

# Introduction to Multivariate Image Analysis (MIA)

©Copyright 1996-2013  
Eigenvector Research, Inc.  
No part of this material may be  
photocopied or reproduced in any form  
without prior written consent from  
Eigenvector Research, Inc.



## *Table of Contents*

- Intro to 3-way arrays and simple visualizations and size/shape analyses
- Practical Multivariate Image Analysis (MIA)
  - PCA, SIMCA, PLSDA and clustering
- Variance Filtering for Images:
  - Maximum Autocorrelation Factors, Maximum Difference Factors, Generalized Least Squares Weighting (MAF, MDF, GLSW)
- Multivariate Image Regression and Quantitative Analyses
  - Partial Least Squares, Classical Least Squares and Multivariate Curve Resolution Models (PLS, CLS, MLR)



## Resources

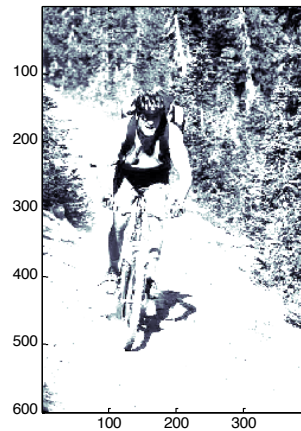
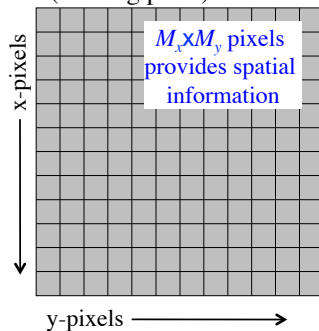
- *Hyperspectral Image Analysis*, eds. P. Geladi and H. Grahn, Wiley (2007), ISBN 978-0-470-01086-0
- *Chemometrics*, M.A. Sharaf, D.L. Illman and B.R. Kowalski, Wiley-Interscience (1986) ISBN 0-471-83106-9
- *Multivariate Analysis*, K.V. Mardia, J.I. Kent and J.M. Bibby, Academic Press, (1979) ISBN 0-12-471252-2
- *Multivariate Calibration*, H. Martens and T. Næs, John Wiley & Sons Ltd. (1989) ISBN 0-471-90979-3
- *Chemometrics: a textbook*, D.L. Massart et al., Elsevier (1988) ISBN 0-444-42660-4
- *Chemometrics: A Practical Guide*, K.R. Beebe, R.J. Pell, M.B. Seasholtz, Wiley (1998) ISBN 0-471-12451-6
- *Multivariate Data Analysis In Practice*, Kim H. Esbensen, CAMO ASA (2000), ISBN 82-993330-2-4
- *A user-friendly guide to Multivariate Calibration and Classification*, T. Næs, T. Isaksson, T. Fearn, T. Davies, NIR Publications(2002), ISBN 0-9528666-2-5
- Journal of Chemometrics
- IEEE Trans. on Geosci. and Remote Sensing
- 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
- ....

3



## Univariate Image

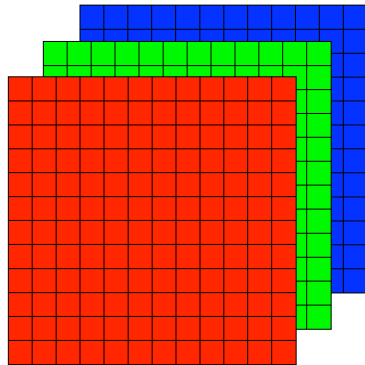
- Grey scale
  - each pixel is an number defining an intensity level e.g.,
    - integer (0 to 255) unsigned 8-bit
    - integer (0 to 4095)
    - double (floating point)



4

## Multivariate Image (3 Variables)

- Red/Green/Blue (RGB) (e.g. JPEG)
  - each layer defines color intensity level
  - much more information-rich



 **EIGENVECTOR**  
RESEARCH INCORPORATED

5

## Image Analysis

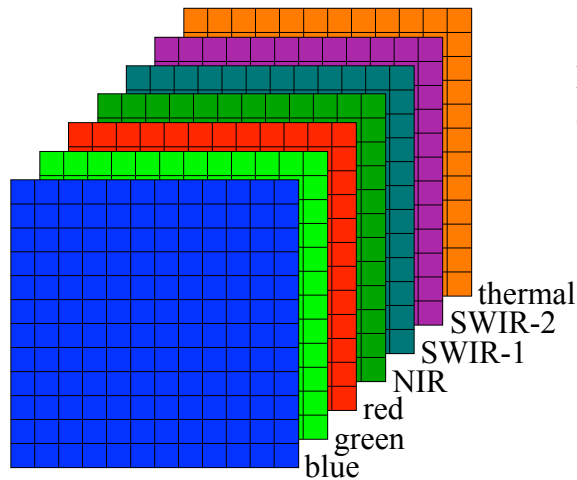
- Many methods have been developed to examine the spatial structure w/in an image
  - the methods recognize spatial patterns within an image
    - based on the light / dark contrast and continuity of regions
  - edge detection, image sharpening, wavelets
  - particle size distributions, machine vision, medical applications, security, ...
- MIA has been traditionally applied to the spectral dimension first followed by spatial analysis
  - some methods that examine both are appearing

