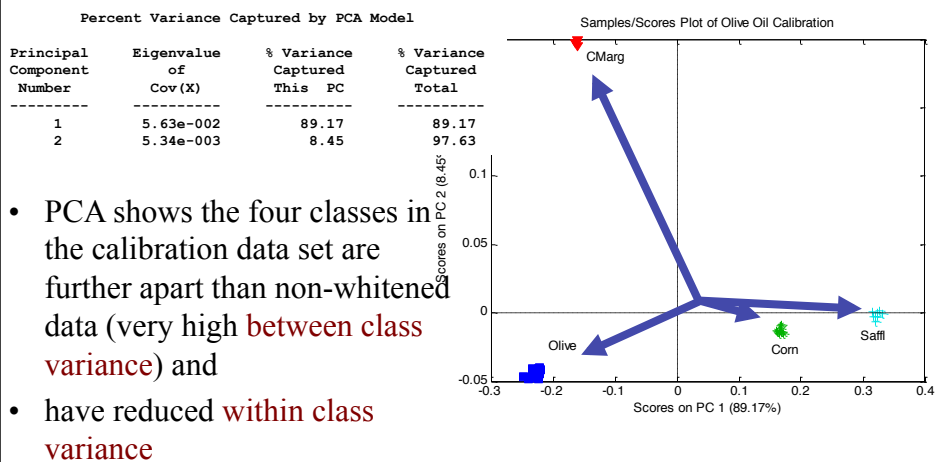
187

Click Preprocessing Shortcut
Autoscale: <-- Remove
GLS Weighting: Add -->
x-block classes
Mean Center: Add-->
OK

188

EIGENVECTOR
RESEARCH INCORPORATED

PCA of Whitenized Spectra



189



Test Set:

Corn Oil*	9 samples	(#1-9)
Olive Oil*	15 samples	(#10-24)
Safflower Oil*	8 samples	(#25-32)
Corn Margarine*	4 samples	(#33-36)
Corn Oil in Olive Oil 5, 10, 20, 30 & 40%	5 samples	(#37-41)
Almond Oil	1 sample	(#42)
Peanut Oil	1 sample	(#43)
Sesame Oil	1 sample	(#44)

* New Samples not included in the calibration

190



Analysis - PCA 2 PCs - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

Click Load X Validation Shortcut
Workspace/Mat file
OliveOilwClasses.mat
xtest

View: SSO Table

Number PCs: 2

Percent Variance Captured by PCA Model

Principal Component	Eigenvalue	% Variance of Cov(X)	% Variance This PC	% Variance Cumulative
1	5.63e-002			
2	5.34e-003			
3	3.11e-004			
4	1.44e-004			
5	1.08e-004			
6	1.03e-004			
7	9.48e-005			
8	9.20e-005			
9	8.47e-005			
10	7.86e-005			
11	6.90e-005			
12	6.29e-005			
13	6.23e-005			
14	4.00e-005			

Warning: This model appears to have been trained using the Scores removed. If these are not errors, d

Import

Import from file type:

- Workspace/MAT file
- Delimited Text File (CSV, TXT)
- XY... Delimited Text Files (TXT, XY)
- Excel File (XLS, CSV, TXT)
- Hamilton Sundstrand ASF File (ASF, AIF, BKH)
- Thermo Galactic File (SPC)
- JCAMP [general] (DX, JDX)
- Extensible Markup Language (XML)
- Paste XML from Clipboard
- AdventacT MTF File (MTF)
- Camcra Ion-ToF BIF Image (BIF)
- Physical Electronics RAW Image (RAW)
- Image (Workspace/MAT file)
- Image (Other...)
- Other...

OK Cancel

Create DataSet from

Look In: C:\Users\Public\Repository\courses\EVRI\HW

Files

- 11...svn
- arylateds.mat
- beer.mat
- caustic.mat
- column_MSPC.mat
- demo_image_pure_pixel.mat
- demo_image_pure_pixel_reduce
- demo_image_pure_variable.mat
- MACHINE.mat
- Machine_aligned.mat
- machine_etch.mat
- nir_shootout_2002.mat
- Oil.mat
- OilforSIMCA.mat
- OliveOil.mat
- OliveOilwClasses.mat
- pipepca.mat
- PipePls2.mat
- PHMA_PS.mat
- SBBdata_EU.mat
- two_comp_mcr.mat
- varile.mat
- verk.mat

Variables

- xcal (dataset)
- xtest (dataset)
- > All... <

Prediction Validation

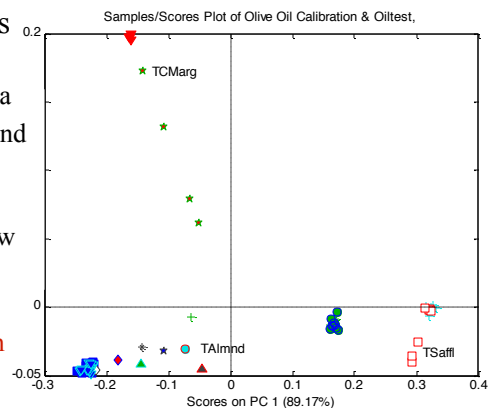
File: OliveOilwClasses.mat Variable: xtest

From Workspace Open Load Cancel

191

Check Test Data

- test samples are filled symbols
 - Olive Oil and Corn appear to cluster on top of calibration data
 - This is less true for Safflower and especially Margarine, why?
 - The Margarine samples might suggest a magnitude effect. How is the effect of magnitude removed?
 - multiplicative scatter correction
 - standard normal variate
 - spectral normalization



Analysis - PCA (No Model) - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

X Y Model Calibration Validation Prediction

View: SSO Table

Number PCs: Percent V Principal Eig Component of

Data has been loaded, Preprocess and T and edited with th

Click Preprocessing Shortcut
GLS Weighting (RHS)
Settings...
change alpha: 0.001

Preprocessing X-block

Available Methods

- Transformations ---
 - Absolute Value
 - Log10
- Filtering ---
 - Baseline (Specified points)
 - Baseline (Weighted Least Squares)
 - Derivative (SavGol)
 - Detrend
 - EMSC (Extended Scatter Correction)
 - EPO Filter
 - GLS Weighting
 - OSC (Orthogonal Signal Correction)
 - Smoothing (SavGol)
- Normalization ---
 - MSC (mean)
 - Normalize
 - SNV
- Scaling and Centering ---
 - Autoscale
 - Group Scale
 - Log Decay Scaling
 - Mean Center
 - Median Center
 - Multiway Center
 - Multiway Scale

Selected Methods

- GLS Weighting (classes, alpha 0.02)
- Mean Center
- <end>

GLSW Settings

Clutter Source

- ☐ automatic
- ☐ y-block gradient
- ☒ x-block classes
- ☐ external data

Load Edit Size: <empty>

Algorithm

- ☒ GLSW
 - alpha: 0.001
- ☐ EPO
 - Number of PCs: 1

OK Cancel Help

193

EIGENVECTOR
RESEARCH INCORPORATED

Preprocessing X-block

Available Methods

- Transformations ---
 - Absolute Value
 - Log10
- Filtering ---
 - Baseline (Specified points)
 - Baseline (Weighted Least Squares)
 - Derivative (SavGol)
 - Detrend
 - EMSC (Extended Scatter Correction)
 - EPO Filter
 - GLS Weighting
 - OSC (Orthogonal Signal Correction)
 - Smoothing (SavGol)
- Normalization ---
 - MSC (mean)
 - Normalize
 - SNV
- Scaling and Centering ---
 - Autoscale
 - Group Scale
 - Log Decay Scaling
 - Mean Center
 - Median Center
 - Multiway Center
 - Multiway Scale

Selected Methods

- GLS Weighting (classes, alpha 0.02)
- Mean Center
- <end>

Norm Settings

Type: 2-Norm (length = 1)

Window: (all variables)

Select All Variables

OK Cancel

Preprocessing X-block

Available Methods

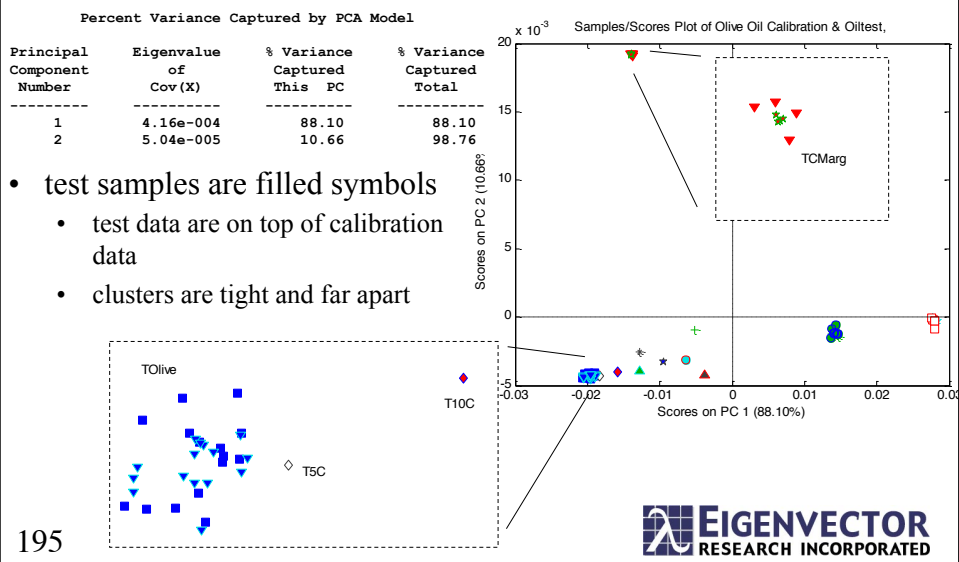
- Transformations ---
 - Absolute Value
 - Log10
- Filtering ---
 - Baseline (Specified points)
 - Baseline (Weighted Least Squares)
 - Derivative (SavGol)
 - Detrend
 - EMSC (Extended Scatter Correction)
 - EPO Filter
 - GLS Weighting
 - OSC (Orthogonal Signal Correction)
 - Smoothing (SavGol)
- Normalization ---
 - MSC (mean)
 - Normalize
 - SNV
- Scaling and Centering ---
 - Autoscale
 - Group Scale
 - Log Decay Scaling
 - Mean Center
 - Median Center
 - Multiway Center
 - Multiway Scale

Selected Methods

- Normalize (2-Norm, Length = 1)
- GLS Weighting (classes, alpha 0.001)
- Mean Center
- <end>

194

With Row Normalization



- test samples are filled symbols
 - test data are on top of calibration data
 - clusters are tight and far apart

Analysis - PCA (No Model) - Olive Oil Calibration

File Edit Preprocess Analysis Tools Help FigBrowser

X Y Model Calibration Prediction Validation

View: SSQ Table

Number PCs: Percent Variance of Principal Component of

Preprocessing X-block

Available Methods

- Transformations ---
- Absolute Value
- Log10
- Filtering ---
- Baseline (Specified points)
- Baseline (Weighted Least Squares)
- Derivative (SavGol)
- Detrend
- EMSC (Extended Scatter Correction)
- EPO Filter
- GLS Weighting
- OSC (Orthogonal Signal Correction)
- Smoothing (SavGol)
- Normalization ---
- MSC (mean)
- Normalize
- SNV
- Scaling and Centering ---
- Autoscale
- Group Scale
- Log Decay Scaling
- Mean Center
- Median Center
- Multiscale Center
- Multiscale Scale

Selected Methods

- Normalize (2-Norm, Length = 1)
- GLS Weighting (classes, alpha 0.001)
- Mean Center
- <end>

Savitsky-Gol...

Filter Width: 15

Polynomial Order: 2

Derivative Order: 2

OK Cancel

Preprocessing Y-block

Selected Methods

- 2nd Derivative (order: 2, window: 15 pt)
- Normalize (2-Norm, Length = 1)
- GLS Weighting (classes, alpha 0.001)
- Mean Center
- <end>

2nd Derivative (order: 2, window: 15 pt)

Normalize (2-Norm, Length = 1)

GLS Weighting (classes, alpha 0.001)

Mean Center

<end>

196

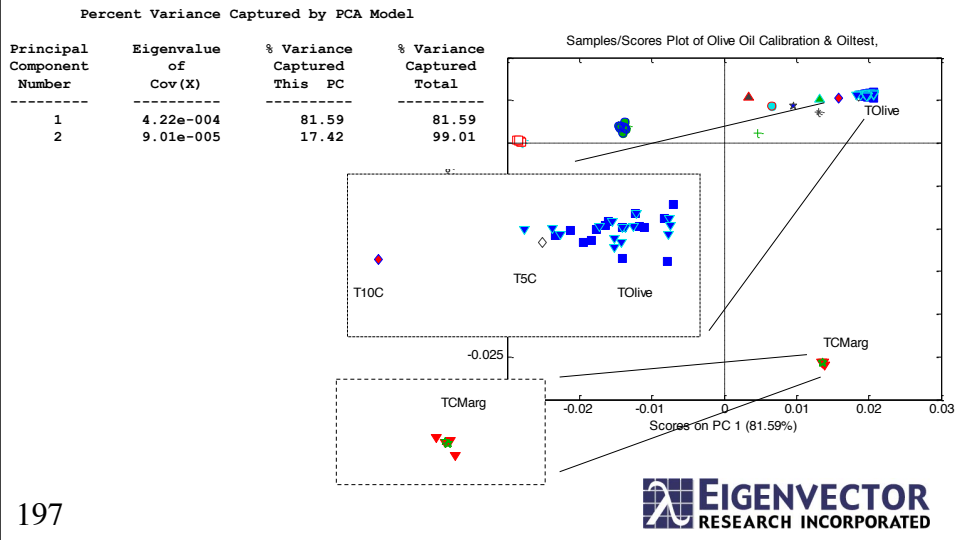
Click Preprocessing Shortcut

Derivative SavGol: Add -->

Derivative Order 2

change ensure derivative is at top of list (use up / down)

Derivative Followed by Row Normalization



GLS Weighting

- GLS Weighting of the spectral data accounted for some of the clutter observed in the spectra.
- The result was
 - clusters that were further apart and
 - clusters that were tighter
 - the ratio of between-class to within-class variance was increased making discrimination easier
 - clusters were so tight and far apart that confidence bounds defining each class could be wider

198

Extended Multiplicative Scatter Correction (EMSC)

- EMSC attempts to account for
 - clutter by using an **extended mixture model** and
 - multiplicative effects like multiplicative scatter correction (MSC)
 - The **extended mixture model** is a classical least squares-like model that is used to explicitly account for clutter (a.k.a. extended least squares).

