

# Orthogonalization Approaches for Data Preprocessing with Pharmaceutical, Petrochemical and Remote Sensing Applications

Barry M. Wise, Jeremy M. Shaver and  
Neal B. Gallagher

Eigenvector Research, Inc.

# Abstract

Over the past dozen years, a number of powerful spectral analysis methods have been published which make use of orthogonalization (*i.e.* projection followed by weighted subtraction) of interferences or "clutter." These filtering methods provide a means to mitigate the effect of interferences arising from background chemical or physical species, instrumental artifacts, systematic sampling errors and instrument or system drift. They have been used very effectively with complex biological systems, remote sensing applications, chemical process monitoring and calibration transfer problems.

This class of methods includes Orthogonal Partial Least Squares (O-PLS), External Parameter Orthogonalization (EPO), Dynamic Orthogonal Projection (DOP), Orthogonal Signal Correction (OSC), Constrained Principal Spectral Analysis (CPSA), Generalized Least Squares Weighting (GLSW), and Science Based Calibration (SBC) among others. All are based on the orthogonalization premise and each touts a unique ability to improve model performance, robustness, and/or interpretability.

Some relationships between these methods are noted, along with ties to older work. Examples are given of the use of the methods in calibration and classification problems in pharmaceutical, petrochemical and remote sensing applications.

# What is an Orthogonalization Filter?

- Removes spectral patterns from data which are "interfering" with signal of interest
- The interfering species are historically called "clutter" (backgrounds, noise, interferents)
- Filters return spectra with features "removed"
- Weighted subtraction of one or more vectors
- "Soft" orthogonalization is deweighting but not outright complete subtraction

# Some Examples Using Orthogonalization Filters (by Eigenvector)

- *In vivo* Tissue identification with NIR probe
- Cancer detection using *in vivo* fluorescence
- Identification of **arthlesclerosis** in artery walls using NIR
- Determination of **hydroxide concentration** in high-concentration aqueous ion solutions using Raman spectroscopy
- Identification of chemical species in **remote sensing**

# SOME Orthogonalization Filters

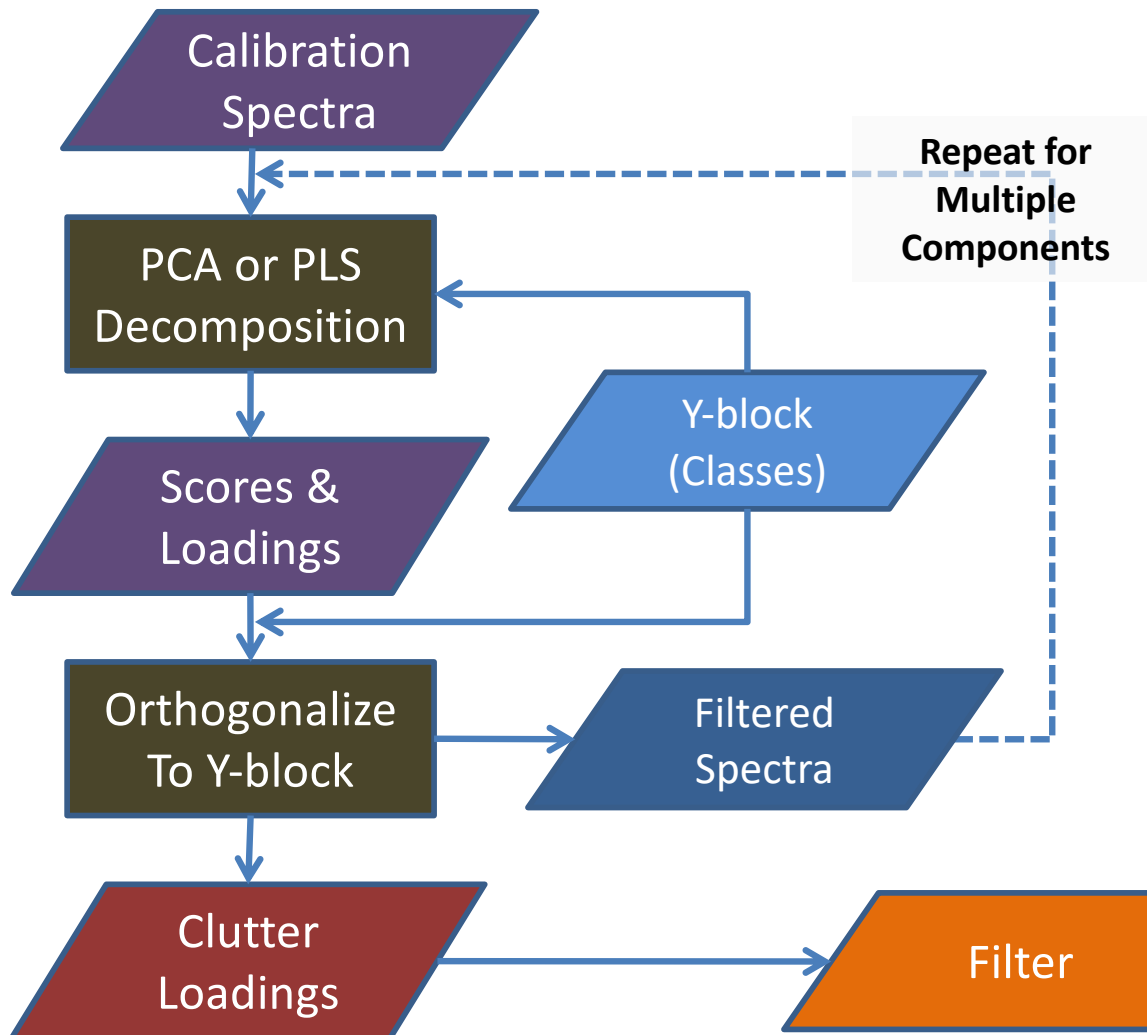
## Method 1: Orthogonalization of Model

- **OSC – Orthogonal Signal Correction** (Wold et. al. 1998)
- **OPLS – Orthogonal PLS** (Trygg, Wold 2002 , patented)
- **MOSC – Modified OSC** (POSC - Feudale, Tan, S. Brown 2003)
- **CPSA - Constrained Principal Spectral Analysis**  
(J. Brown 1990 , patented)
- **EPO – External Parameter Orthogonalization**  
(Roger, Chauchard, Bellon-Maurel 2003)
- **GLS – Generalized Least Squares**  
(Aitken 1935, Martens et. al. 2003)
- **SBC – Science Based Calibration**  
(Marbach 2005, patented (?))

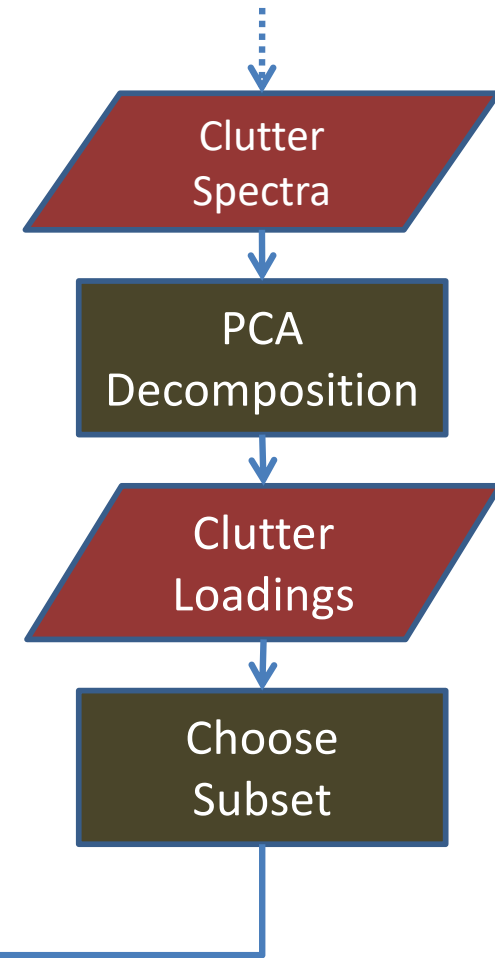
## Method 2: Pre-selection of "clutter"

# Two General Approaches

## Method 1: Orthogonalization of Model



## Method 2: Pre-selection of "clutter"



# Orthogonal Signal Correction (OSC)

- Introduced by Wold in 1998
  - “OSC paper not the clearest thing” – Johan Trygg, June 9, 2011
- OSC objective function

$$\max_{\mathbf{r}} [\text{var}(\mathbf{t}) \mid \mathbf{t} = \mathbf{X}\mathbf{r} \wedge \mathbf{t} \perp \mathbf{y}]$$

$$\mathbf{p} = \mathbf{X}^T \mathbf{t} (\mathbf{t}^T \mathbf{t})^{-1}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{r}^T \mathbf{X} \mathbf{p}$$

# OSC Issues

- To the extent the objective function is optimized OSC doesn't work
  - Only works if you don't try too hard!
- Many algorithms (at least 5) with various problems
  - Factors not orthogonal to  $\mathbf{y}$
  - Factors don't capture maximum variance in  $\mathbf{X}$
  - Filtered  $\mathbf{X}$  not in same subspace as original  $\mathbf{X}$
- Often implemented prior to cross validation—totally misleading!



# O-PLS

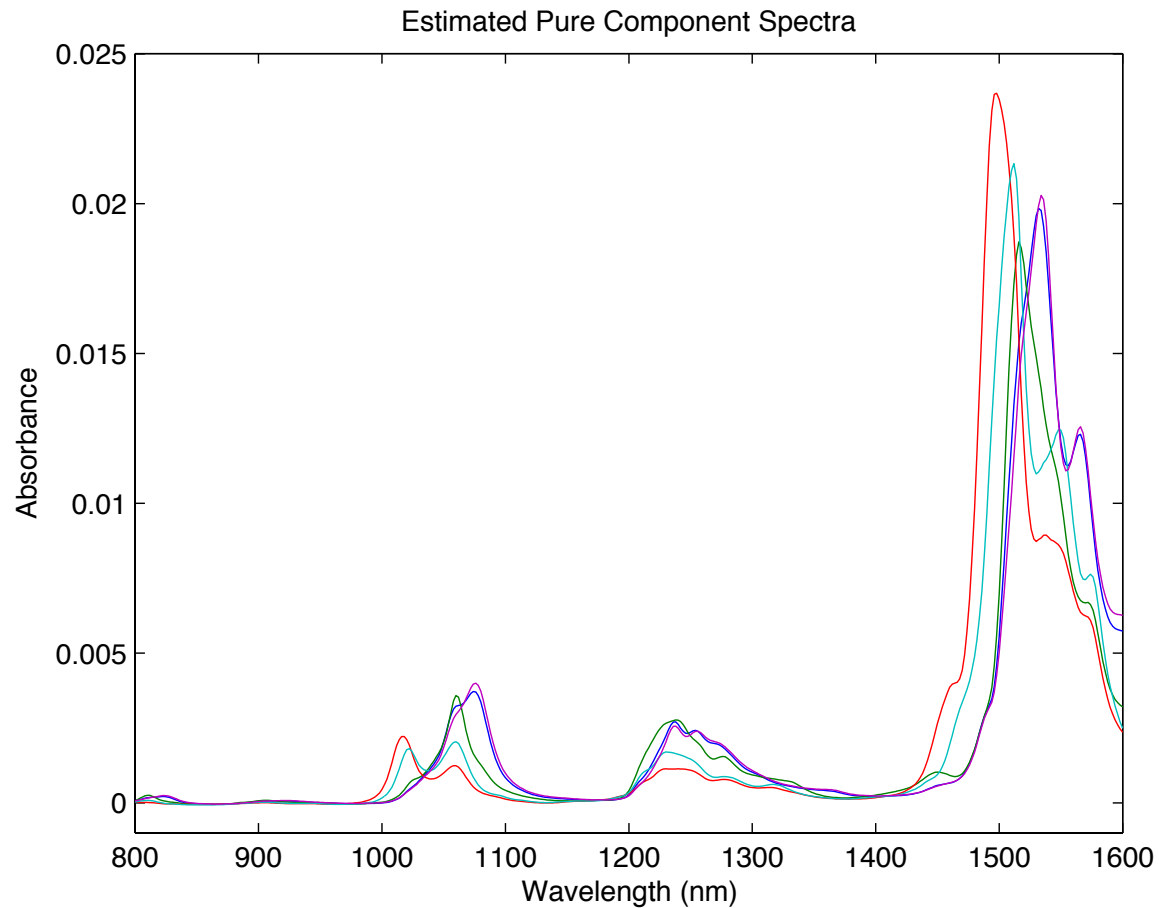
- Originally formulated as sequential algorithm (NIPALS based)
- Since shown to be obtainable from post-processing conventional PLS model
- Does not improve prediction
- Claim is that model is more interpretable