199



EMSC

Provide spectra of:

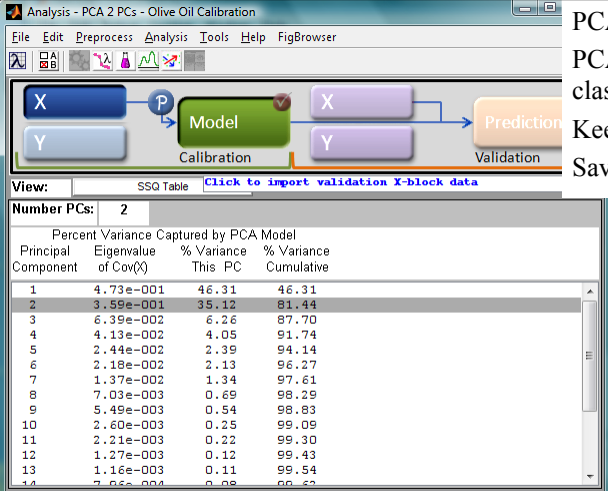
- Known **target analytes** **S**
- Polynomial baselines **P**
- Known **interferences** **Q**
 - e.g., loadings from a PCA model of clutter
 - the coefficients for each linear effect are estimated using least-squares (indicated by "hat")

$$\mathbf{s}_{2, \text{measured}} = \mathbf{s}_{\text{ref}} \mathbf{c}_{\text{ref}} + \mathbf{S} \mathbf{c}_S + \mathbf{P} \mathbf{c}_P + \mathbf{Q} \mathbf{c}_Q \quad \mathbf{P} = \begin{bmatrix} \mathbf{L} & \mathbf{v}^2 & \mathbf{v} & 1 \end{bmatrix}$$

$$\mathbf{s}_{2, \text{corrected}} = (\mathbf{s}_2 - \mathbf{P} \hat{\mathbf{c}}_P - \mathbf{Q} \hat{\mathbf{c}}_Q) / \hat{c}_{\text{ref}} \quad \mathbf{Q} = \text{loadings}$$

200





Analysis - PCA 2 PCs - Olive Oil Calibration

View: SSO Table [Click to import validation X-block data](#)

Number PCs: 2

Principal Component	Eigenvalue	% Variance of Cov(X)	% Variance This PC	% Variance Cumulative
1	4.73e-001	46.31	46.31	46.31
2	3.59e-001	35.12	81.44	81.44
3	6.39e-002	6.26	87.70	87.70
4	4.13e-002	4.05	91.74	91.74
5	2.44e-002	2.39	94.14	94.14
6	2.18e-002	2.13	96.27	96.27
7	1.37e-002	1.34	97.61	97.61
8	7.03e-003	0.69	98.29	98.29
9	5.49e-003	0.54	98.83	98.83
10	2.60e-003	0.25	99.09	99.09
11	2.21e-003	0.22	99.30	99.30
12	1.27e-003	0.12	99.43	99.43
13	1.16e-003	0.11	99.54	99.54
14	7.65e-004	0.08	99.62	99.62

A model has been calibrated from the data. Review the model using the toolbar button(s) save the model (File)

```

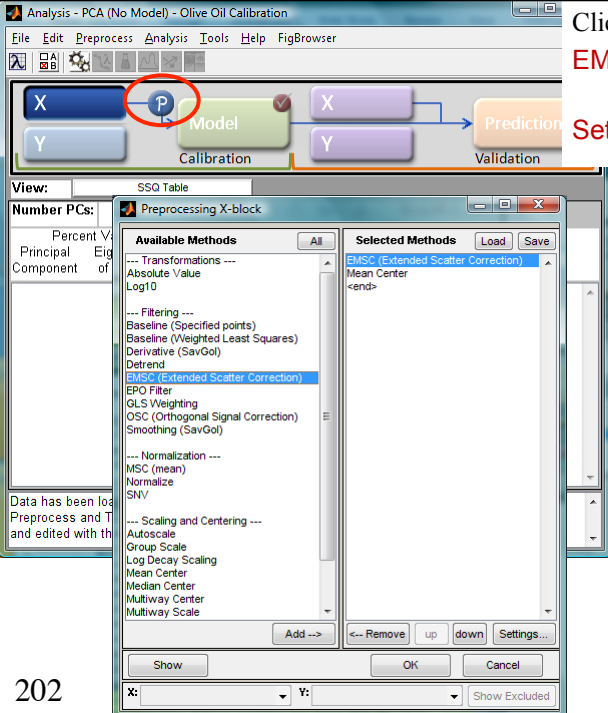
preprocess
>> z = xcal;
>> for i1=1:4
    z.data(find(xcal.class{1}==i1), :) = mncn(z.data(find(xcal.class{1}==i1), :));
end
>> z.description = char(z.description, 'Each class center to its own mean. ');
>> p = zeros(2, 518);
>> p(:, z.include{2}) = pcam.loads{2}';

```

201

EIGENVECTOR
RESEARCH INCORPORATED

PCA of clutter.
PCA of calibration data with
classes centered to class mean.
Keep 2 PCs to model the clutter.
Save model to pcam



Analysis - PCA (No Model) - Olive Oil Calibration

View: SSO Table

Number PCs: 2

Preprocessing X-block

Available Methods

- Transformations ---
 - Absolute Value
 - Log10
- Filtering ---
 - Baseline (Specified points)
 - Baseline (Weighted Least Squares)
 - Derivative (SavGol)
 - Detrend
 - EMSC (Extended Scatter Correction)
 - EPO Filter
 - GLS Weighting
 - OSC (Orthogonal Signal Correction)
 - Smoothing (SavGol)
- Normalization ---
 - MSC (mean)
 - Normalize
 - SNV
- Scaling and Centering ---
 - Autoscale
 - Group Scale
 - Log Decay Scaling
 - Mean Center
 - Median Center
 - Multway Center
 - Multway Scale

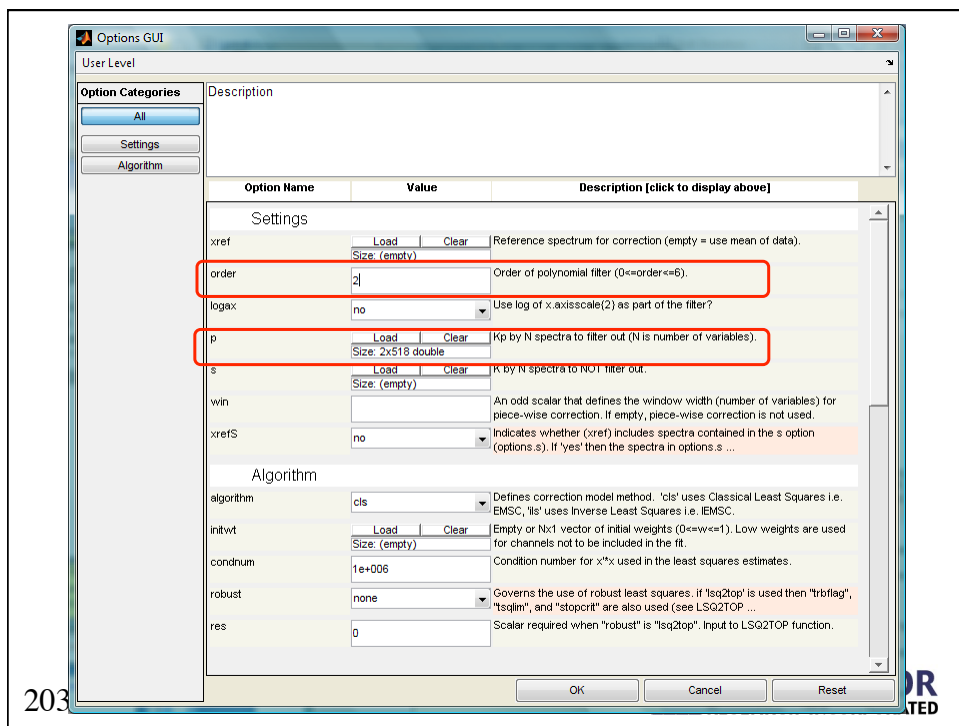
Selected Methods

- EMSC (Extended Scatter Correction)
- Mean Center
- <end>

202

EIGENVECTOR
RESEARCH INCORPORATED

Click Preprocessing Shortcut
EMSC (Extended Scatter
Correction)
Settings...

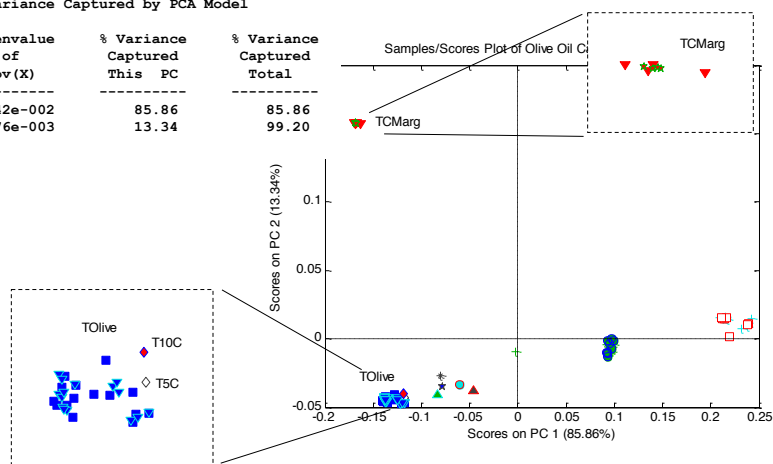


203

EMSC

Percent Variance Captured by PCA Model

Principal Component Number	Eigenvalue of Cov(X)	% Variance Captured This PC	% Variance Captured Total
1	2.42e-002	85.86	85.86
2	3.76e-003	13.34	99.20



204

EMSC Summary

- EMSC attempts to account for clutter in an explicit way
 - e.g., model clutter with basis vectors (e.g., PCA loads)
 - analyst takes control of the model
 - requires good use of measurements: clutter and target spectra
 - use what you know!
 - interpretable
 - analyst control is more daunting than using simple SavGol and MSC, but
 - the results are much more interpretable than 2nd derivative spectra

205



Analysis - PCA 2 PCs - Olive Oil Calibration (2)

File Edit Preprocess Analysis Tools Help FigBrowse

X-block
Y-block
Load Preprocessing
Save Preprocessing
Plot Preprocessed Data
Calibration
Validation
X-block
Y-block
Prediction
Validation

View: SSQ Table

Number PCs: 2

Principal Component	Eigenvalue of Cov(X)	% Variance This PC	% Variance Cumulative
1	2.42e-002	85.86	85.86
2	3.76e-003	13.34	99.20
3	9.36e-005	0.33	99.53
4	5.83e-005	0.21	99.74
5	1.73e-005	0.06	99.80
6	1.17e-005	0.04	99.84
7	1.00e-005	0.04	99.88
8	8.27e-006	0.03	99.91
9	4.42e-006	0.02	99.92
10	3.90e-006	0.01	99.94
11	3.30e-006	0.01	99.95
12	2.23e-006	0.01	99.96
13	1.77e-006	0.01	99.96
14	1.48e-006	0.01	99.97

Warning: This model appears to have some unusual Q contributions using the Scores plot and determine if the removed. If these are not errors, consider consider adding them back.

Click Preprocessing Menu
Preprocess:Plot Preprocessed Data: Calibration: X-block
Plot:Rows
View:Classes:Oil

Plot Controls

File Edit View Plot FigBrowse

Fig 3: Calibration

X: Wavenumber

Y: Row 29
Row 30
Row 31
Row 32
Row 33
Row 34
Row 35
Row 36

Z: none

Color By...

Plot [x] auto-update

Select Tool

Table Ctrl+T
Numbers Ctrl+U
Labels Ctrl+L
Classes Ctrl+Z
Excluded Data Ctrl+E
Declobber Labels Ctrl+K
Label Angle Ctrl+A
Axis Lines
Diagonal 1:1 line
Log Scales
Auto Y-Scale Ctrl+F
Subplots
Duplicate Figure Ctrl+D
Spawn Static View
Dock Controls
Settings...

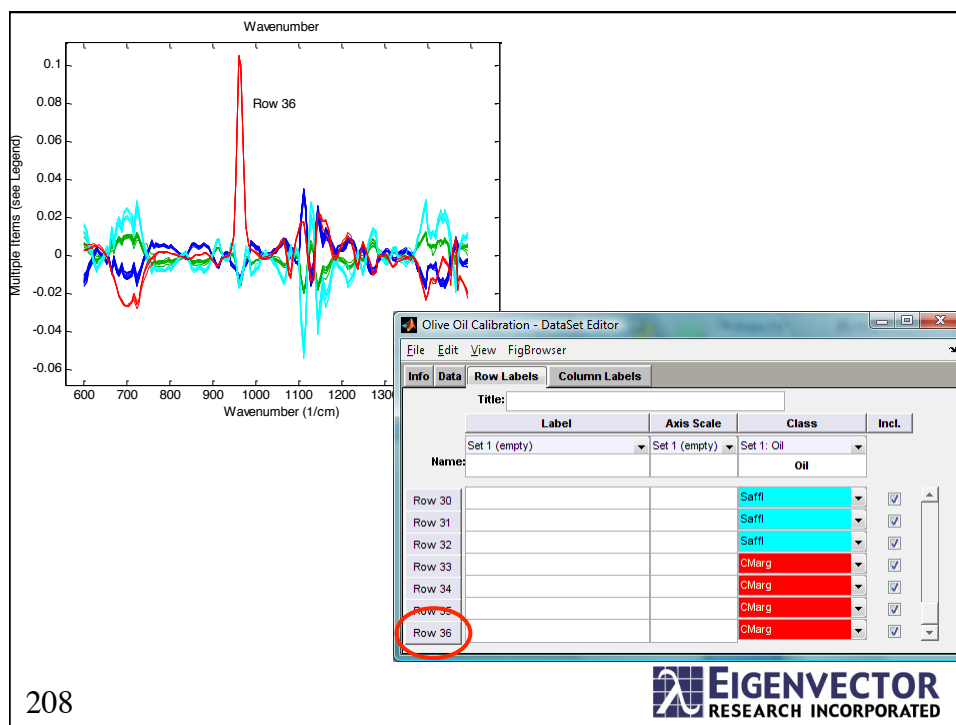
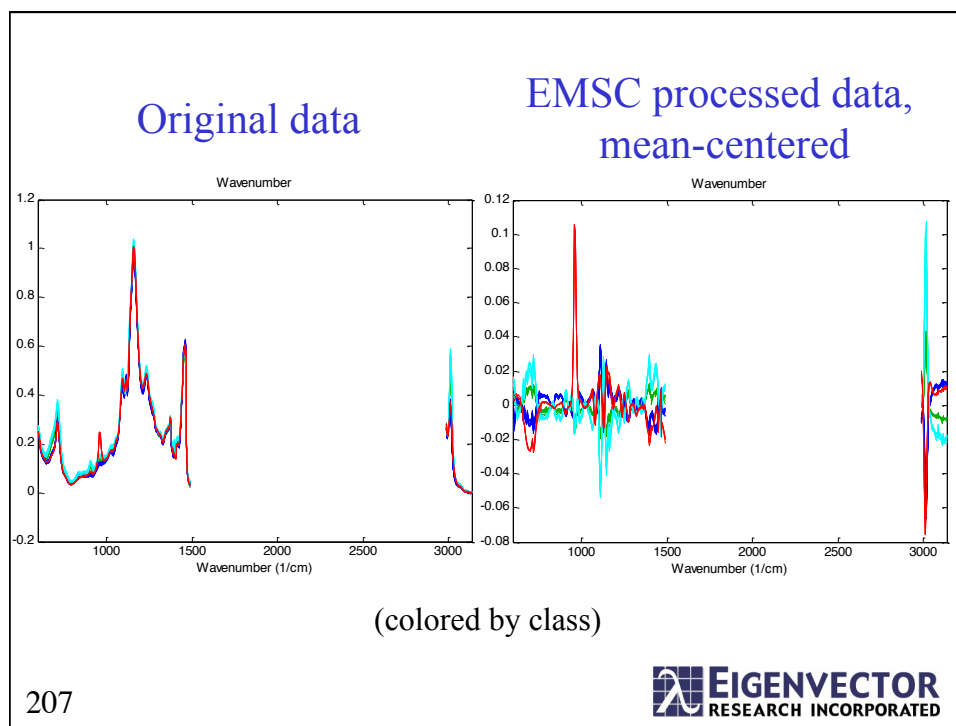
View: SSQ Table

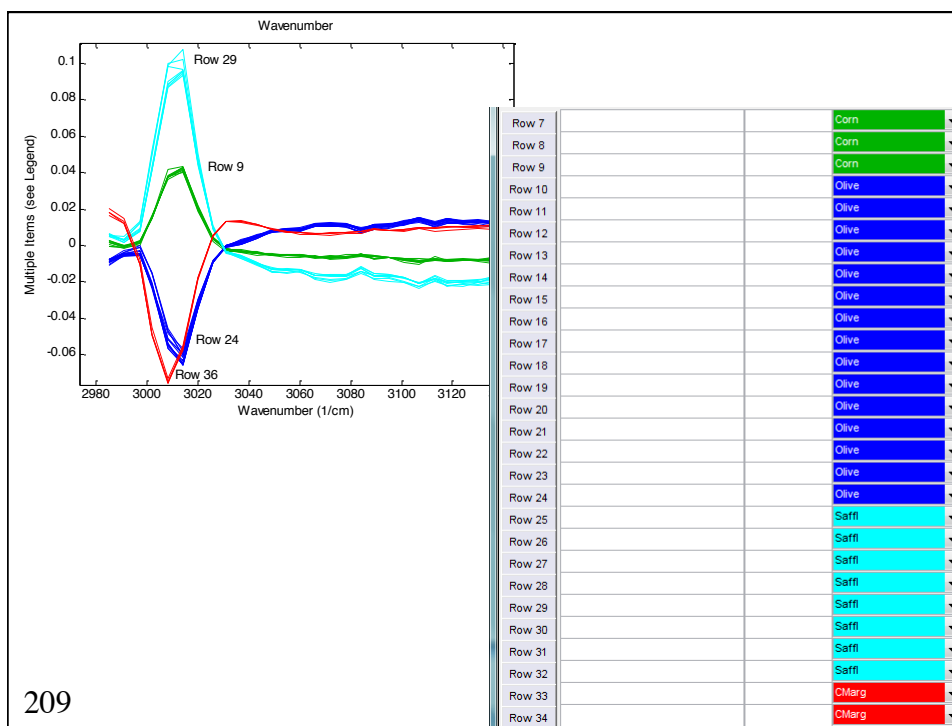
Number PCs: 2

Principal Component	Eigenvalue
4	5.83e-005
5	1.73e-005
6	1.17e-005
7	1.00e-005
8	8.27e-006
9	4.42e-006
10	3.90e-006
11	3.30e-006
12	2.23e-006
13	1.77e-006
14	1.48e-006

Warning: This model appears to have some unusual Q contributions using the Scores plot and determine if the removed. If these are not errors, consider consider adding them back.

206





Section Definitions

- **Between class variance:** The sum-of-squares of the class means centered to the global data set mean divided by number of classes.
- **Within class variance:** The sum-of-squares of each class centered to the class mean divided by number of samples in the class.
- **Multiplicative scatter correction (MSC):** (a.k.a. Multiplicative Signal Correction) Data pretreatment that removes multiplicative effects and baseline offset based on a reference *e.g.* a reference spectrum.
- **Savitzky-Golay Smoothing and Differentiation:** Numerical method for calculating the derivative of a spectrum that uses windowed polynomials.
- **Normalization:** the 2-norm divides a spectrum by the square root of the sum-of-squared signal in each frequency channel. This removes magnitude information from the spectrum.
- **Extended mixture model:** a classical least squares-like model that is used to explicitly account for clutter.
- **Extended multiplicative scatter correction (EMSC):** a model that combines MSC and the extended mixture model to explicitly account for clutter.

210

Outline

- Introduction
- Advanced Preprocessing
 - Clutter and characterizing clutter
 - Generalized least squares weighting
 - Extended multiplicative scatter correction
 - Interval PLS (iPLS)
 - Model Robustness
- Multivariate image analysis
- Multi-way Analysis
- Summary

©Copyright 2008-2012
Eigenvector Research, Inc.
No part of this material may be photocopied or
reproduced in any form without prior written
consent from Eigenvector Research, Inc.

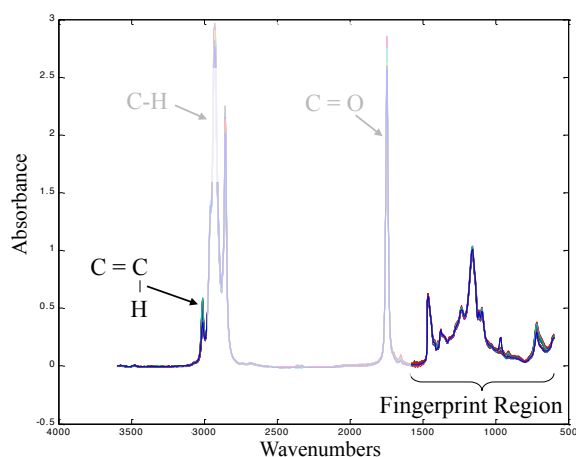


Why Variable Selection?

- Improvement of the model
 - Remove irrelevant, unreliable or noisy variables (clutter)
 - Improve predictions
 - Improve statistical properties
- Interpretation
 - Obtain a model that is easier to understand
- Costs
 - Use fewer measurements to replace expensive or time-consuming one
- Development of fast instruments/routines for on-line control
 - Find wavelength ranges for a filter-based instrument



Already done some based
on *a posteriori* knowledge...



213



Variable Selection Methods

- ***a priori***
 - Choose measurements
- ***a posteriori***
 - Use chemical/physical insight
- **Model based**
 - Look at loadings
- **"Random based"**
 - Genetic algorithms
 - Simulated annealing
- **"Spectral"**
 - i-PLS
 - fullsearch
- **Classical**
 - Forward, backward selection
 - Best subset selection
 - Significance tests
 - Significance based on Jack-knife
 - GOLPE
- **Other**
 - Pure variables
 - Principal variables
 - Iterative weighting with regression vector
 - ...

214

(see the Variable Selection Course at EigenU)



Variable Selection Methods

- How to choose which method?!?
- Different methods work in different situations
- Interval-PLS is a good “example” method to understand the considerations of variable selection. Simple to implement and use.

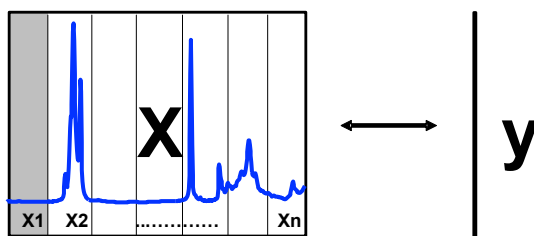
215



i-PLS Theory

iPLS: Interval PLS

Build local models using “intervals” of X-block variables. Very intuitive and useful approach that can be easily combined with variable selection.



216

L. Nørgaard, A. Saudland, J. Wagner, J. P. Nielsen, L. Munck, S. B. Engelsen. Interval partial least-squares regression (iPLS). *Appl.Spectrosc.* 54 (3):413-419, 2000.



Fit Criteria

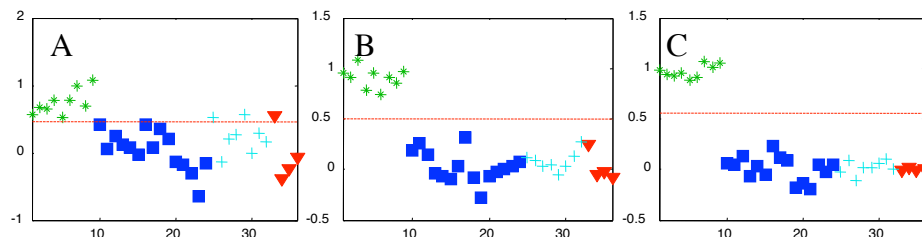
- RMSECV is used to determine “best” interval
- For Olive Oil data: perform discriminant analysis using PLSDA (use logical y-block with PLS to separate classes).
- Note: Always validate afterwards! Variable selection methods have a tendency to give over-optimistic RMSE results.