E.K. Kemsley and H.S. Tapp, "OPLS filtered data can be obtained directly from non-orthogonalized PLS1," *J. Chemo*, **23**, 263-264, 2009

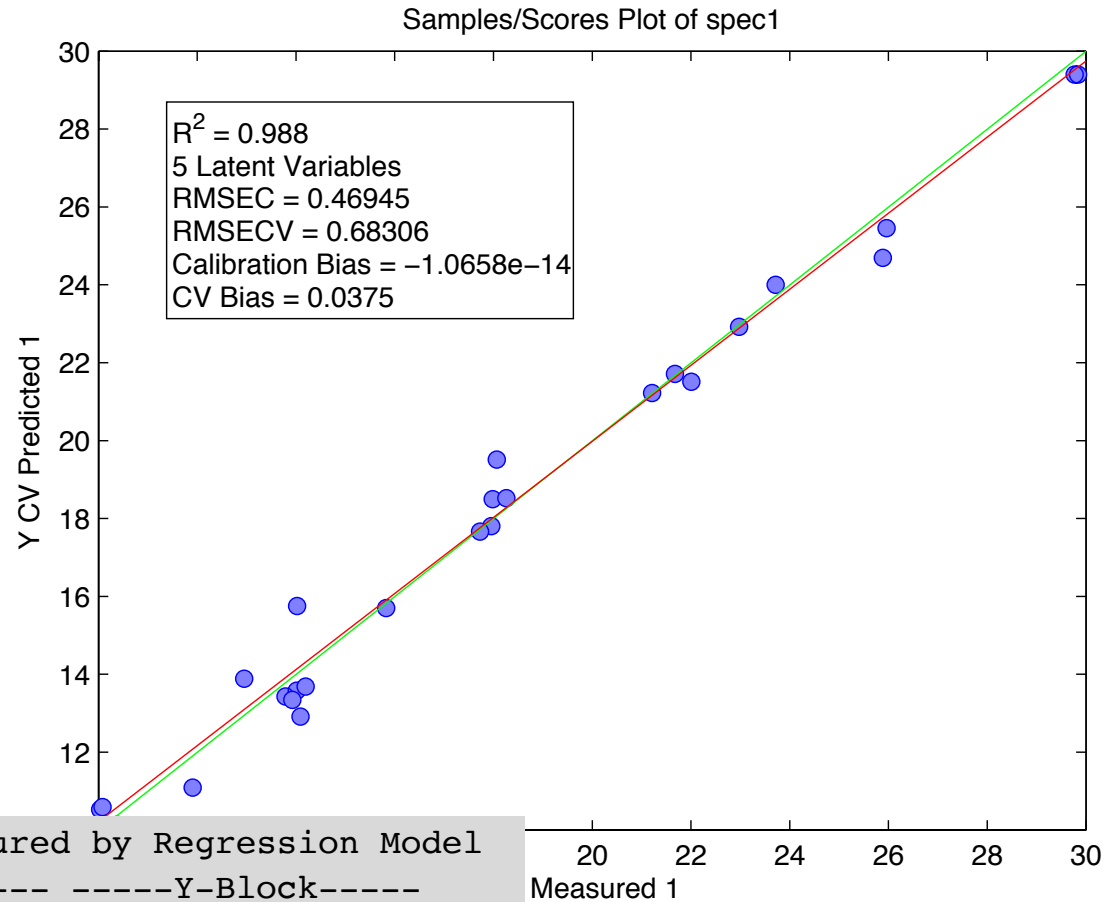
R. Ergon, "PLS post-processing by similarity transformation (PLS+ST): a simple alternative to OPLS," *J. Chemo*, **19**, 1-4, 2005

J. Trygg and S. Wold, "Orthogonal Projections to Latent Structures (O-PLS)," *J. Chemo*, **16**, 119-128, 2002.

# NIR of Pseudo-gasoline Samples



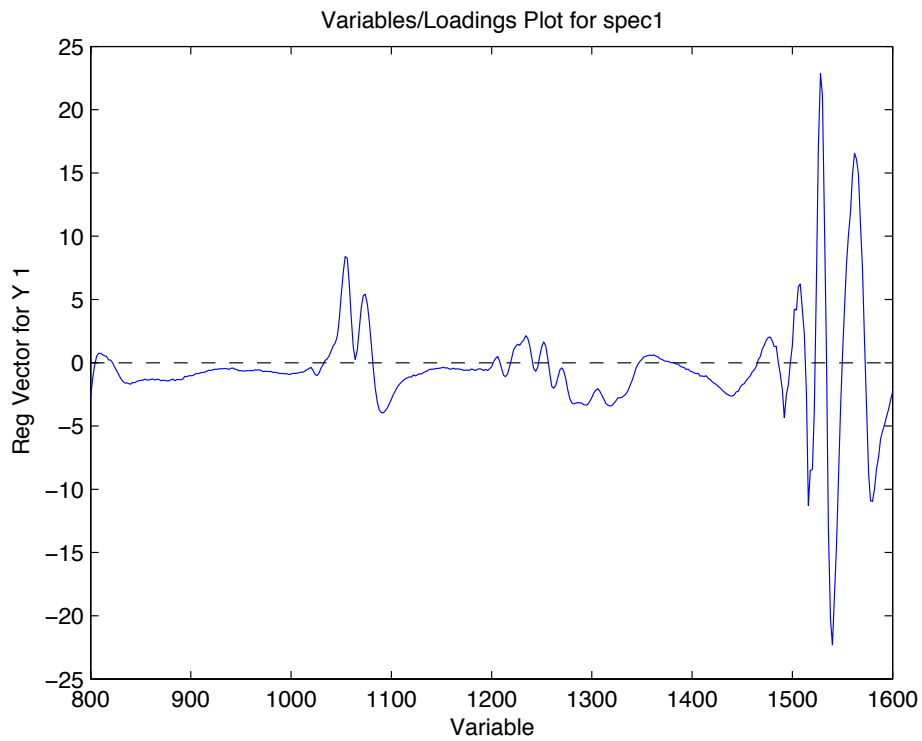
# PLS Model on Component 1



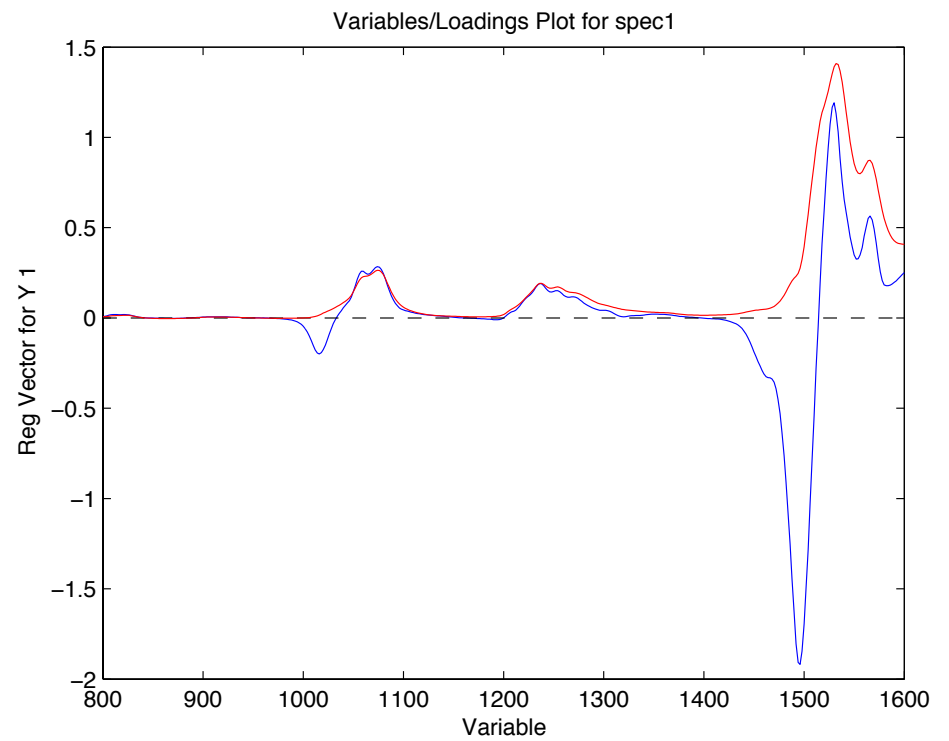
Percent Variance Captured by Regression Model

Comp	-----X-Block-----		-----Y-Block-----	
	This	Total	This	Total
1	91.17	91.17	8.36	8.36
2	7.40	98.57	7.19	15.55
3	0.93	99.50	32.81	48.36
4	0.46	99.96	26.18	74.54
5	0.02	99.98	24.90	99.44

# Regular and O-PLS Filtered Regression Vectors



Regular



O-PLS Filtered

# Interpretation

- Better in previous example, but partly because we know what spectra should look like
- What if problem has discrete variables with signal that could be positive, negative or zero?
- Much harder! (see “*On the Interpretability of O-PLS Models*”)
- Working on developing better understanding of when it will work and when it won't

# Orthogonalize Model

Analysis - PLS 6 LVs - m5spec, propvals

File Edit Preprocess Analysis Tools Help FigBrowser

Tools

- ✓ Cross-Validation
- Orthogonalize Model**
- Show Details
- Report Writer
- Test Model Robustness
- Permutation Test
- Correlation Map
- Estimate Factor SNR
- View Cache
- Toolbar

View: SSQ Table

Number LVs: 6 Auto Select

Latent Variable	Percent Variance X-Block		Percent Variance Y-Block	
	LV	Cum	LV	Cum
1	99.08	99.08	39.05	39.05
2	0.76	99.84	19.26	58.31
3	0.06	99.90	23.49	81.79
4	0.03	99.93	14.25	96.04
5	0.03	99.96	2.24	98.28
6	0.01	99.98	1.00	99.28 ← Suggested
7	0.01	99.98	0.31	99.59
8	0.01	99.99	0.09	99.68
9	0.00	99.99	0.16	99.83
10	0.00	100.00	0.02	99.85
11	0.00	100.00	0.09	99.94
12	0.00	100.00	0.02	99.96
13	0.00	100.00	0.01	99.97
14	0.00	100.00	0.00	99.97
15	0.00	100.00	0.00	99.98
16	0.00	100.00	0.01	99.98



Analysis - PLS 6 LVs - m5spec, propvals

File Edit Preprocess Analysis Tools Help FigBrowser

Tools

- ✓ Cross-Validation
- Orthogonalize Model**
- Show Details
- Report Writer
- Test Model Robustness
- Permutation Test
- Correlation Map
- Estimate Factor SNR
- View Cache
- Toolbar

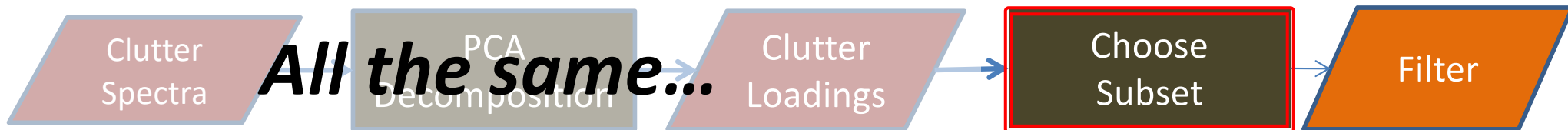
View: SSQ Table iPLS Variable Selection

Number LVs: 6 Auto Select

Latent Variable	Percent Variance Captured by Model X-Block		Percent Variance Captured by Model Y-Block	
	LV	Cum	LV	Cum
1	96.86	96.86	99.28	99.28
2	2.73	99.60	0.00	99.28
3	0.20	99.80	0.00	99.28
4	0.10	99.90	0.00	99.28
5	0.06	99.95	0.00	99.28
6	0.02	99.98	0.00	99.28 ← Suggested
7	0.01	99.98	0.31	99.59
8	0.01	99.99	0.09	99.68
9	0.00	99.99	0.16	99.83
10	0.00	100.00	0.02	99.85
11	0.00	100.00	0.09	99.94
12	0.00	100.00	0.02	99.96
13	0.00	100.00	0.01	99.97
14	0.00	100.00	0.00	99.97
15	0.00	100.00	0.00	99.98
16	0.00	100.00	0.01	99.98