217



RMSE vs. Misclassification Rate in PLSDA

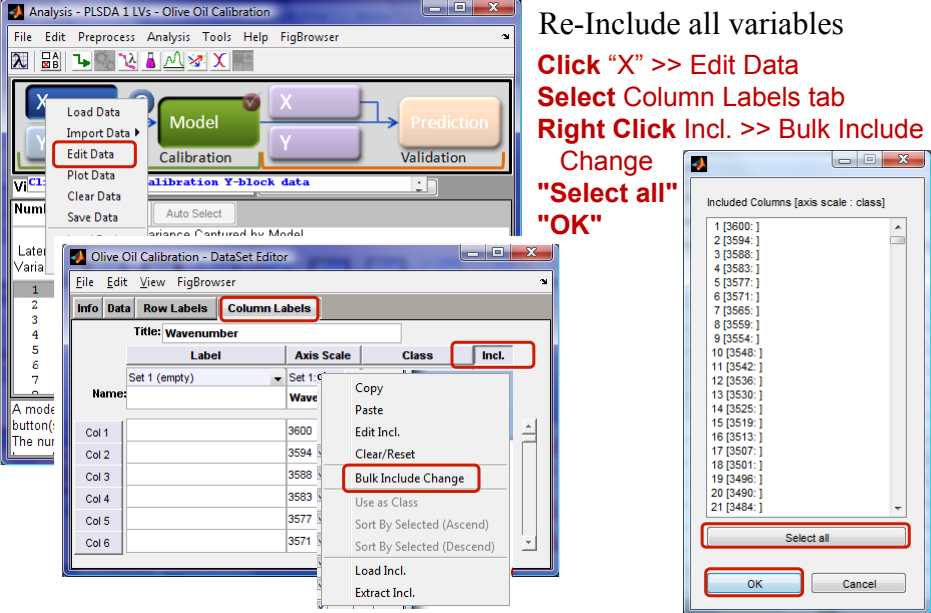
- RMSE shows deviation from predicting 0 or 1.
- Misclassification Rate shows prediction on "wrong side" of decision line.
- RMSE: $A > B > C$ Misclassification: $A > B = C$



218



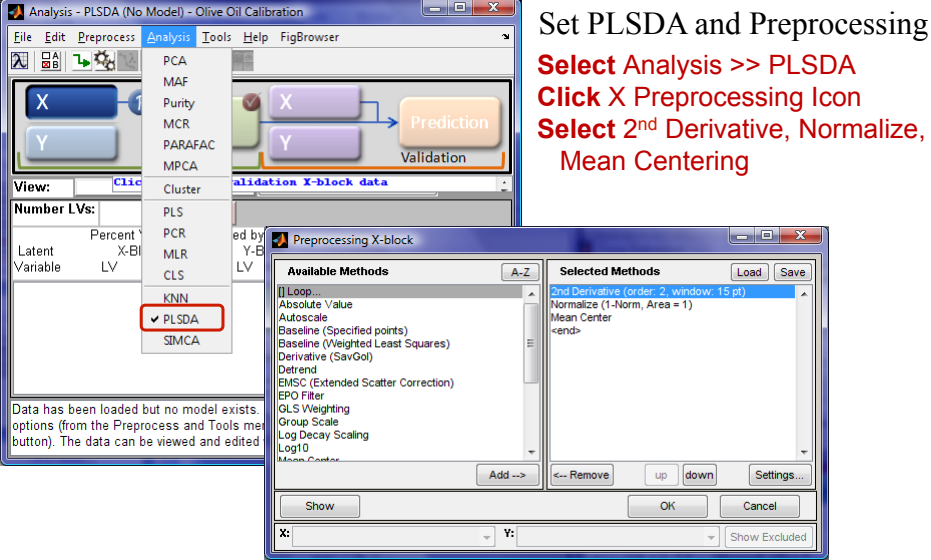
Re-Include all variables
Click "X" >> Edit Data
Select Column Labels tab
Right Click Incl. >> Bulk Include
Change "Select all"
"OK"



219

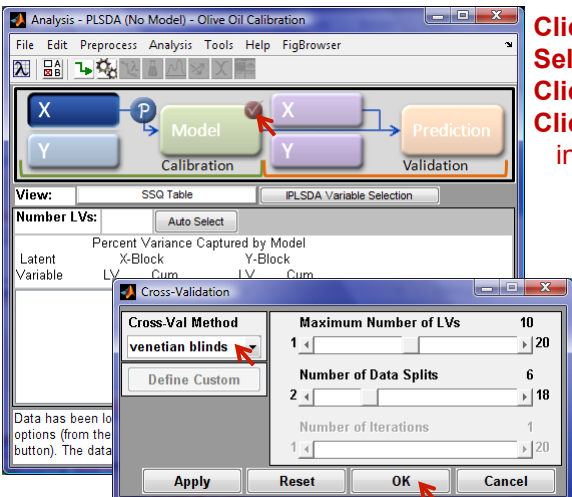
EIGENVECTOR
RESEARCH INCORPORATED

Set PLSDA and Preprocessing
Select Analysis >> PLSDA
Click X Preprocessing Icon
Select 2nd Derivative, Normalize, Mean Centering



220

EIGENVECTOR
RESEARCH INCORPORATED

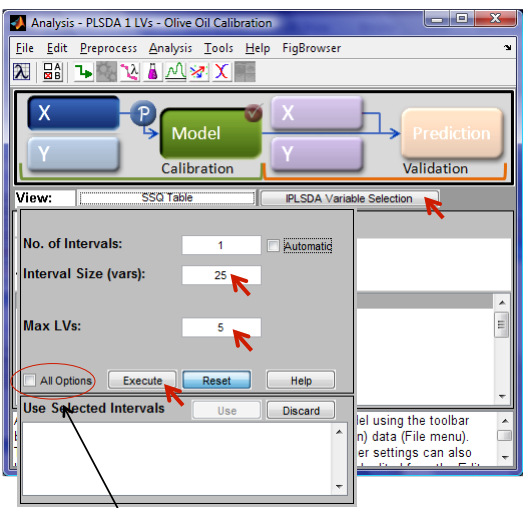


Click Cross Validation Icon
Select "Venetian Blinds"
Click OK
Click "Model" (to build model initial PLSDA model)

Normally, should now review scores for outliers!

221

EIGENVECTOR
RESEARCH INCORPORATED



Click "PLSDA Variable Selection" button
Enter 25 for Interval Size
Enter 5 for max LVs
Click "Execute"

Interval Size defines how many variables to group together.

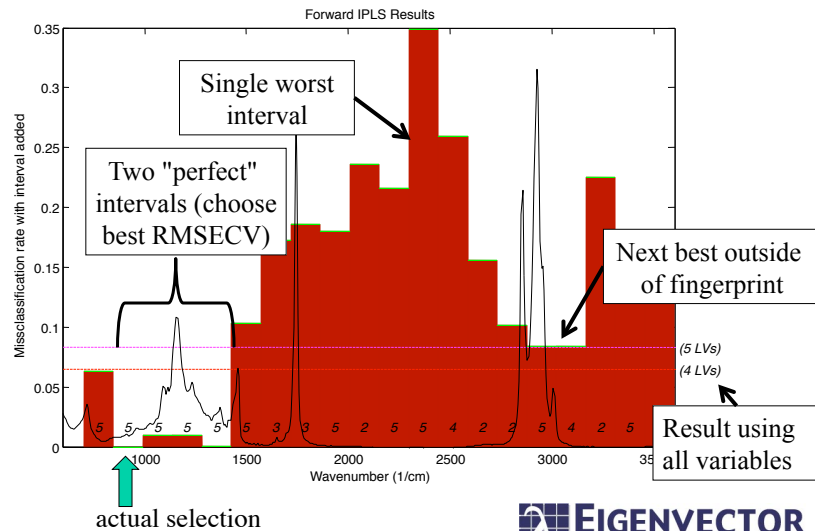
Spectral data: use size>1
 Non-contiguous data: use size=1 (single variable intervals)

provides access to other i-PLS settings

222

EIGENVECTOR
RESEARCH INCORPORATED

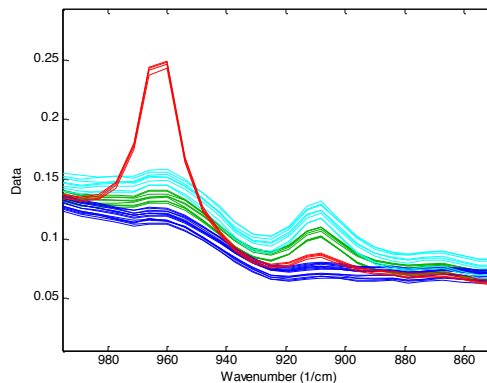
First Interval Result



223

Why the Low-Frequency Intervals? Information Content

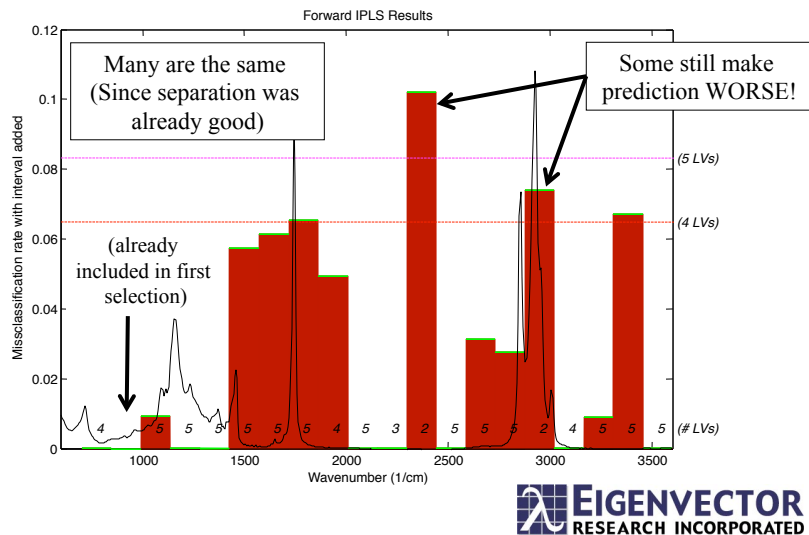
- Best intervals contain useful signal from all four classes.
- Original data shown; Preprocessing only makes this separation better!



EIGENVECTOR
RESEARCH INCORPORATED

224

Repeat Execution and “Add To” Existing



225

Number of Intervals

- Can choose a pre-set number of intervals to find
- Can also use “Automatic” to continue selecting intervals until RMSECV/misclassification does not improve
- This is **not** the same as exhaustive combinatorial search (fullsearch). It is sequential (choose one, “lock” it in, choose a second, “lock” it in...)
- For very complex data, may not give actual “best” windows, but probably not a bad one.

226

EIGENVECTOR
RESEARCH INCORPORATED

What is the Result?

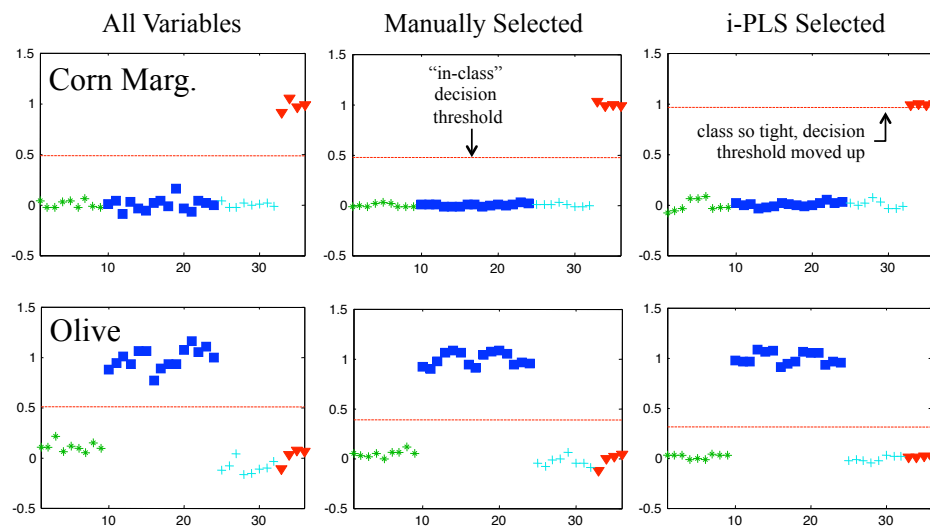
- For PLSDA, lower RMSECV should indicate better class separation in predicted Y values
- Selecting additional intervals gives little improvement in RMSECV (on this data)
- Use ONLY first selected interval and build new model...



227



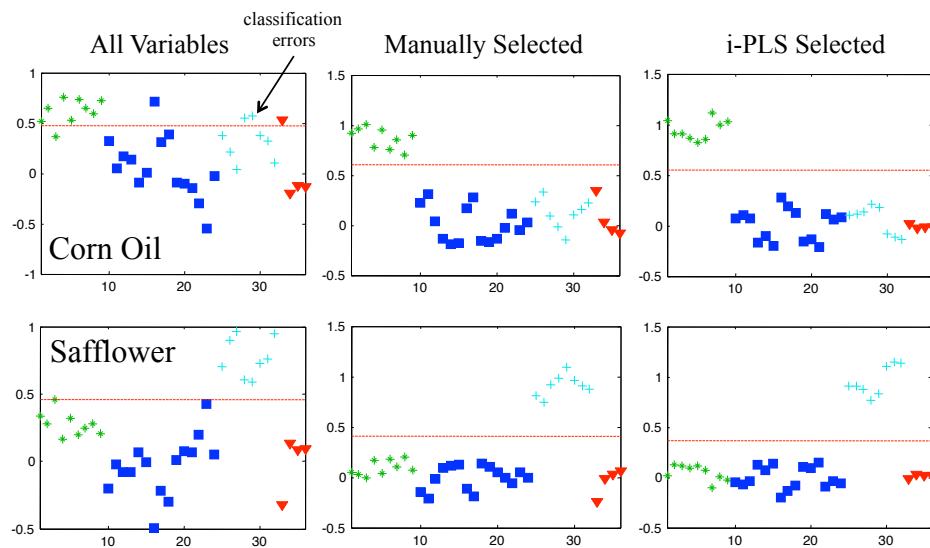
Cross-validated Predictions



228 Mild improvement... (note: LVs = 5, 5, 3)



Cross-validated Predictions



229 i-PLS: ~tighter clusters than manual



Reverse i-PLS

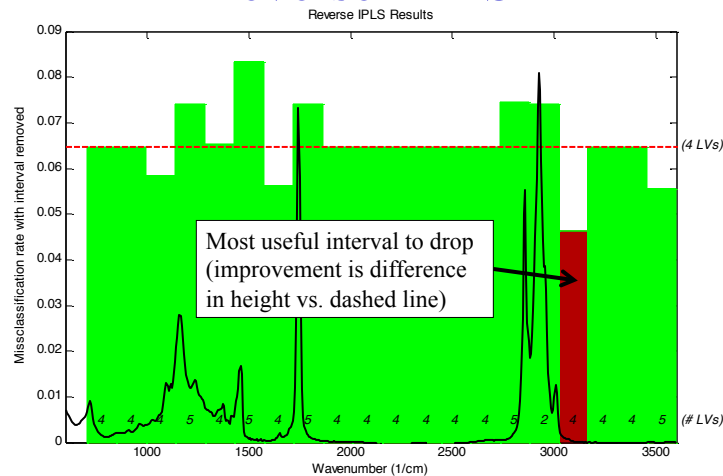
- **Principle**
 - Make full model and select the variable contributing the least to the fit (exclude regions which contain **more clutter**)
 - Repeat as desired/needed
- **Good**
 - Takes interactions into account
 - Reasonably fast
- **Bad**
 - "Random" removal for large data sets
 - Often works bad for many irrelevant variables

(see "All options" checkbox on i-PLS controls)

230



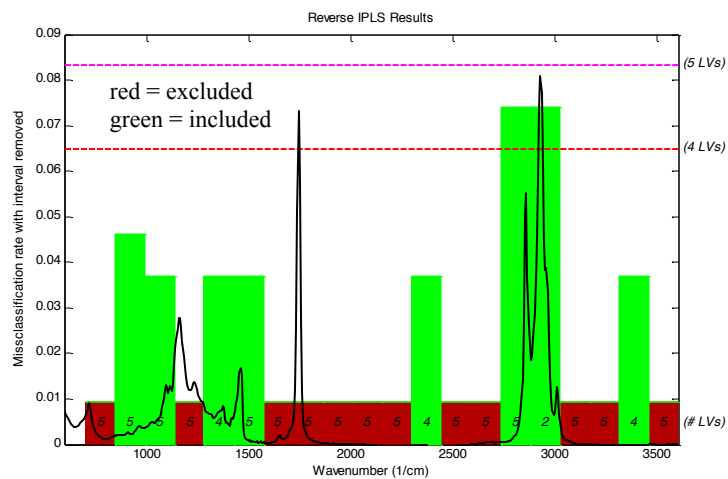
Reverse i-PLS



231



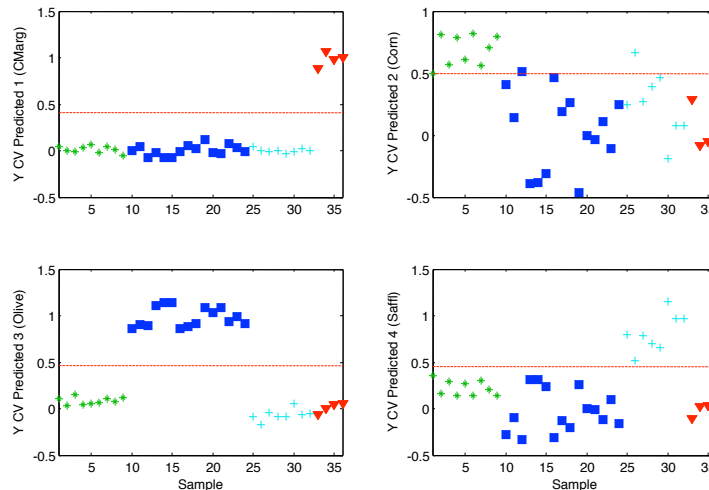
Automatic Reverse i-PLS



232



Cross-validated Predictions after Automatic Reverse i-PLS



233 8 intervals retained. Class 2 and 4 not as good as forward i-PLS!



What is Model Robustness?

- When developing calibration models focus is generally on improving prediction error
- Models often developed with small amount of data taken over relatively short time
- Prediction errors over long term often dominated by artifacts not represented in calibration data
 - Changes in spectrometer / sensor
 - Changes in sample

234



Typical Changes in the System

- Sample
 - New analyte(s)
 - Changes in physical properties (e.g. scattering)
 - Temperature
 - Pressure
- Instrument (spectrometers)
 - Wavelength/Frequency registration shift
 - Stray light
 - Resolution
 - Noise

235



What Constitutes a Good Model?

- Acceptable prediction error (not necessarily the best achievable)
- Longevity, i.e. **robustness** to minor changes
- Once you have built a model, you should exercise it with expected changes
- Use real data or simulate typical instrumental and multivariate errors

236



Robustness Testing

- Develop model with desired preprocessing, #LVs, etc.
- “Perturb” test data set
- Apply calibration model to “perturbed” data
- Look at prediction error as function of perturbations
- Test and compare multiple models

237



Perturbations

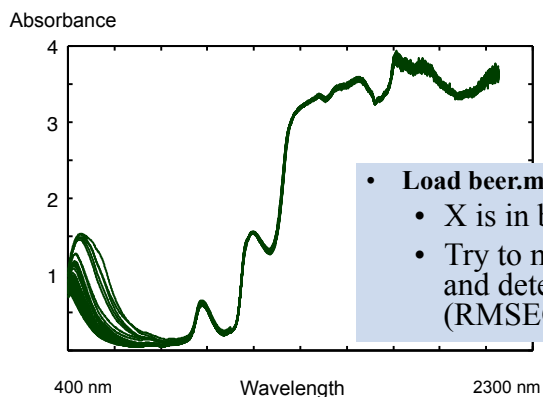
- New analyte – add Gaussian peak of variable width across wavelength range
- Wavelength registration shift – shift spectra left-right as well as expand and contract
- Baseline shift – change offset and slope
- Stray light – add fraction of signal before log transform
- Temperature – decrease resolution and vary path length
- Noise variation – add noise with varying bandwidth

238



VIS/NIR spectra of 61 beers

Purpose: prediction of real extract



- **Load beer.mat**
 - X is in beer and Y in extract
 - Try to make a nice PLS model and determine quality (RMSEC, RMSECV)

239



Exercise Data

Determination of the amount of extract from NIR spectra of beers.

Dispersive visual & near-infrared data collected (at 25 C) NIRSystems Inc. (Model 6500) spectrophotometer. Split detector system – silicon detector 400-1100 nm & (PbS) detector 1100-2500 nm.

VIS-NIR transmission recorded directly on undiluted degassed beer in 30 mm quartz cell. Spectral data collected at 2 nm intervals 400-2250 nm & converted to absorbance units.

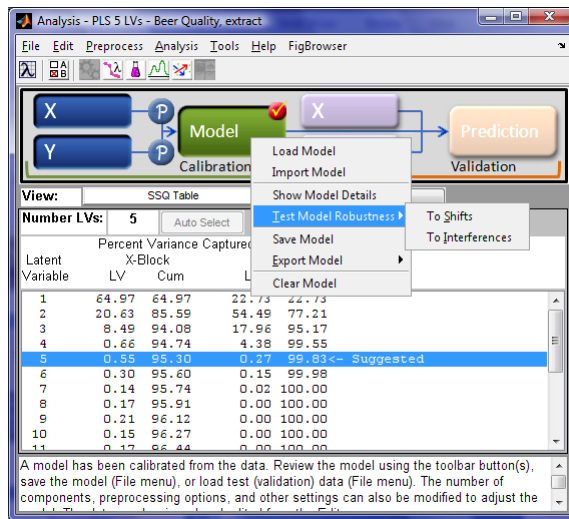
Original *extract* concentration is a quality parameter in the brewing industry, indicating the substrate potential for the yeast to ferment alcohol and serving as a taxation parameter. Original extract concentration determined by Carlsberg A/S in the range of 4.23-18.76% plato.

Data sorted by extract value, and a model independent test set was constructed by selecting every third sample of this full data set. There are thus two data sets: one for calibration (40 samples) and one for independent estimation of prediction error (20 samples).

240



Test Robustness



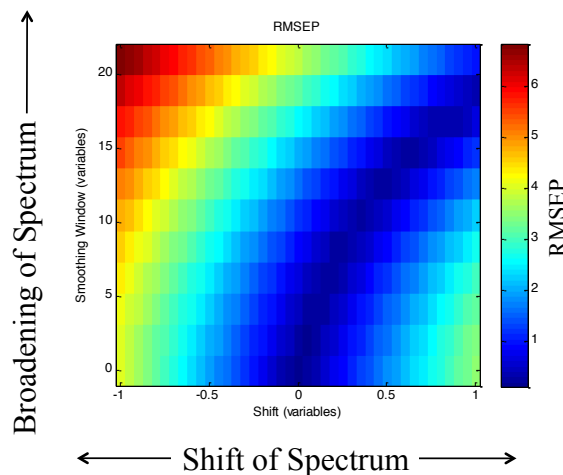
Two built-in tests:

- Shift / Resolution Test
- Interference Test

241



Shift / Resolution Test

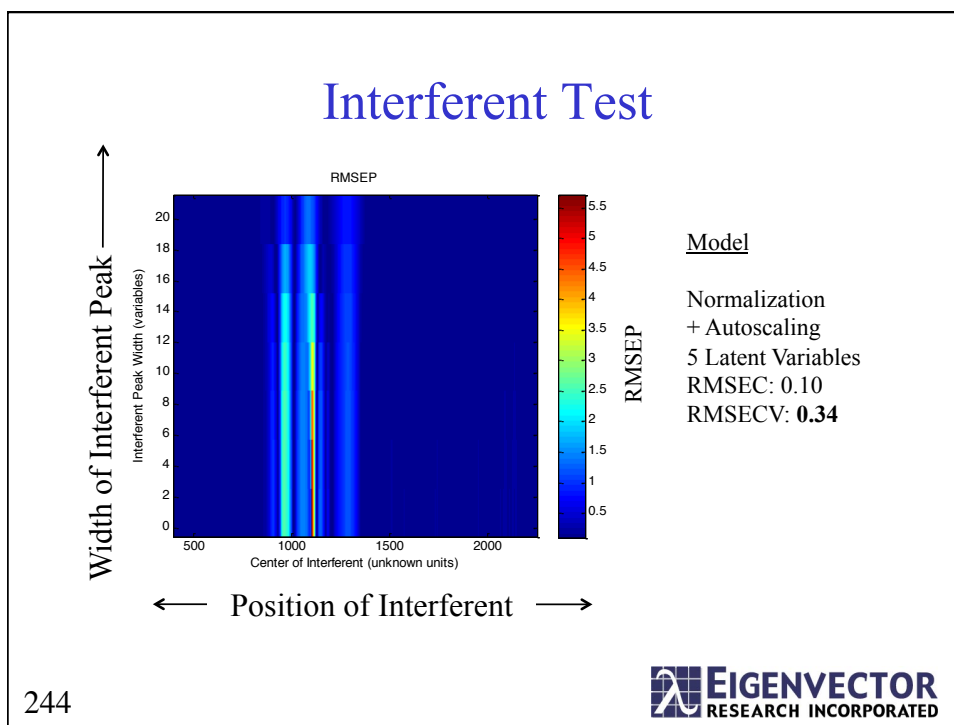
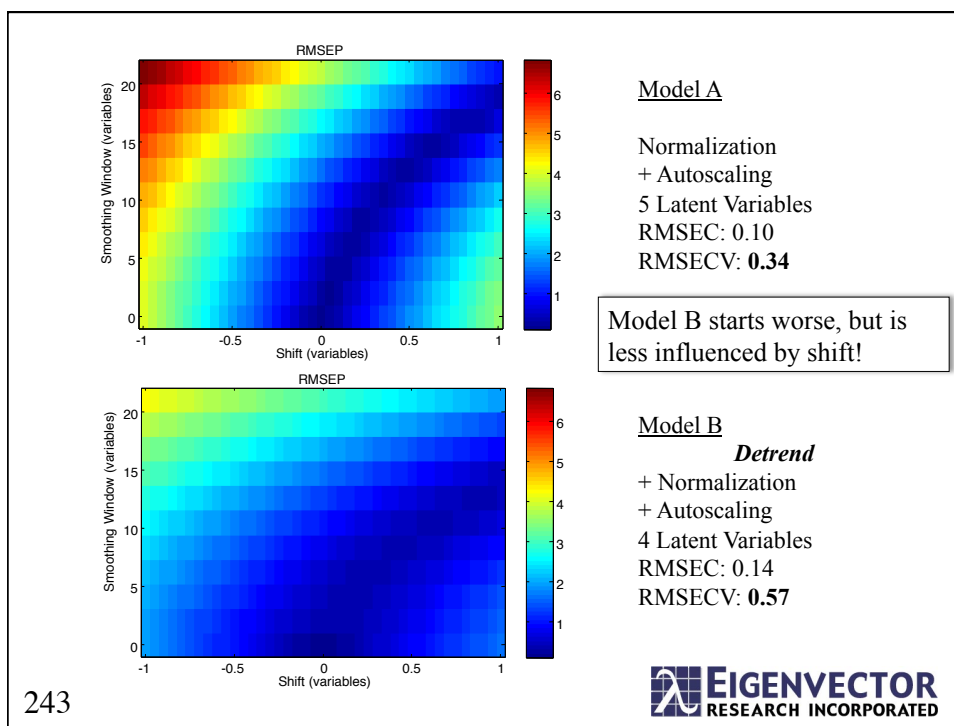