# Pre-selection Methods...

- CPSA - Constrained Principal Spectral Analysis  
(J. Brown 1990, patented)
  - EPO – External Parameter Orthogonalization  
(Roger, Chauchard, Bellon-Maurel 2003)
- Identical  
• Choose # of PCs



- GLS – Generalized Least Squares  
(Aitken 1935)
  - SBC – Science Based Calibration  
(Marbach 2005, patented (?))
- Quite similar  
• Down-weight by scale of eigenvalues

# Clutter Covariance

Clutter source 1

Clutter source 2

$$\mathbf{X}_c = (\mathbf{X}_{1,c} - \bar{\mathbf{x}}_{1,c}) + (\mathbf{X}_{2,c} - \bar{\mathbf{x}}_{2,c}) + \dots$$

$$\mathbf{C} = \frac{\mathbf{X}_c^T \mathbf{X}_c}{N - 1}$$



# Covariance to GLS Weighting Matrix

$$\mathbf{C} = \mathbf{V}\mathbf{S}^2\mathbf{V}^T$$

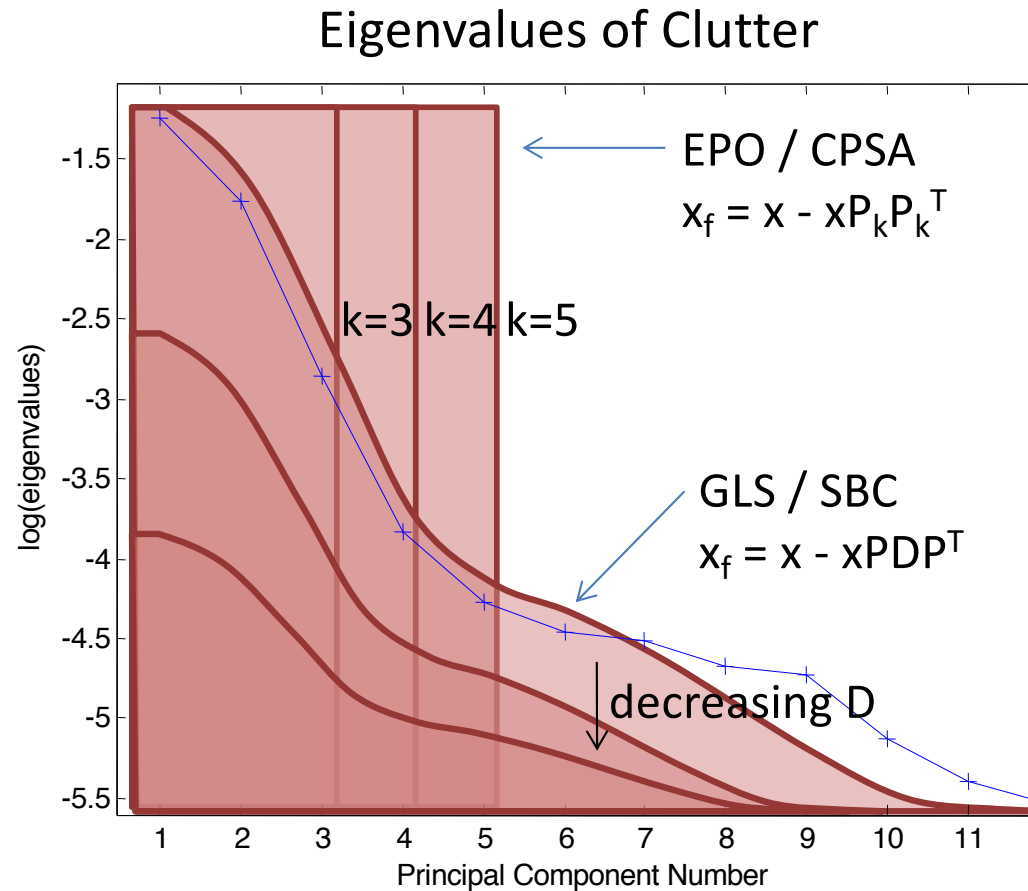
$$\mathbf{G} = \mathbf{V}\mathbf{D}^{-1}\mathbf{V}^T$$

with

$$d_{i,i}^{-1} = \frac{1}{\sqrt{\frac{s_{i,i}^2}{g^2} + 1}}$$

Large  $g \rightarrow 1$ ,  
dimension  
unaffected  
Small  $g \rightarrow 0$ ,  
dimension eliminated

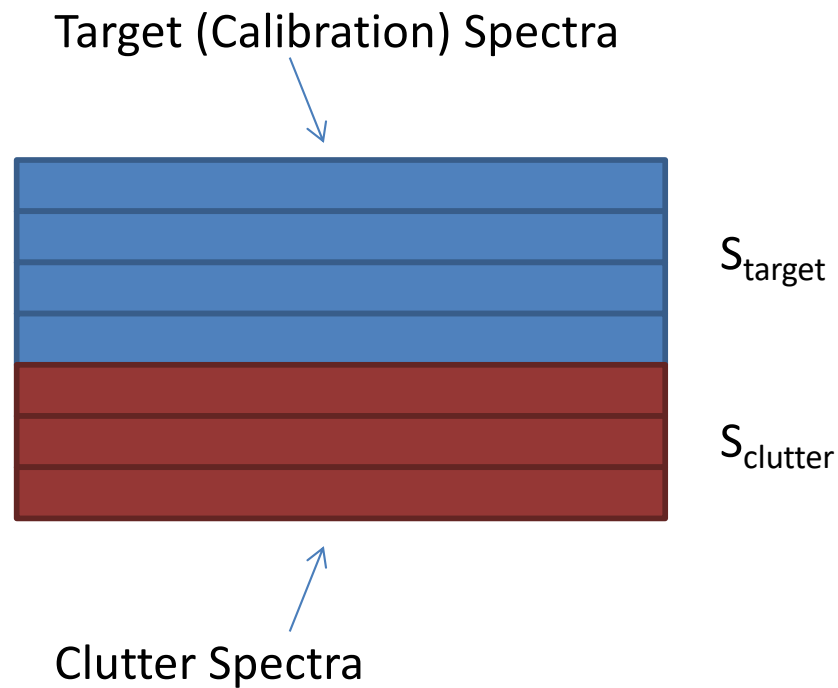
# Choosing Components



One adjustable parameter in each method

# Other Similar Pre-selection Filters...

- Extended Mixture Model (Extended Least Squares) orthogonal filtering for Classical Least Squares (CLS) models!



$$c = xS(S^T S)^{-1}$$

Pseudo-inverse is an orthogonalization!

Equivalent to full-rank EPO / CPSA model

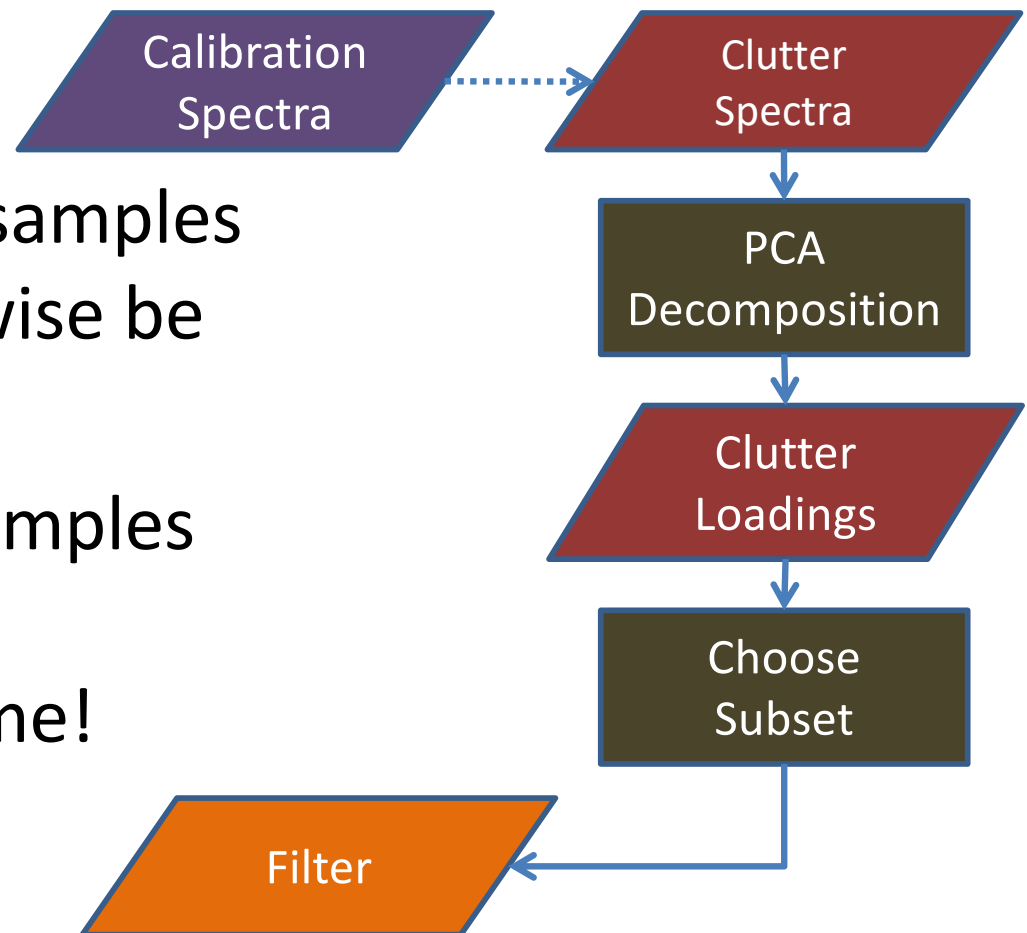
# Pre-selecting Clutter

How to get clutter?

Look at differences in samples which should otherwise be the same.

In classification – all samples within a class should nominally be the same!

Use Calibration itself!



# More on How to Get Clutter

- Pure component spectra of known interferences
- Subspace spanned by
  - samples where analyte of interest is not present
  - variation in data that is all of the same class
  - differences between samples where analyte of interest is (nearly) the same, *e.g.*  $y$ -gradient
  - repeat measurement of blanks
- Make it up! *e.g.* polynomial baseline shapes

# Y-gradient Method

- Sort samples by  $y$  (reference) values
- Take differences between adjacent samples
- Weight  $X$ -differences by inverse of difference in  $y$  values
- Deweight by covariance of differences (GLS) or orthogonalize against some number of PCs (EPO, ELS, EMM, PA-CLS)

# Orthogonalization Filters

Filter	Soft/ Hard	Adj. Params	Clutter source	Improves Prediction?
OSC	Hard	# LVs	Part of $\mathbf{X}$ orthogonal to $\mathbf{y}$	No, but reduces models complexity
O-PLS	Hard	# LVs	Part of $\mathbf{X}$ -model space orthogonal to $\mathbf{X}'\mathbf{y}$	No, but improves interpretation
MOSC	Hard	# PCs	Part of $\mathbf{X}$ orthogonal to $\mathbf{y}$	Maybe
CPSA	Hard	# PCs	A priori, includes pathlength adj.	Yes
EPO	Hard	# PCs	Classes, $\mathbf{y}$ -gradient or a priori	Yes
DOP	Hard	# PCs	Synthetic reference samples	Yes
GLS	Soft	Shrinkage $\alpha$	Classes, $\mathbf{y}$ -gradient or a priori	Yes
SBC	Soft	# PCs (20?)	Repeat samples or blanks	Yes
EMM	Hard	None	A priori from known interferents, clutter subspace	Yes, CLS model
ELS	Hard	# PCs	Clutter subspace	Yes
PA-CLS	Hard	None/# PCs	Baseline shapes, residuals	Yes, CLS model
WLS	Soft	Regularization	Noise measurements	Yes

# We think it is useful to use Clutter!

The screenshot shows the FigBrowser software interface. A red arrow points to the 'Clutter' button in the workflow diagram. The workflow diagram illustrates the process: X and Y data are processed through 'Calibration' (using a 'Model') to produce 'Prediction' results. The 'Clutter' button is located between the 'Calibration' and 'Prediction' stages.

View: SSQ Table | iPLS Variable Selection

Number LVs: 2 | Auto Select

Latent Variable	X-Block		Y-Block	
	LV	Cum	LV	Cum
1	78.09	78.09	98.12	98.12
2	10.12	88.20	0.57	98.70
3	1.57	89.77	0.21	98.91
4	0.89	90.67	0.14	99.05
5	0.75	91.41	0.08	99.13
6	0.56	91.98	0.07	99.20
7	0.49	92.46	0.05	99.24
8	0.41	92.87	0.04	99.29
9	0.37	93.24	0.03	99.32
10	0.26	93.50	0.03	99.35
11	0.41	93.91	0.02	99.37

The screenshot shows the Declutter Settings dialog box. The 'Clutter Source' section has 'y-block gradient' selected. The 'Algorithm' section has 'Ignore Means (mean center)' checked and 'GLSW' selected. The 'Declutter Threshold' is set to 0.002. The 'Number of PCs' is set to 1.

Clutter Source: automatic, y-block gradient, x-block classes, external data

Algorithm:  Ignore Means (mean center),  GLSW,  EPO,  EMM / ELS,  None (disable filter)

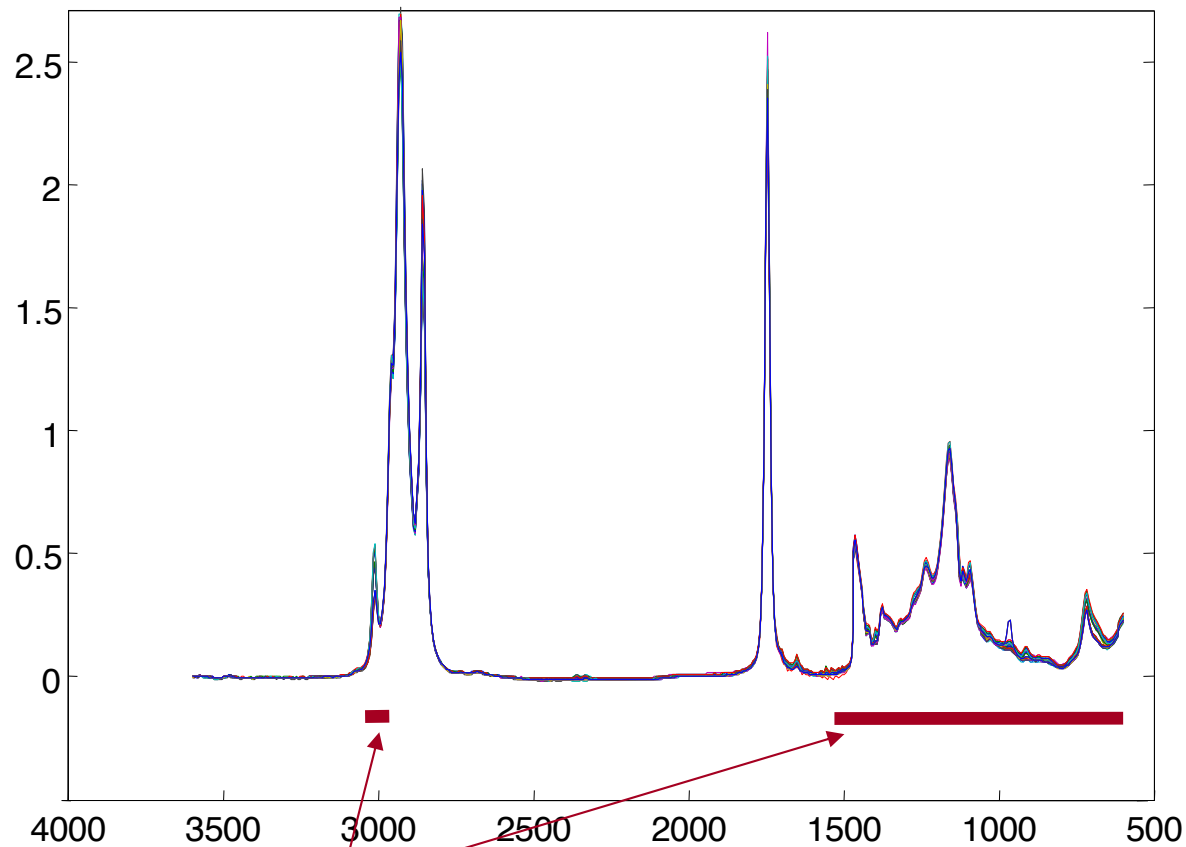
Declutter Threshold: 0.002

Number of PCs: 1



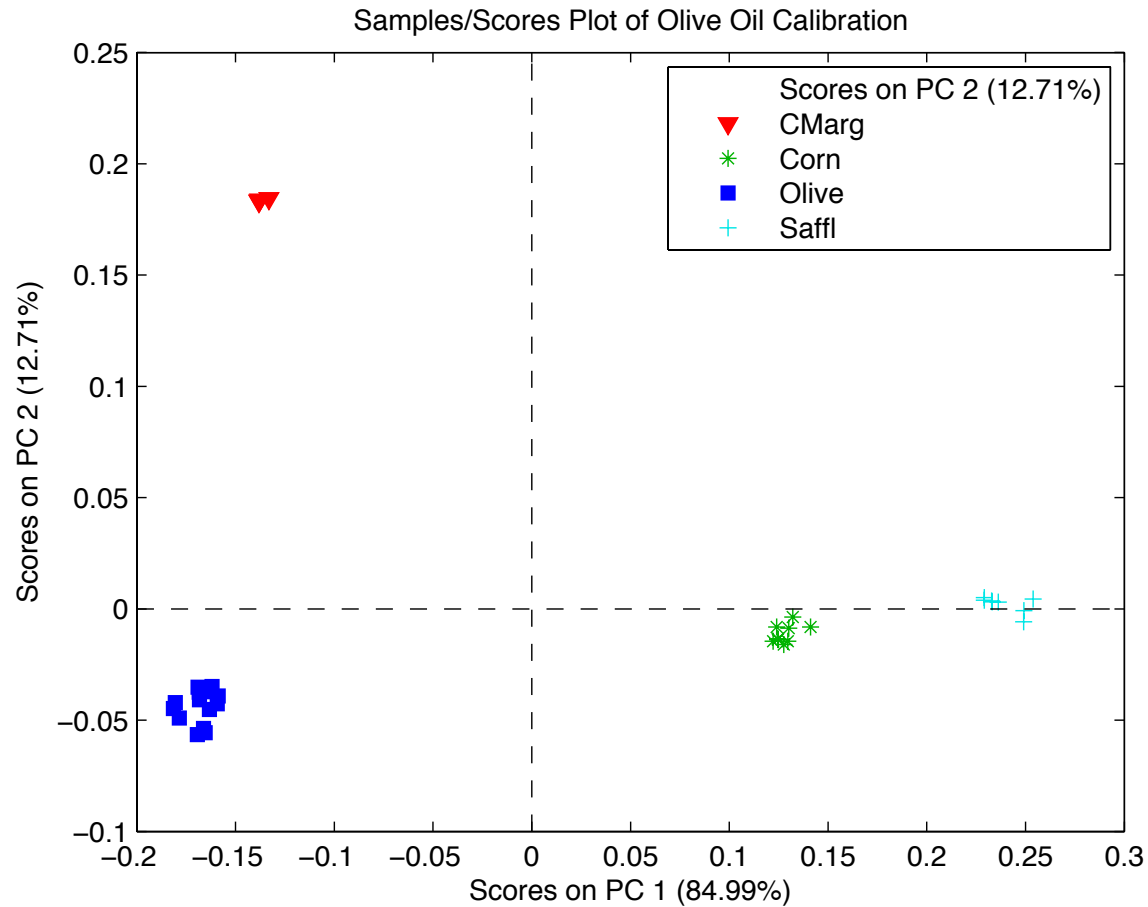
# Example Classification Data

- Mid-IR spectra of food grade oils
- Classify oils, detect adulterated olive oil