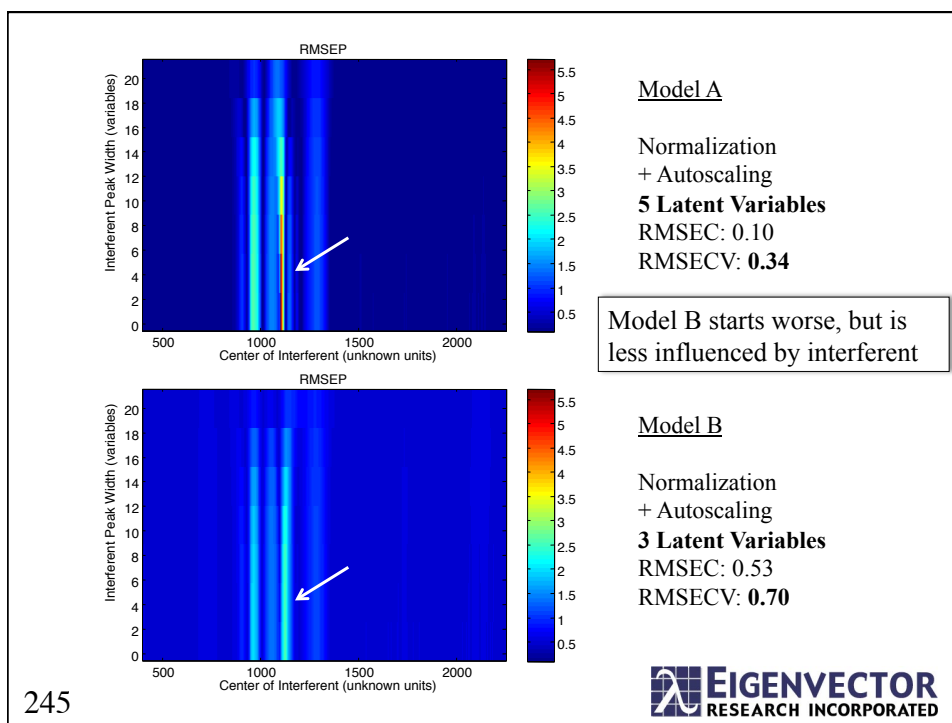
Model
 Normalization
 + Autoscaling
 RMSEC: 0.10
 RMSECV: 0.34

Hint: use "caxis"
 command to get color
 range and
 "caxis([min max])"
 command to set color
 range

242







Outline

- Introduction
- Advanced Preprocessing
- **Multivariate image analysis**
- Multi-way Analysis
- Summary

©Copyright 2008-2012
Eigenvector Research, Inc.
No part of this material may be photocopied or
reproduced in any form without prior written
consent from Eigenvector Research, Inc.

EIGENVECTOR
RESEARCH INCORPORATED

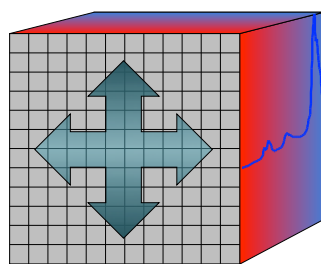
Multivariate Images

A data array of *dimension three* (or more) where the first two dimensions are *spatial* and the last dimension(s) is a function of another variable.

247



Multivariate Images



Spatial Information
between pixels

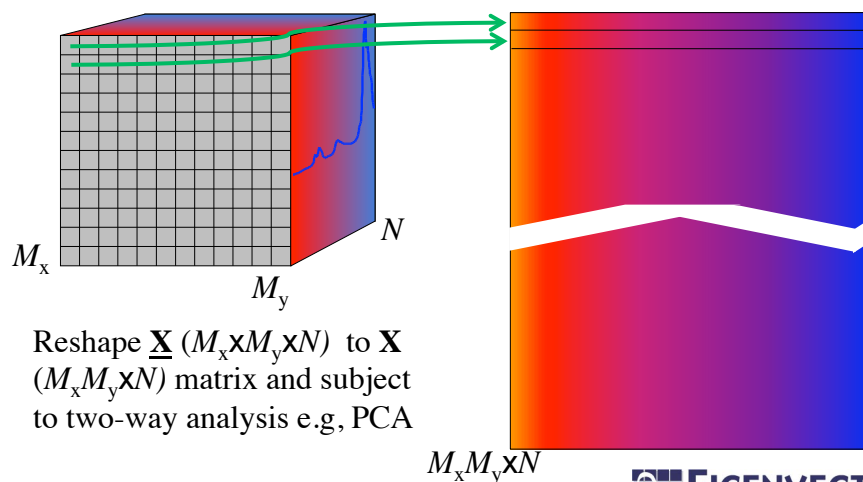
Spectral Information
between channels
(chemical information)

**Spatial distribution of
chemical analytes, physical
features, and other
properties**

248



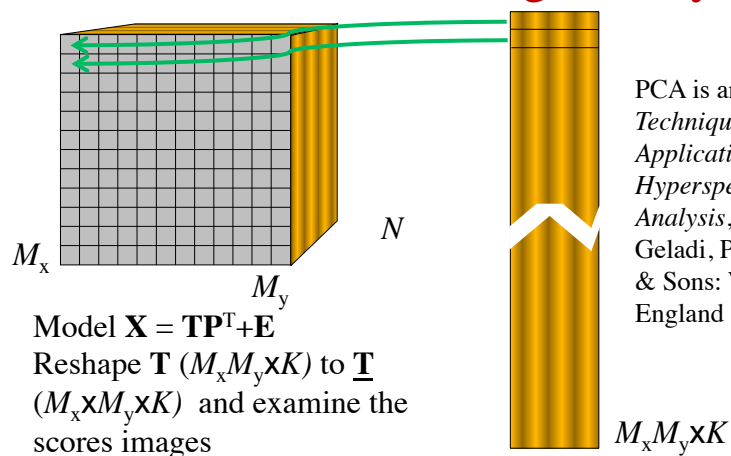
Reshaping Images for Analysis



249

 **EIGENVECTOR**
RESEARCH INCORPORATED

PCA for Multivariate Image Analysis



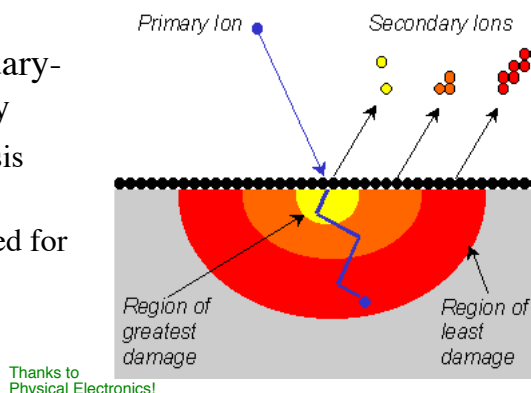
PCA is an example of MIA.
Techniques and Applications of Hyperspectral Image Analysis, Grahn, H. F.; Geladi, P., Eds. John Wiley & Sons: West Sussex, England (2007)

250

 **EIGENVECTOR**
RESEARCH INCORPORATED

Example: TOF-SIMS

- Time-of-Flight Secondary-Ion-Mass Spectrometry
 - common surface analysis technique
 - mass spectrum generated for each pixel



251



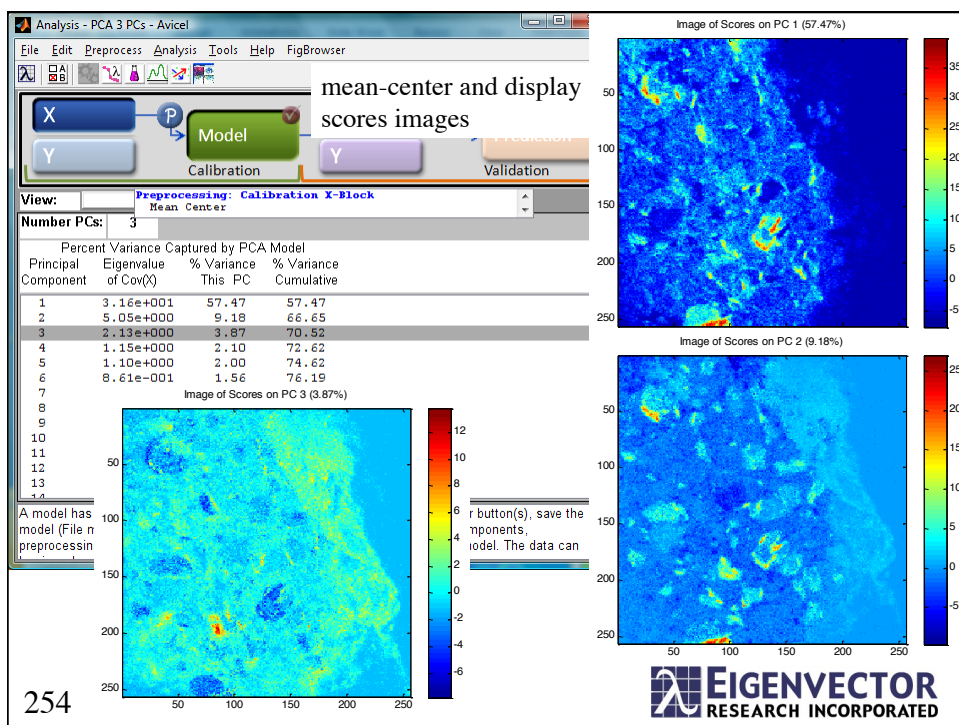
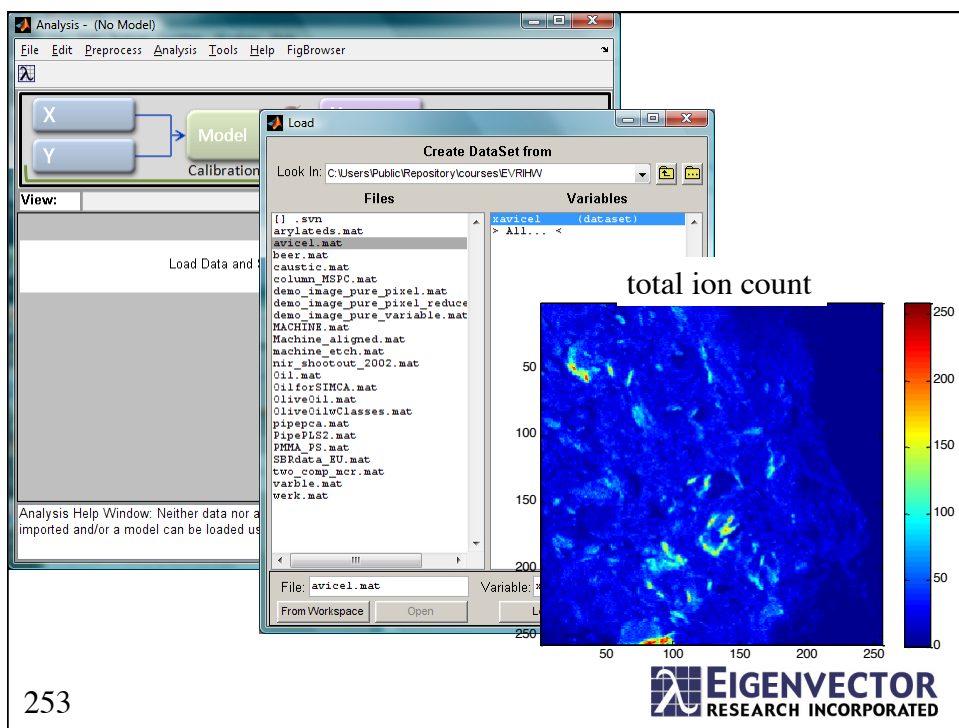
TOF-SIMS of Time Release Drug Delivery System

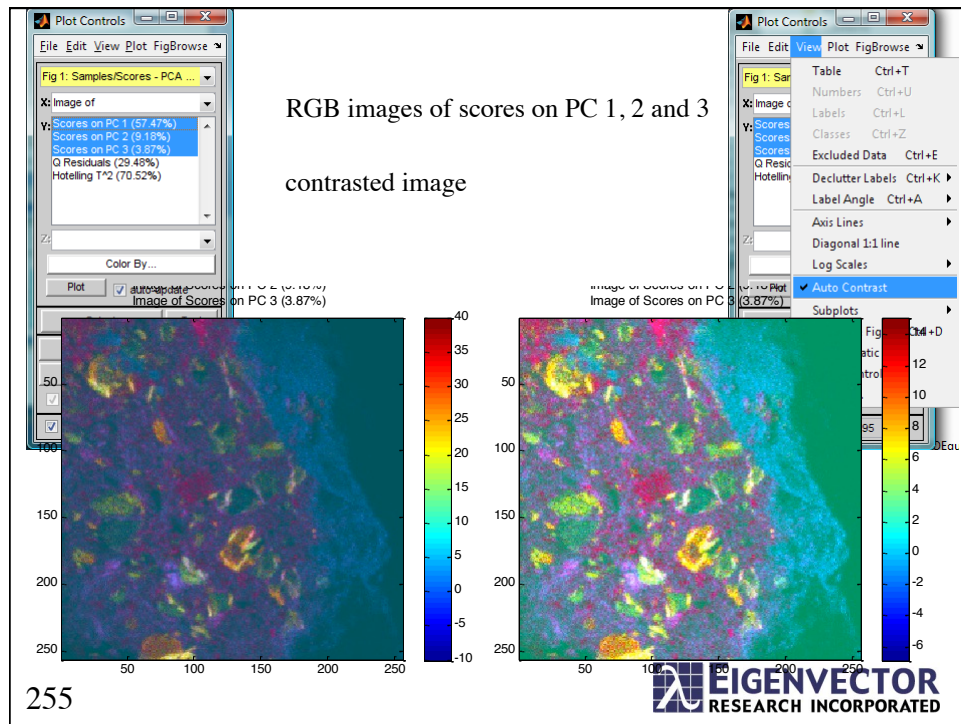
- Multi-layer drug beads serve as controlled release system
- TOF-SIMS of cross section of bead
- Evaluate the integrity of the layers and distribution of ingredients

A.M. Belu, M.C. Davies, J.M. Newton and N. Patel, "TOF-SIMS Characterization and Imaging of Controlled-Release Drug Delivery Systems," *Anal. Chem.*, **72**(22), 5625–5638 (2000).
 Gallagher, N.B., Shaver, J.M., Martin, E.B., Morris, J., Wise, B.M. and Windig, W., "Curve resolution for images with applications to TOF-SIMS and Raman", *Chemometr. Intell. Lab.*, **73**(1), 105–117 (2003).