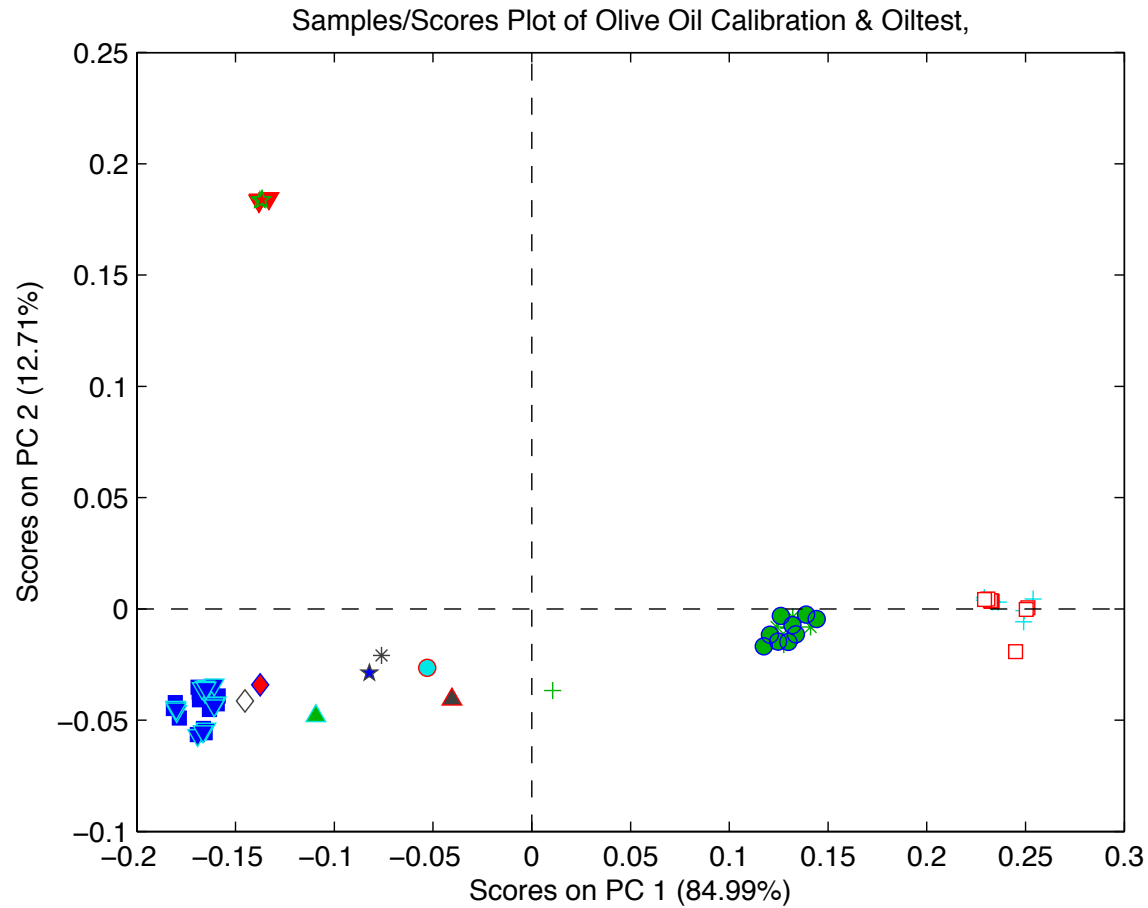


Using these regions only

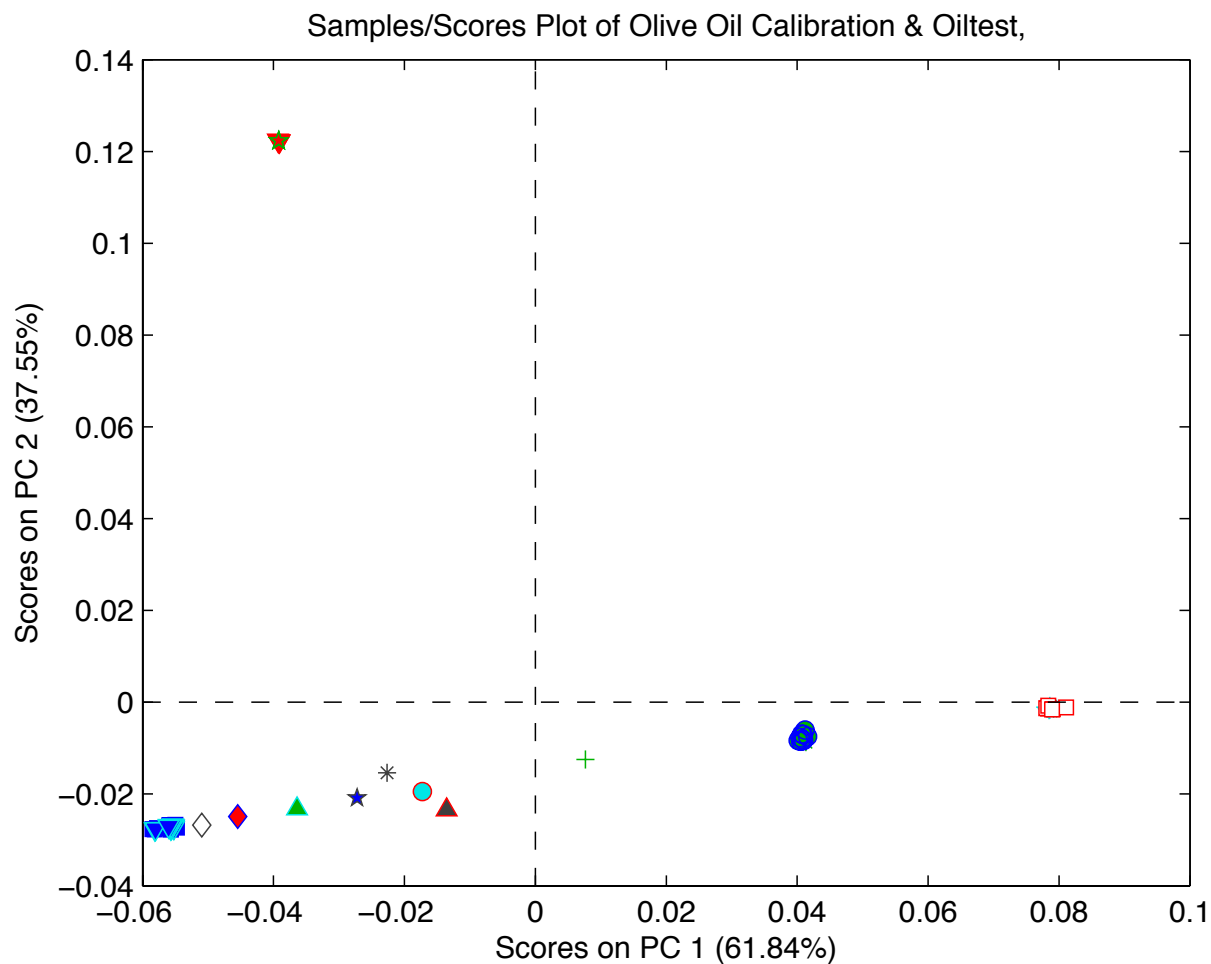
# Calibration with MSC



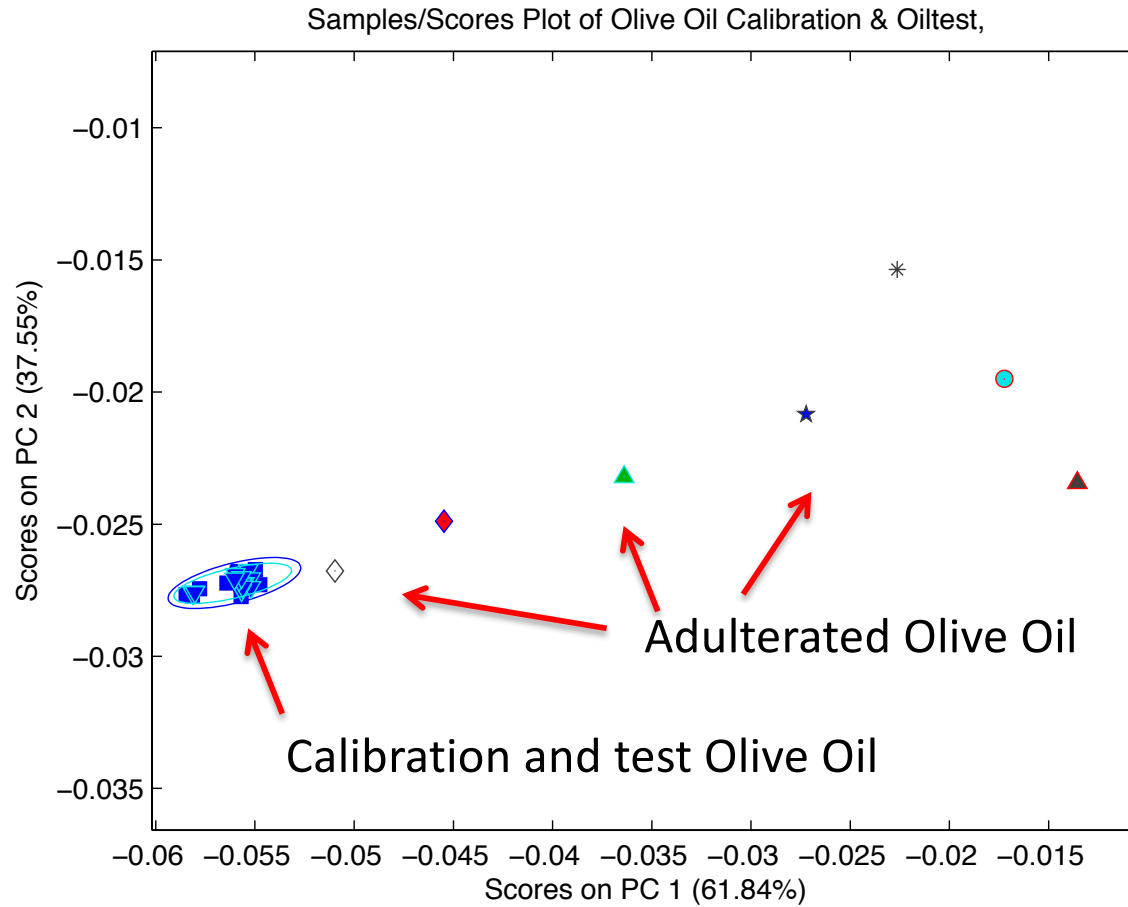
# Cal and Test with MSC



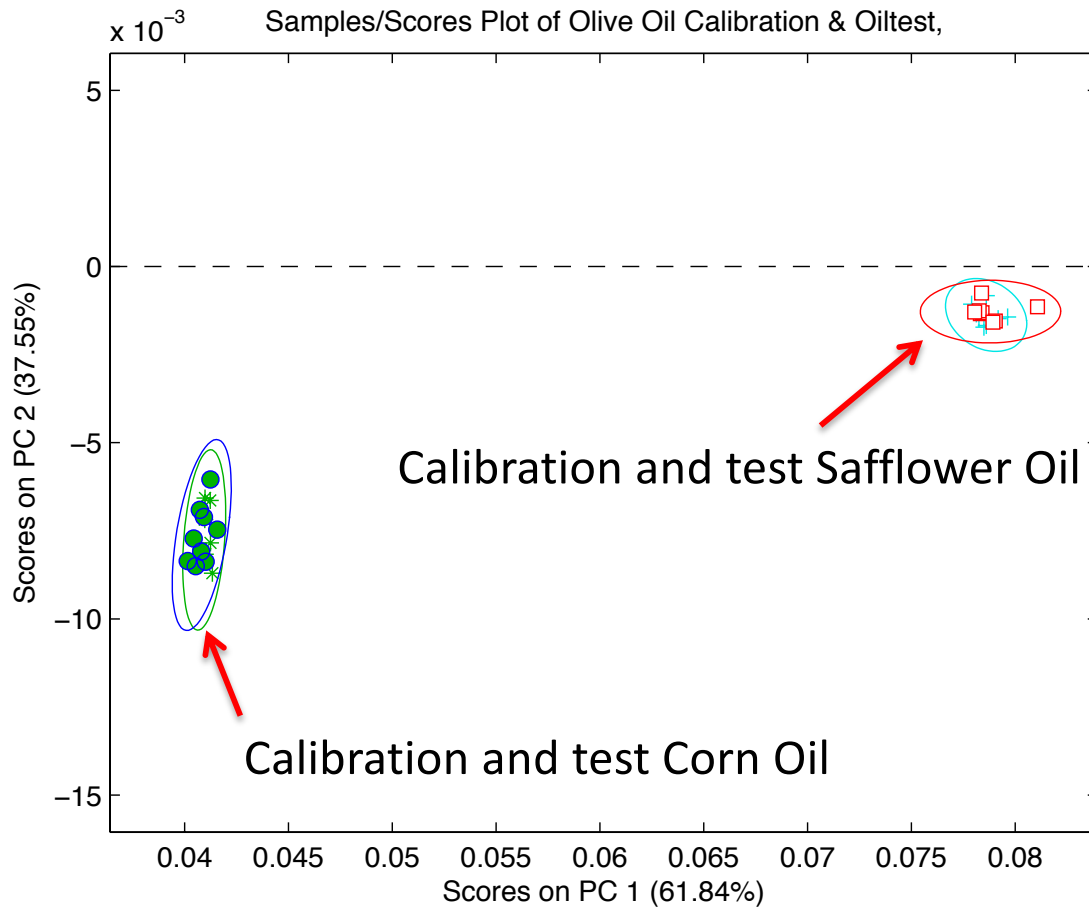
# With MSC and GLS



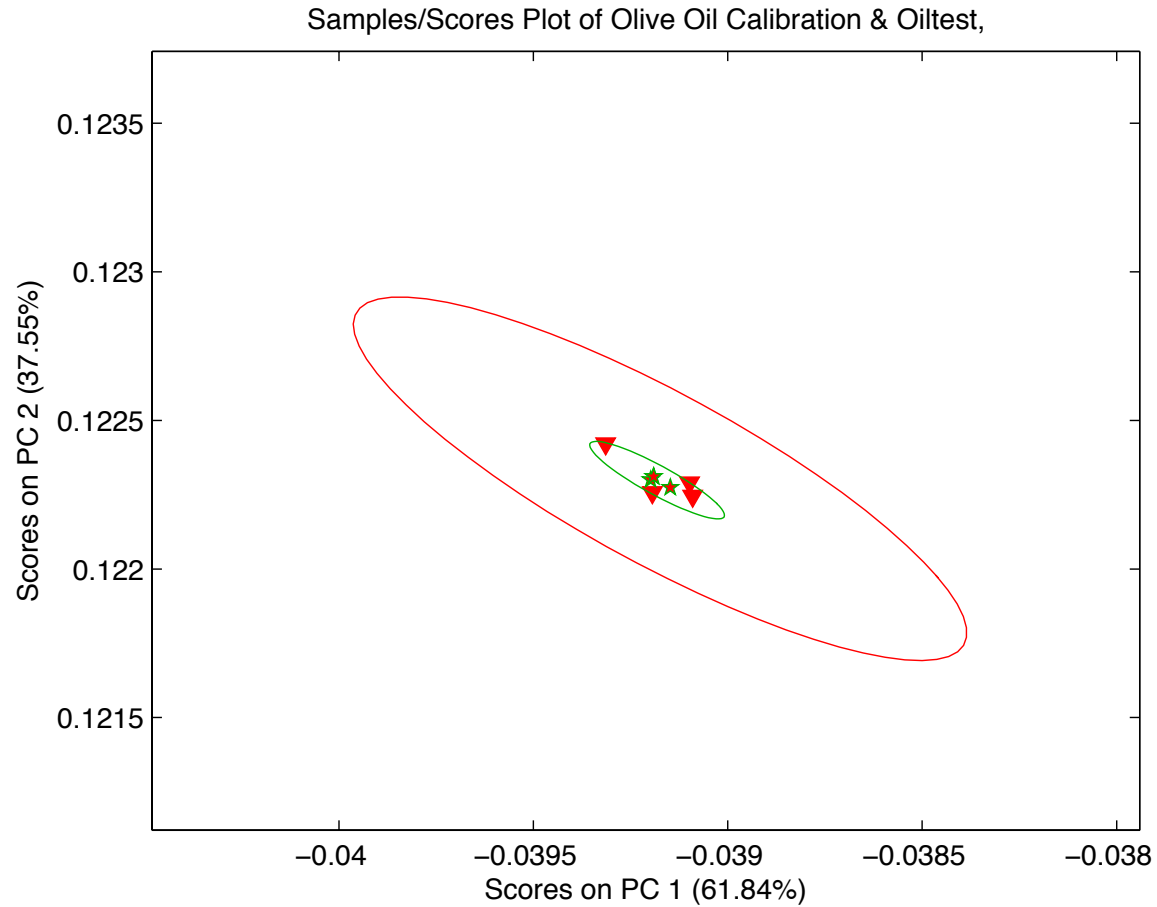
# Zoom on Olive Oil



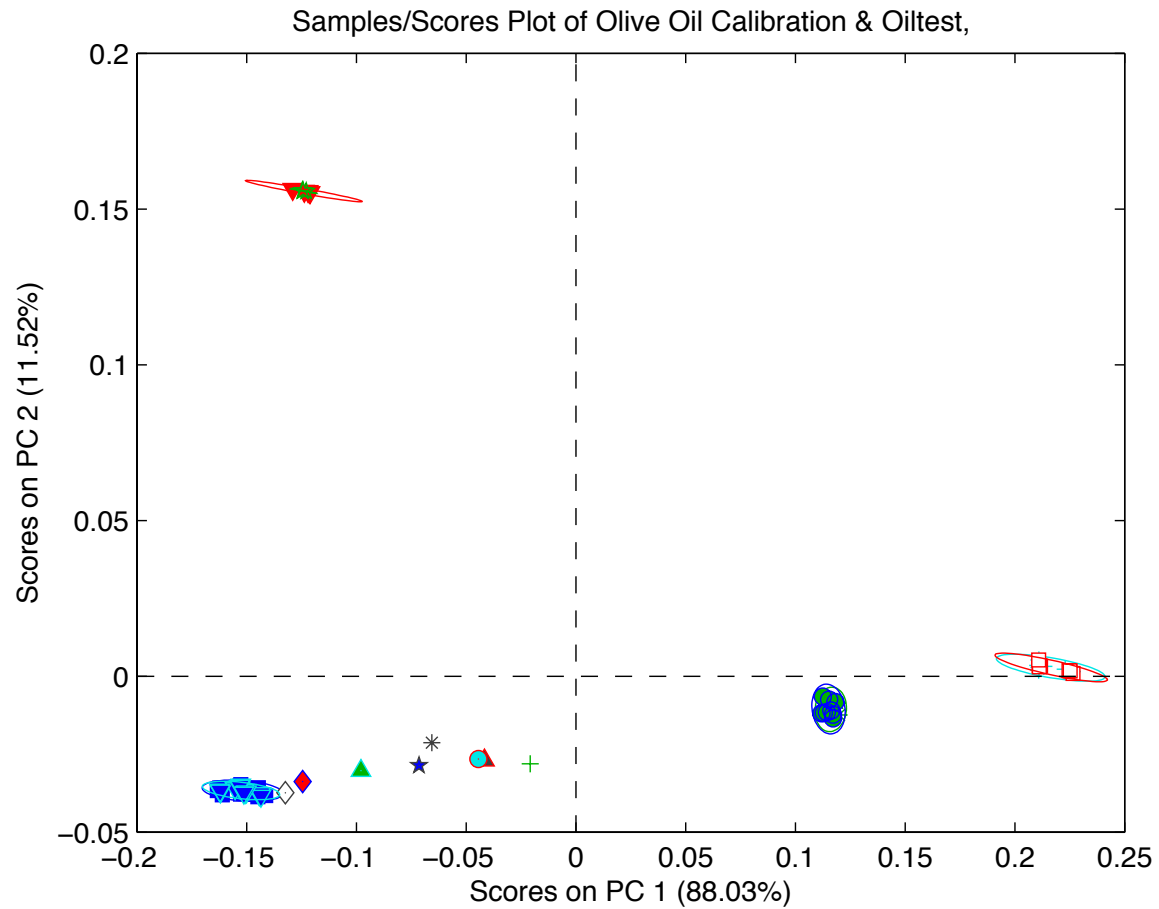
# Zoom on Corn and Safflower Oil



# Zoom on Corn Margarine



# With MSC and EPO





# Indian Pines Data

- Classic image data set used in many publications
- Crop area near West Lafayette, Indiana
- Ground truth identified 16 known crop areas
- Data from AVIRIS: Airborne Visible/Infrared Imaging Spectrometer
- 220 channels, 400-2500nm

# Indian Pines Image

Image of Scores on PC 1 (72.48%) & Scores on PC 3 (1.73%)

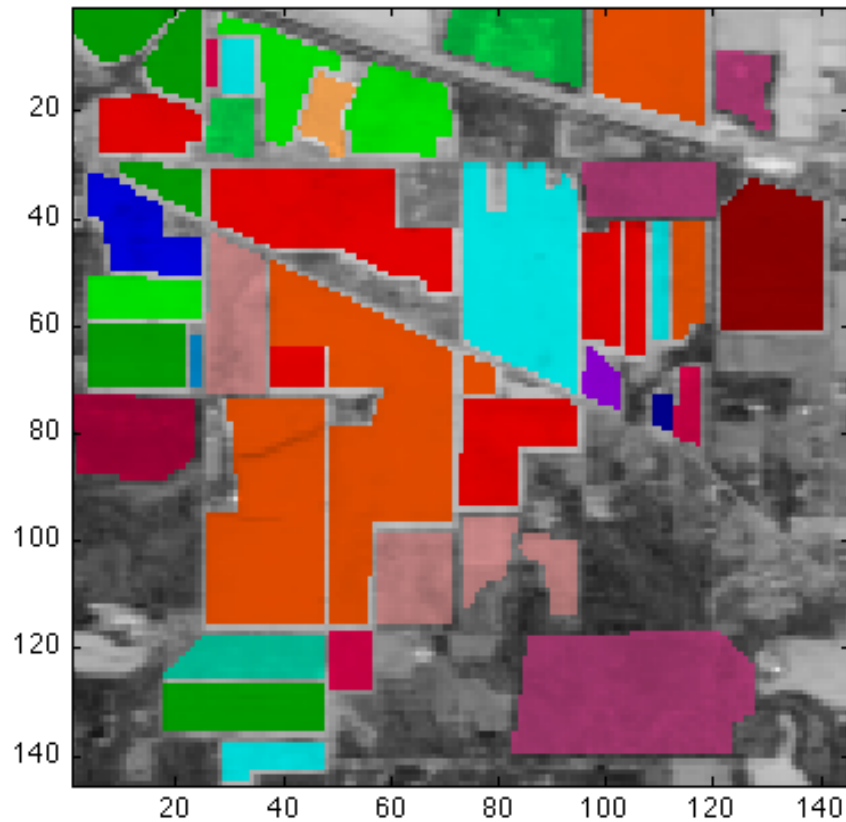
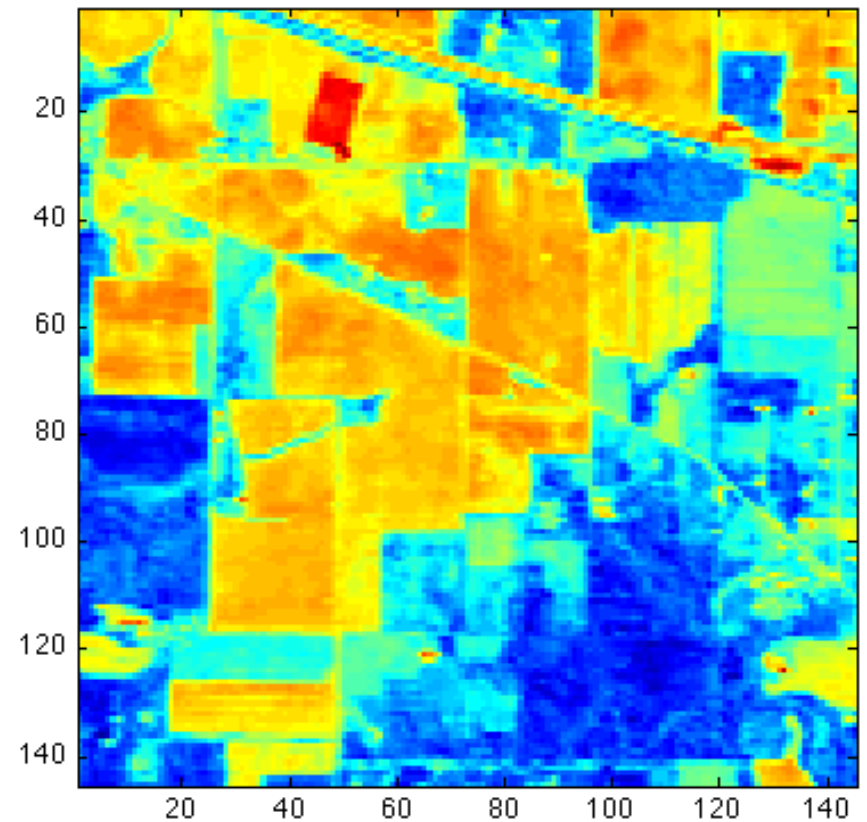
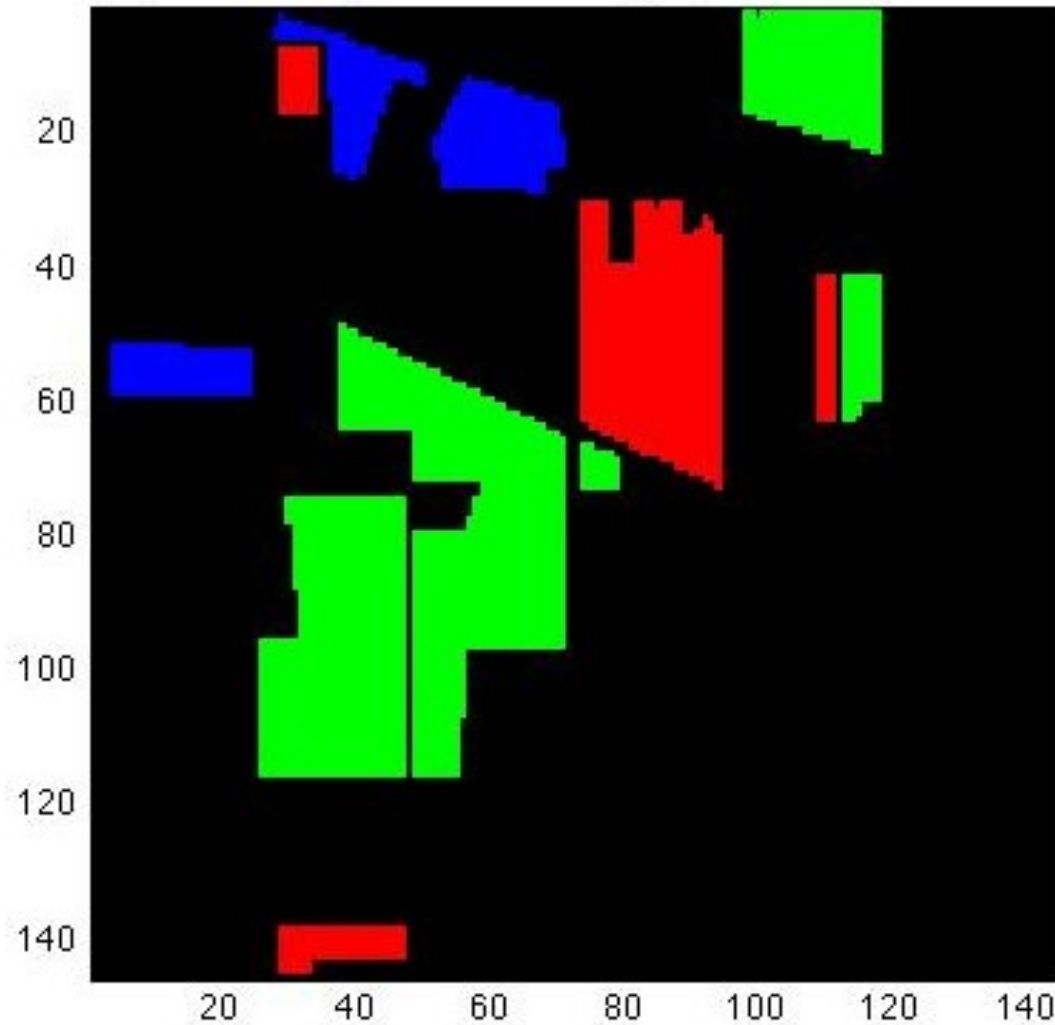


Image of Scores on PC 1 (72.48%)