252







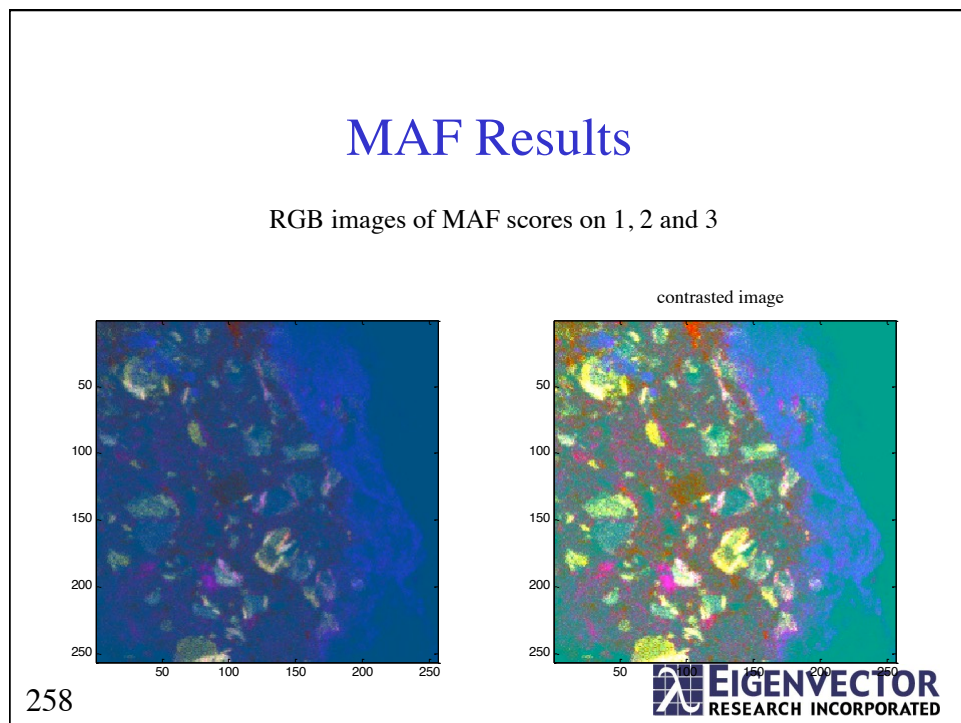
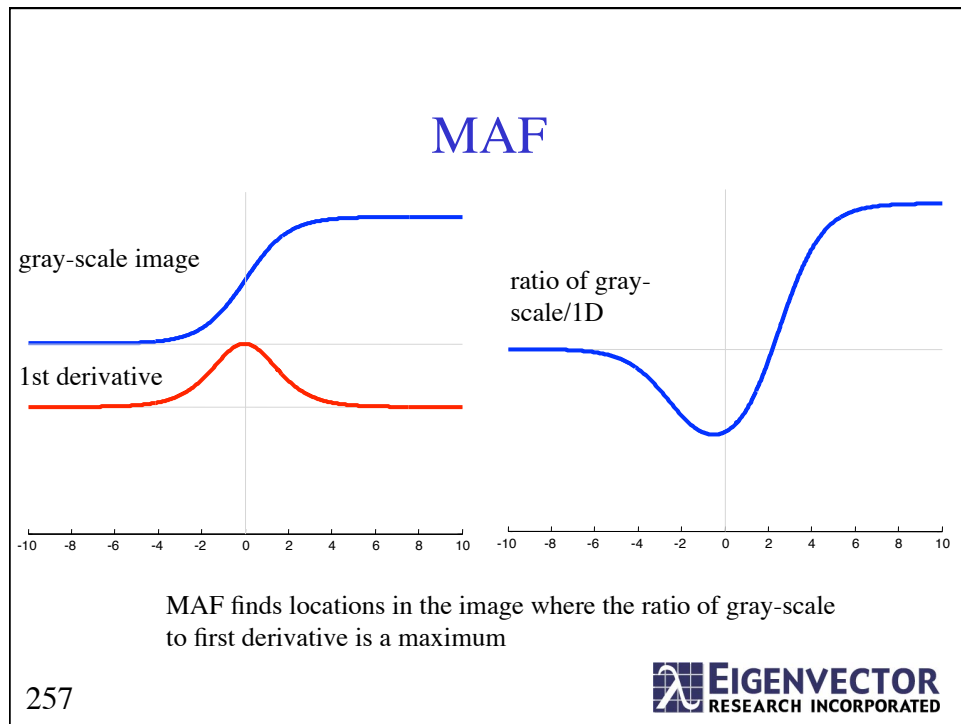
Minimum Noise Factors (MNF)

- MNF attempts find directions in the data that maximize the signal-to-clutter.

$$\max_{\mathbf{v}_i \neq 0} \left(\frac{\mathbf{v}_i^T \Sigma_X \mathbf{v}_i}{\mathbf{v}_i^T \Sigma_C \mathbf{v}_i} \right) \quad \text{the objective function}$$

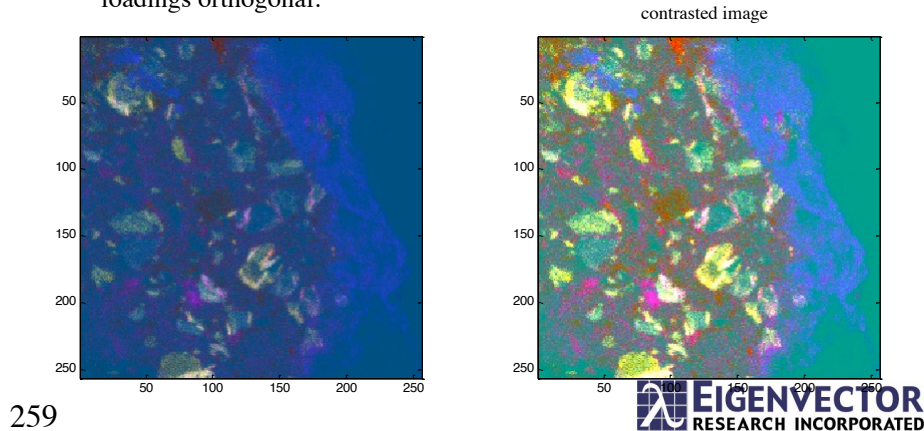
- Result is a PCA-like eigenvector problem
- In **maximum autocorrelation factors (MAF)** clutter is the first difference image (difference between near-by pixels)

256



PCA w/ GLS Weighting for ~MAF

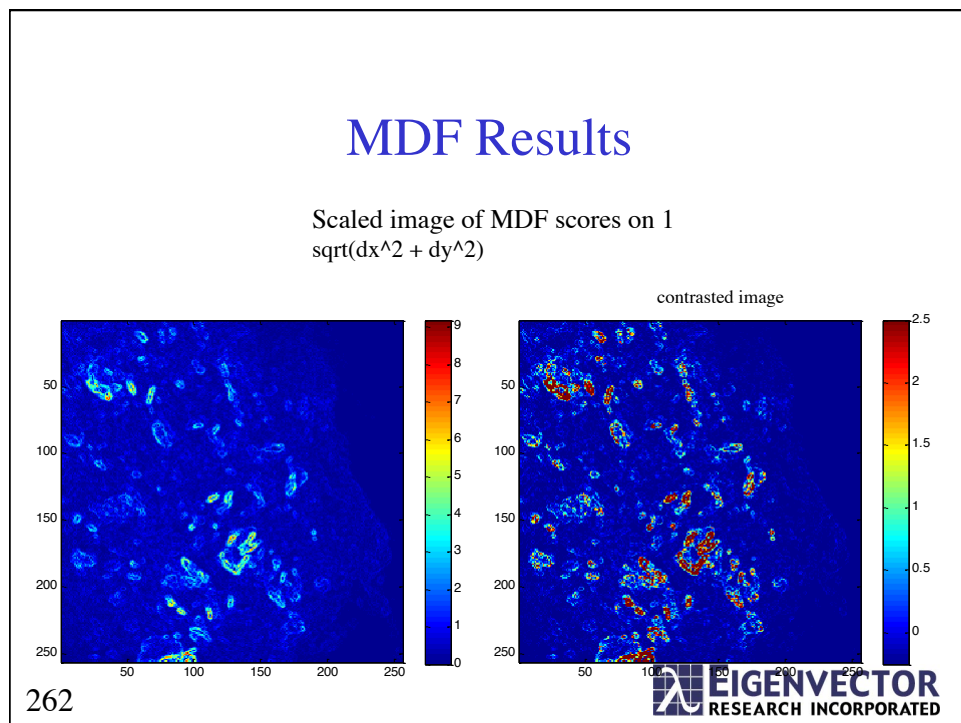
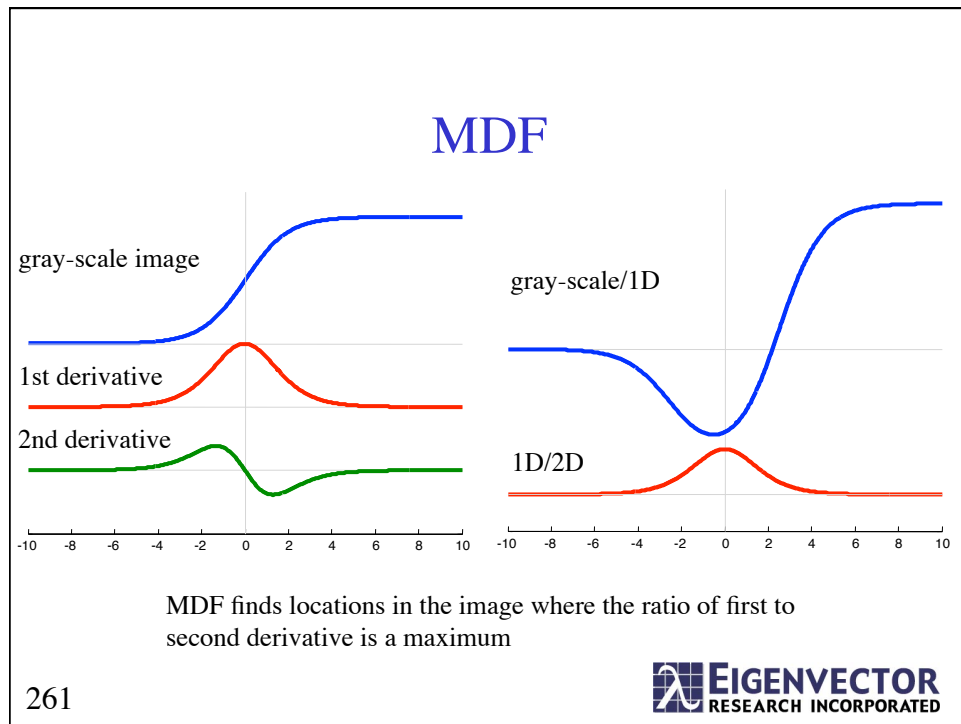
RGB images of PCA w/ GLS weighting scores on 1, 2 and 3. Similar to MAF results. Objective function ~similar, but PCA scores and loadings orthogonal.