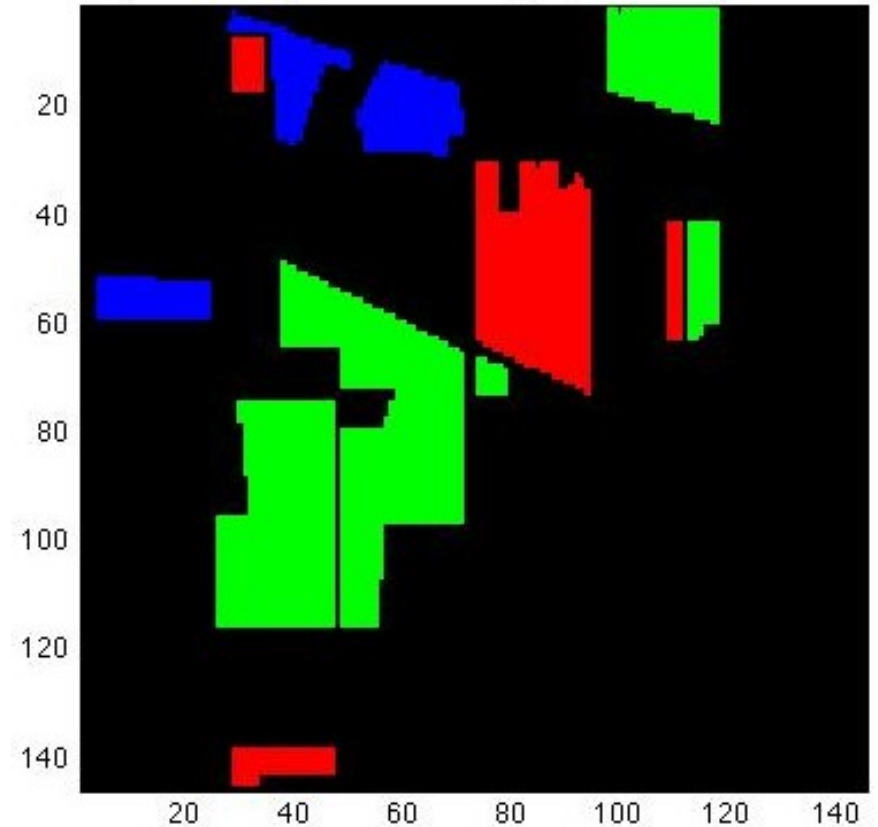
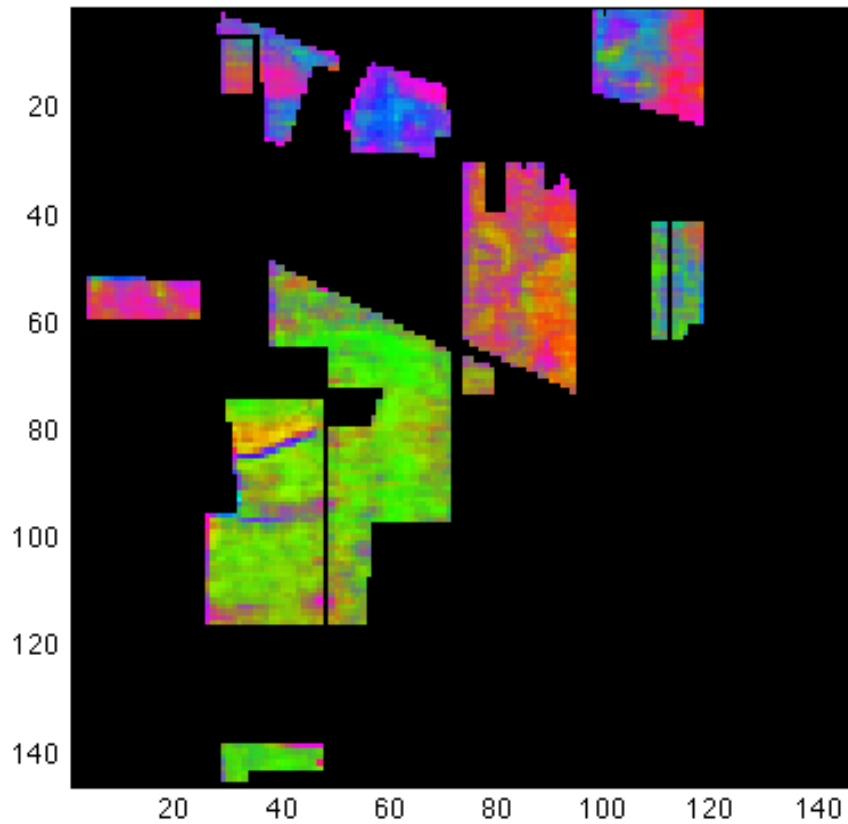


# Soybean Fields



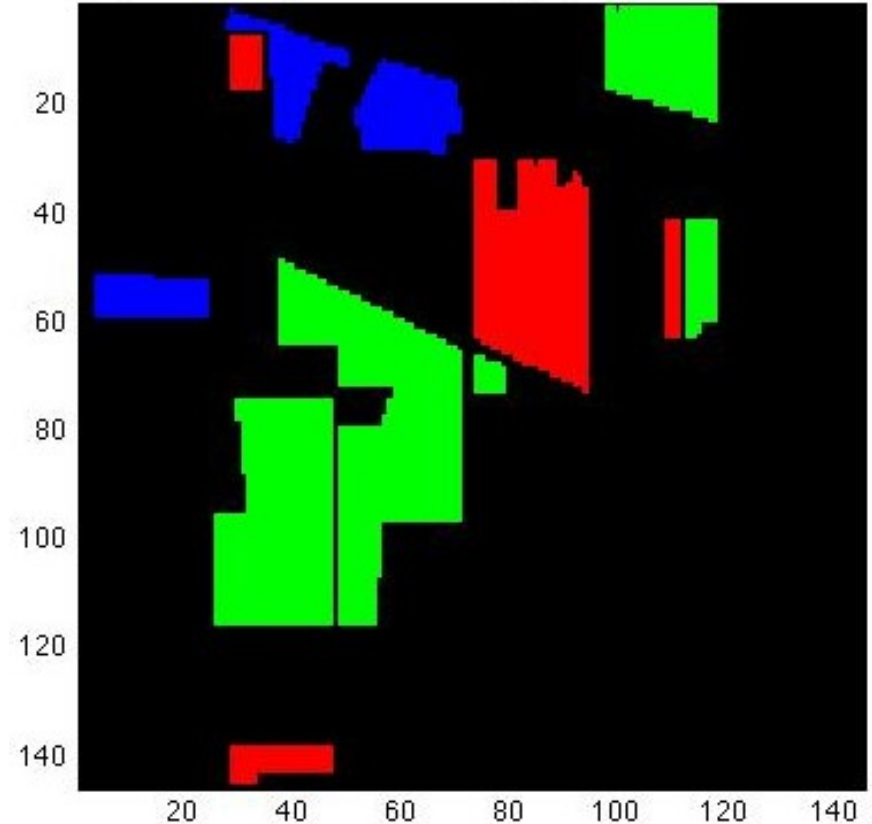
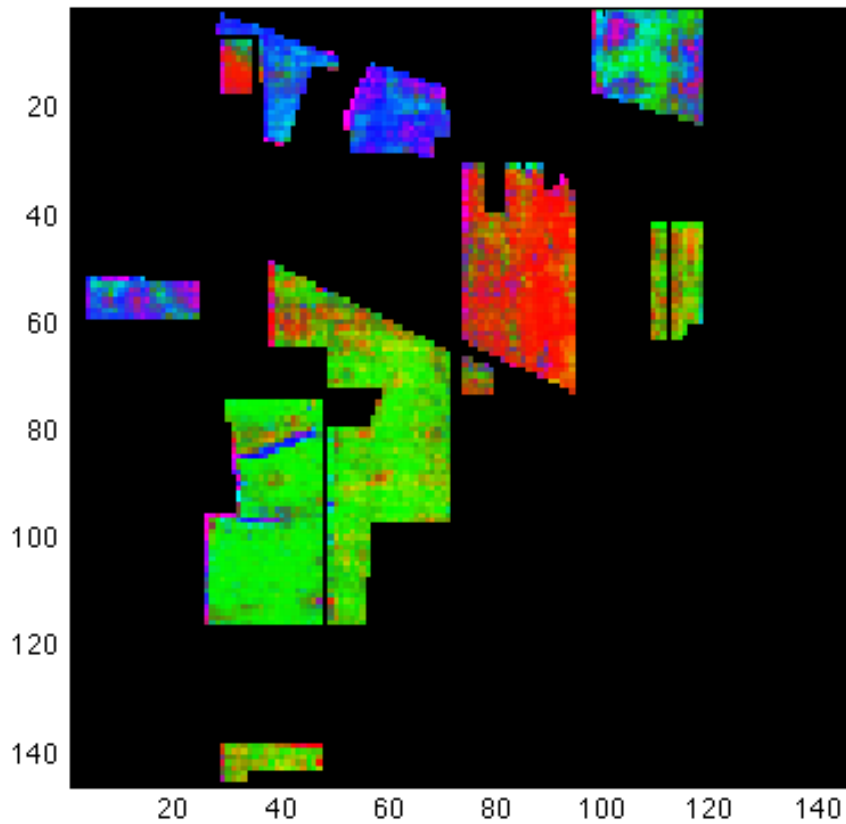
Soybeans no till  
Soybeans min  
Soybeans clean

# PLS-DA, Mean-Center Only



Class Probability Image

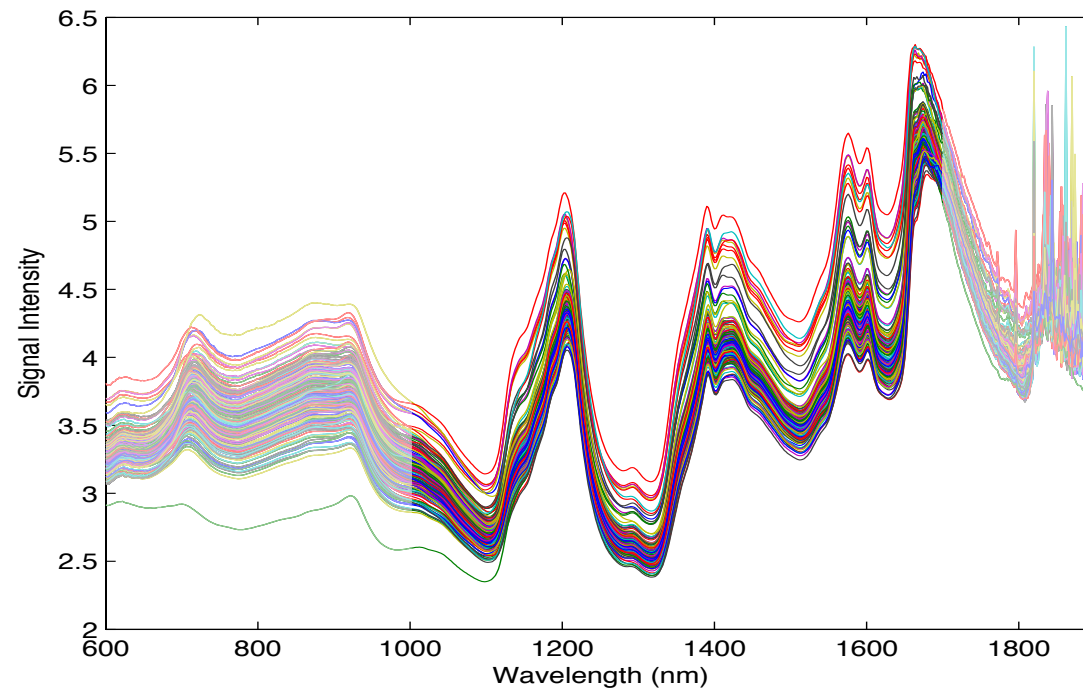
# PLS-DA, EPO 1-PC



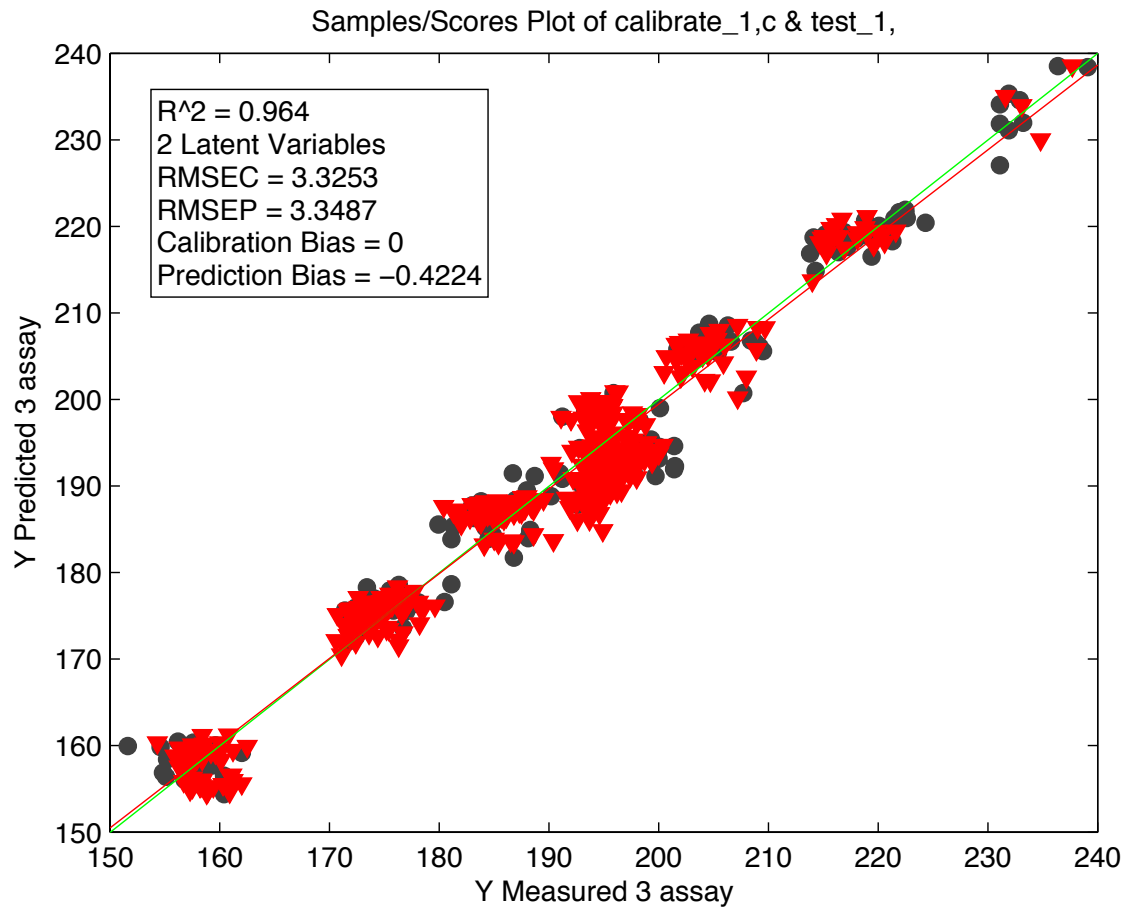
Class Probability Image

# Example Calibration Data

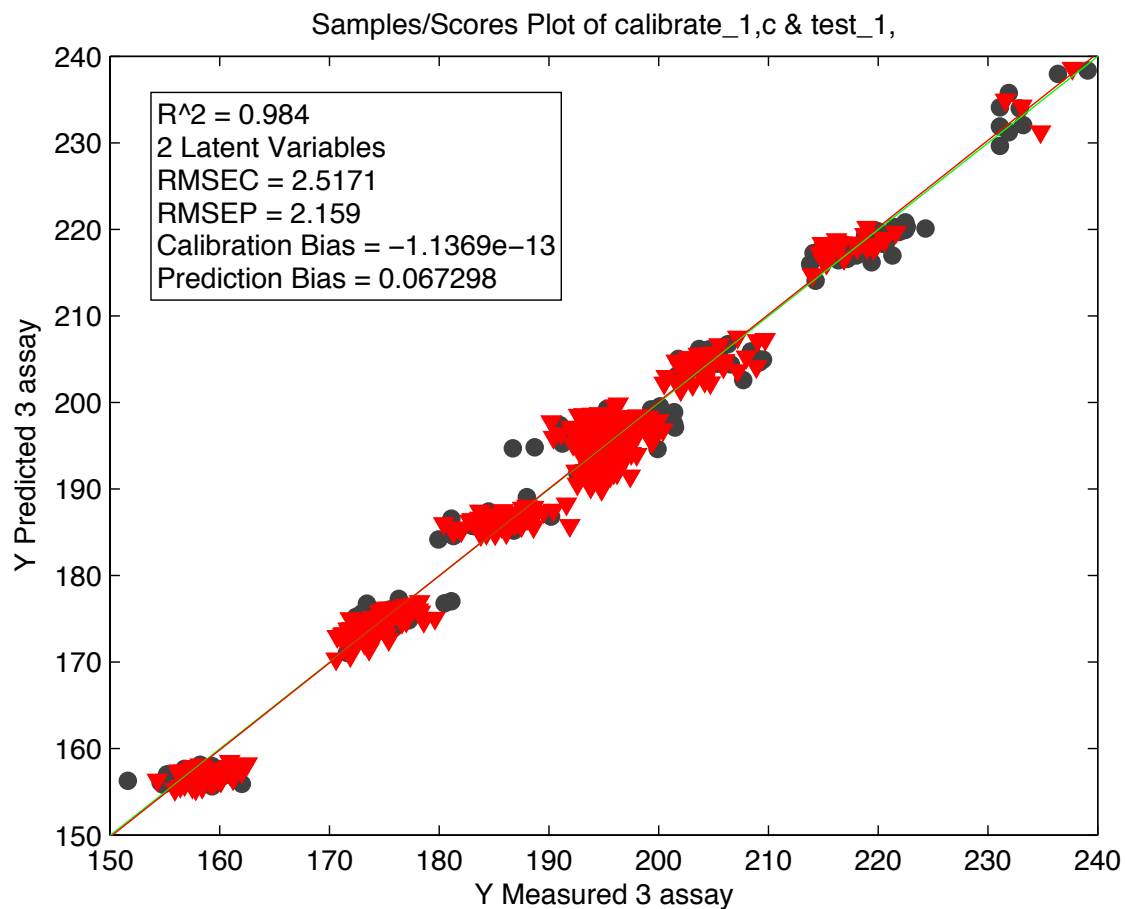
- IDRC-2002 Shootout data
- NIR Transflectance of pharmaceutical tablets
- Goal is to predict assay value



# Calibration and Test with MSC & MC

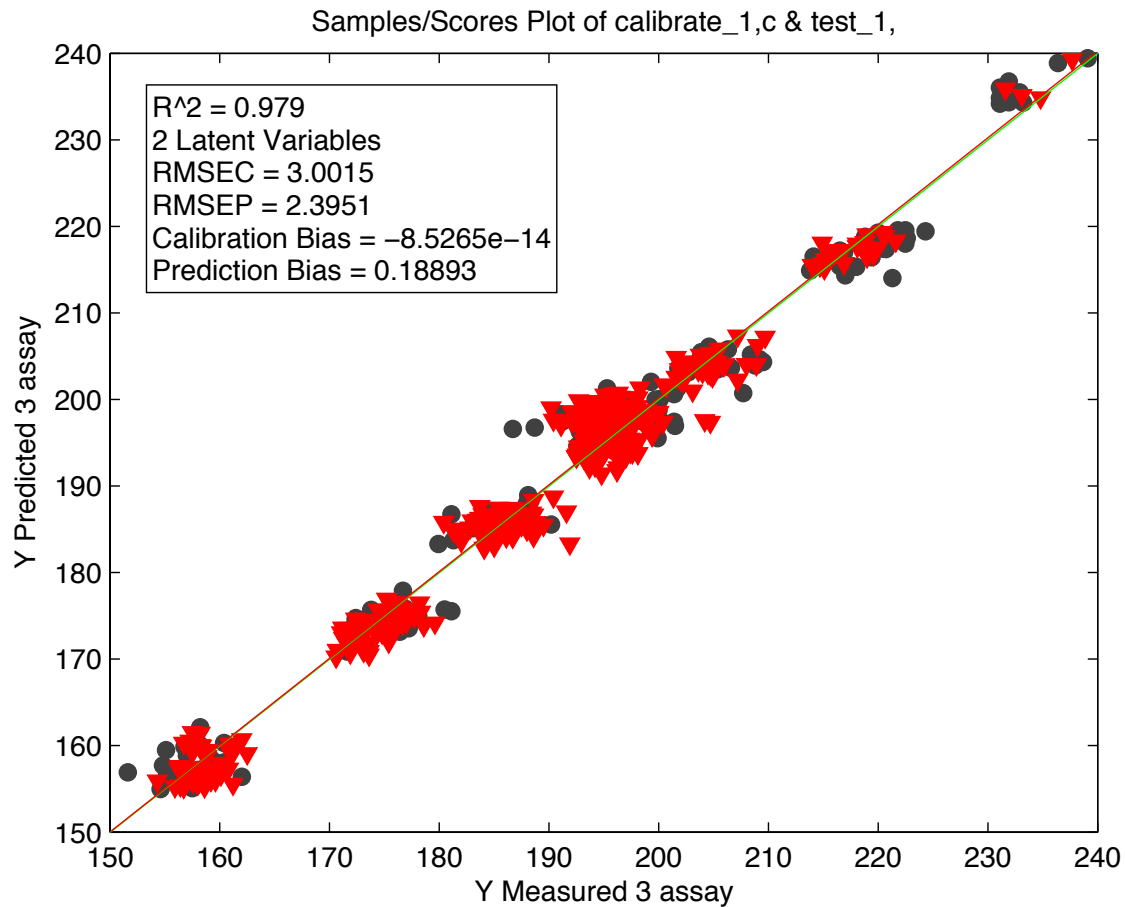


# With MSC, GLS & MC

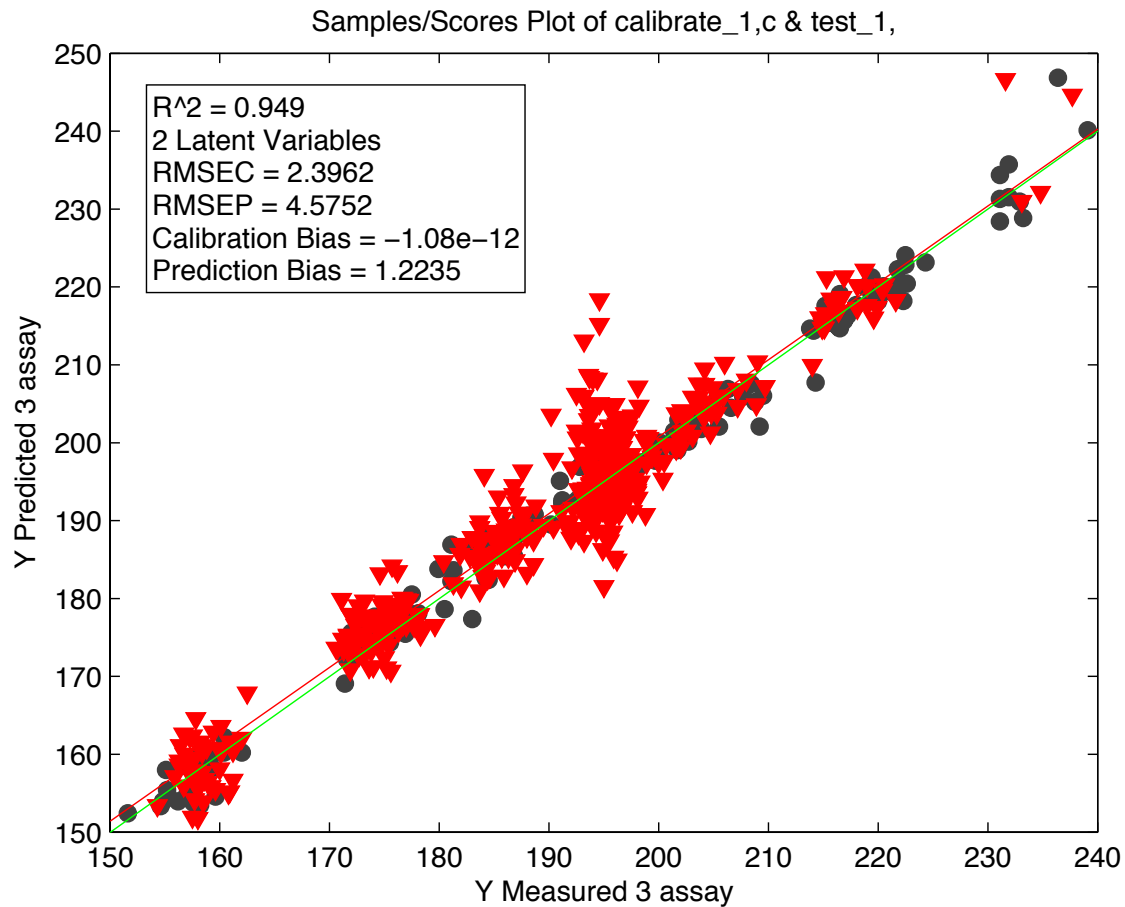




# With MSC, EPO & MC



# With MSC, ELS and MC



# Conclusions

- Main differences between methods are
  - How the clutter is defined
  - Whether the de-weighting is hard or soft
- Filtering methods are more similar than published statements might have you believe
- Methods achieve similar results, model performance generally improved (except O-PLS, OSC)
- Interpretation of filtered results can be challenging – except OPLS (mostly)