Maximum Difference Factors (MDF)

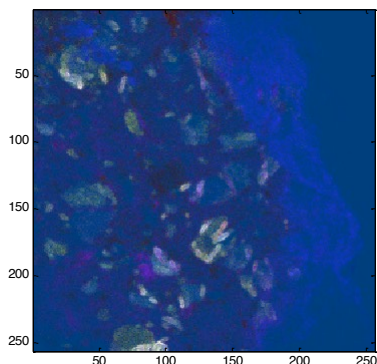
- In MDF the signal covariance corresponds to the first derivative across the spatial dimensions.
 - in MAF the first difference is the clutter
- The clutter corresponds to the second derivative across the spatial dimensions.
- Gives a multivariate analysis estimate of edges in an image.
 - analogous method available for GLS weighting w/ PCA

260

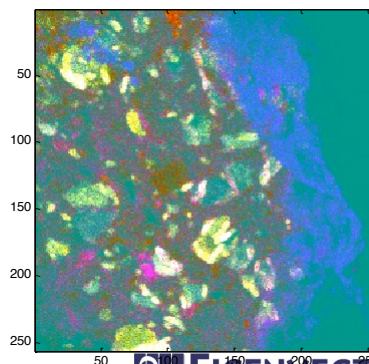


MAF+MDF Results

RGB image of MAF scores 1, 2 and 3 + MDF scores on 1
 $\sqrt{dx^2 + dy^2}$



contrasted image



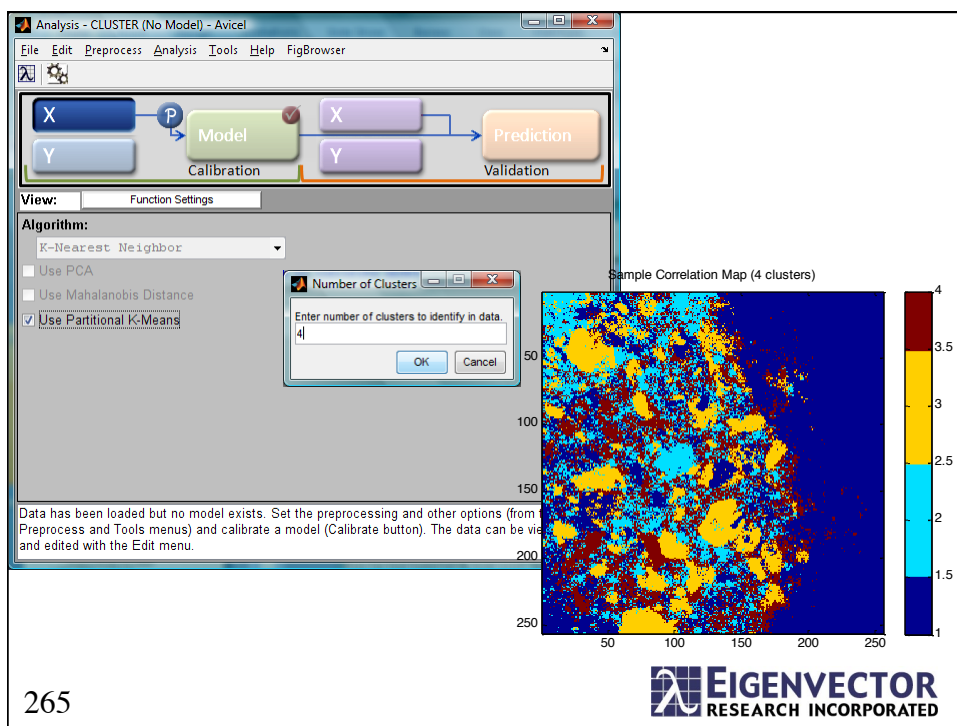
263



cluster analysis

264





MIA ...

- Much more to MIA
 - linked scores plots and density plots
 - interactive exploration of the image(s)
 - image SIMCA and PLS-DA
 - classification
 - curve resolution
 - chemical identification and mapping
 - image statistical process control (ISPC) for multivariate statistical process control (MSPC)
 - ...

266

Outline

- Introduction
- Advanced Preprocessing
- Multivariate image analysis
- Multi-way Analysis
- Summary

©Copyright 2008-2010
Eigenvector Research, Inc.
No part of this material may be photocopied or
reproduced in any form without prior written
consent from Eigenvector Research, Inc.



Definition of Order

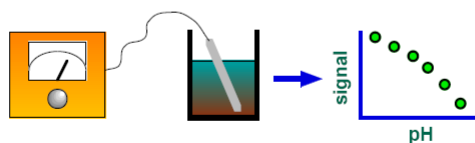
- The **order** of a device is equal to the **dimension** (number of modes) of the data it produces for each sample:
- a **single datum** per sample → **zero order**
- a **vector** (first order tensor) per sample → **first order**
- a **matrix** (second order tensor) sample → **second order**
- **Multi-way analysis** is concerned with data with three or more modes

268



Zero Order Instrument

- The most basic instruments are zero order devices
 - produce a single datum per sample
 - pH, temperature, absorbance at a single channel
 - no way to detect errors or interferences

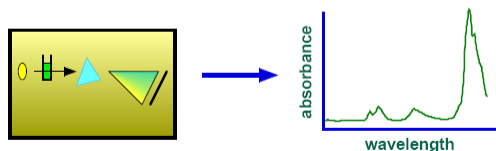


269

 **EIGENVECTOR**
RESEARCH INCORPORATED

First Order Instrument

- Many analytical instruments are first order
 - produce a vector for each sample
 - spectroscopy, LC, GC, sensor arrays
 - the presence of interferences can be detected but not corrected

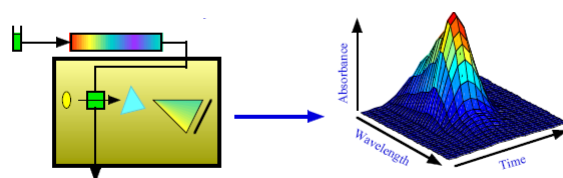


270

 **EIGENVECTOR**
RESEARCH INCORPORATED

Second Order Instrument

- Many analytical instruments are second order
 - produce a matrix for each sample
 - separation followed by spectroscopy, GC-MS, LC-UV
 - interferences can be detected and accounted for

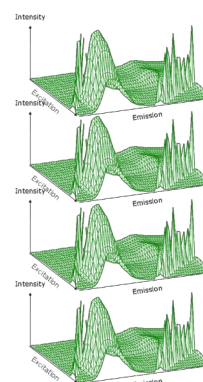
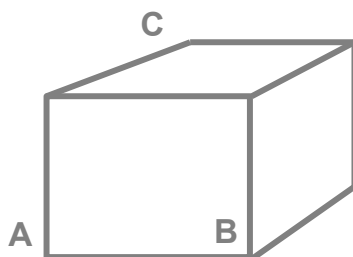


271

EIGENVECTOR
RESEARCH INCORPORATED

Three-way data

- A set of 'equivalent' two-way matrices obtained at different occasions
- Data measured as a function of three 'things' (three different modulations)
 - E.g. samples, variables, times
- x_{ij} is a matrix element and x_{ijk} is a three-way element



272

EIGENVECTOR
RESEARCH INCORPORATED

Examples

- Sensory analysis
 - Score as a function of (Food sample, Judge, Attribute)
- Process analysis
 - Measurement as a function of (Batch, Variable, time)
 - Measurement as a function of (Variable, Lag, Location)
- Image analysis
 - Pixelvalue as a function of (Sample, Image pixel, Variable)
- Experimental design
 - Response as a function of (factor 1, factor2, factor3,...)
- Spectroscopy
 - Intensity as a function of (Wavelength, Retention, Sample, Time, Location , Treatment)
- Environmental analysis
 - Measurement as a function of (Location, Time, Variable)
- Chromatography
 - Measurement as a function of (Sample, Retention time, Variable)

273



Multi-way Algorithms

- Multi-way PCA (weakly multi-way)
- Generalized Rank Annihilation (GRAM)
- Tri-Linear Decomposition (GRAM)
- PARallel FActor Analysis (**PARAFAC**)
- Tucker

274



PARallel FACtor analysis

PARAFAC invented in 1970 by Harshman and independently by Carroll & Chang under the name CANDECOMP. Based on a principle of parallel proportional profiles suggested in 1944 by Cattell

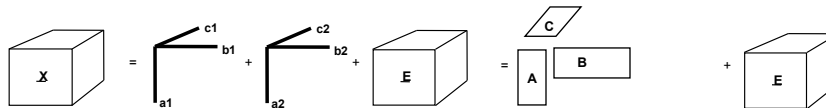
- PCA - bilinear model,

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

*R. A. Harshman. *UCLA working papers in phonetics* 16:1-84, 1970.
*J. D. Carroll and J. Chang. *Psychometrika* 35:283-319, 1970.
*R. B. Cattell. *Psychometrika* 9:267-283, 1944.

- PARAFAC - trilinear model,

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



275



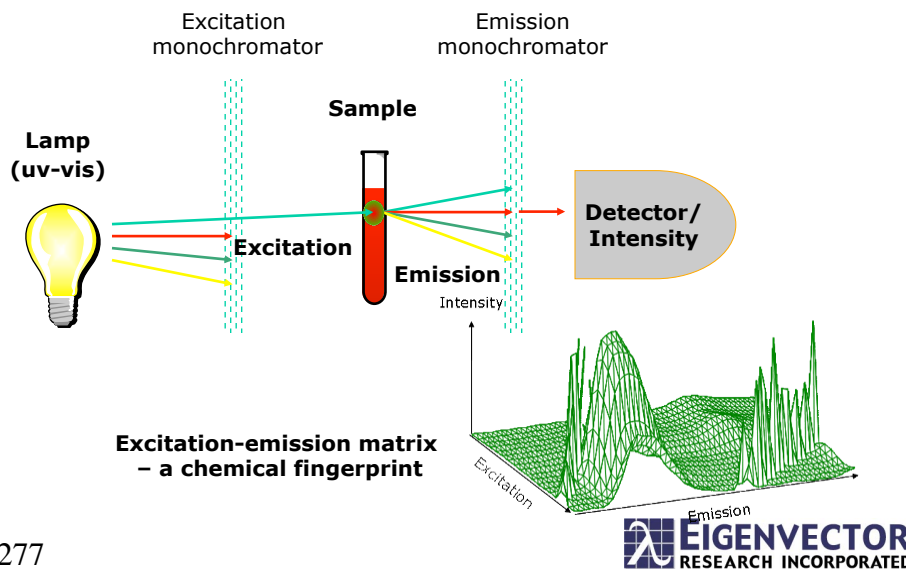
Example: Excitation-Emission Fluorescence

- Use EEM measurements for quantification
 - Measure pure response for target: TRP
 - Measure response of a mixture that includes target + interferences
- Spectra are highly overlapped in both modes
- Example of second order calibration
 - the goal is to detect one analyte in the presence of unknown varying interferences using the entire EEM response
 - use PARAFAC
 - quantification in the presence of previously unseen interferences

276



Excitation Emission Fluorescence Spectroscopy

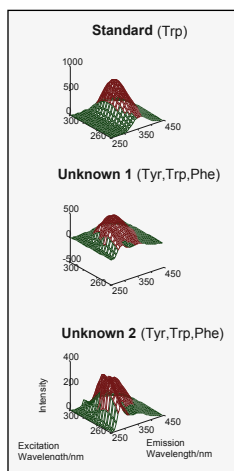


277

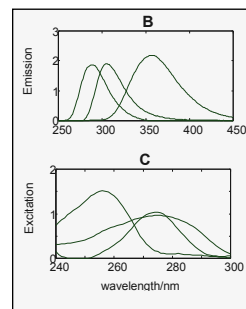
Who Needs Regression?

Calibration set: *One* sample with *one* analyte (2.67 μ M Trp)

Test set: Two samples with three analytes each (Trp, Tyr, Phe)



$$\mathbf{A} = \begin{bmatrix} 2.67 & 0 & 0 \\ 1.52 & 1.30 & 1.29 \\ .86 & 1.12 & 1.06 \end{bmatrix}, \quad \text{reference} = \begin{bmatrix} 2.67 \\ 1.58 \\ .88 \end{bmatrix}$$



278

EigenVECTOR
RESEARCH INCORPORATED

Second Order Calibration

- The PARAFAC model estimated
 - amount of target in the test set in the presence of interferences not seen in the calibration set!
 - this is not possible with PLS
 - estimates of the response in both modes
 - allows potential library searches
- This has enormous potential for environmental sensing and MSPC
 - Smilde, A., Bro, R., and Geladi, P., “Multi-way Analysis with Applications in the Chemical Sciences”, John Wiley & Sons, New York, NY (2004).

279



Multi-way Analysis ...

- Curve resolution
 - PARAFAC needs less futsing than two-way MCR
- MSPC, images, DECRA, ...
- > 3 Modes
 - GCxGCxMS, sensor fusion, ...

Smilde, A., Bro, R. and Geladi, P., “Multi-way Analysis with Applications in the Chemical Sciences”, John Wiley & Sons, New York, NY (2004).

280



Summary

- Data analysis requires knowledge of
 - the system, physics, chemistry and math → ~~black box~~
- Advanced Preprocessing
 - uses knowledge of the clutter (GLS, ELS, etc.)
- Multivariate image analysis
 - spatial and spectral information
- Multi-way Analysis
 - measurements a function of multiple modulations giving ≥ 3 modes per sample → quant w/ unknown interferences

©Copyright 2008-2011
Eigenvector Research, Inc.
No part of this material may be photocopied or
reproduced in any form without prior written
consent from Eigenvector Research, Inc.



Section Definitions 1/2

- **Multivariate image analysis (MIA)**: Analysis of multivariate images (for many variables → hyperspectral image analysis).
- **Multivariate image (MI)**: A data array of dimension three (or more) where the first two dimensions are spatial and the last dimension(s) is a function of another variable.
- **Maximum/minimum noise fractions (MNF)**: Algorithm that maximizes capture of signal relative to a clutter covariance resulting in a generalized eigenvector problem.
- **Maximum autocorrelation factors (MAF)**: MNF with the clutter covariance corresponding to the first difference of image pixels.
- **Maximum autocorrelation factors (MAF)**: MNF with the clutter covariance corresponding to the first difference of image pixels.
- **Maximum difference factors (MDF)**: MNF with the signal corresponding to the covariance of the first spatial derivative and clutter covariance corresponding to the second spatial derivative.
- **Order**: is the dimension of the data produced per sample.
- **Dimension**: number of modes.



Section Definitions 2/2

- **Multi-way analysis:** Analysis of data with three or more modes.
- **Multivariate image (MI):** A data array of dimension three (or more) where the first
- **PARallel FACto Analysis (PARAFAC):** model for multi-way analysis of ≥ 3 mode data.