 **EIGENVECTOR**  
RESEARCH INCORPORATED

6

## Multivariate Image (4-10 Variables)

- Measure at several wavelengths (e.g., Landsat)



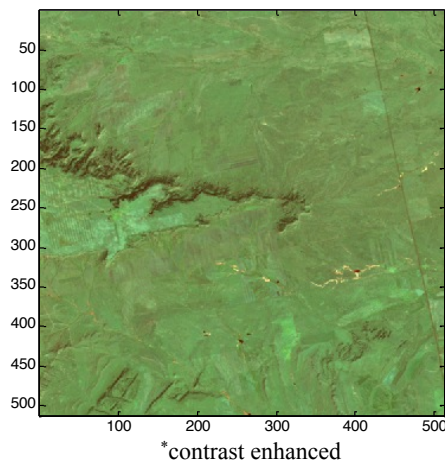
How should we display  
a seven variable image?



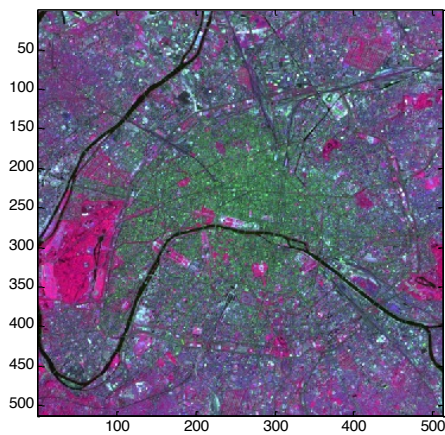
7

## Multivariate Image (4-10 Variables)

- Choose 3 of 7 (Landsat)  
Montana (blue/SWIR-1/thermal)



Paris (NIR/blue/SWIR-1)\*

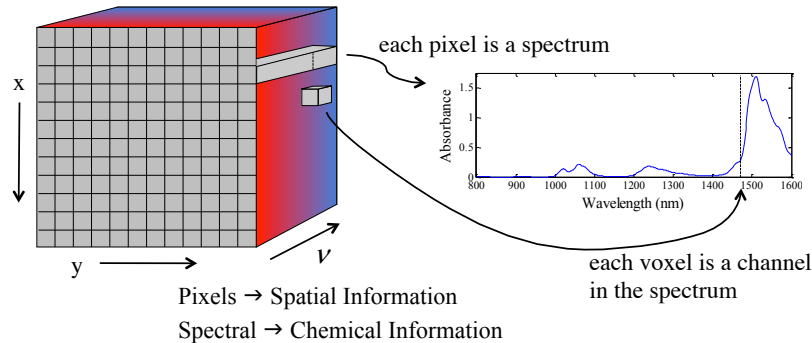


8



## Hyperspectral Image (>10 Variables)

- Spectrum at each pixel
  - could be 100-1000s of variables
  - often floating point double 10-100s Mbytes



9



## Multivariate Images

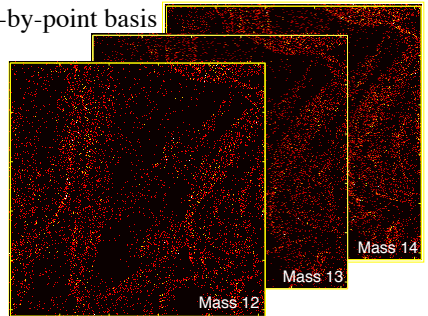
- Data array of *dimension three* (or more)
  - where the first two dimensions are *spatial* and
  - the last dimension(s) is a function of another variable (e.g, spectroscopy).
- Chemical system(s) of interest include
  - microscopic, medical, machine vision, process monitoring crystallization, stand-off and remote sensing, ...
  - vapors, liquids, solids (or combination)
  - visible, infra-red, Raman, mass spectroscopy, ...

10



## Physics of Measurement

- Point scanning
  - spectra measured on a point-by-point basis
  - secondary-ion mass spec
  - atomic force microscopy
  - surface Raman
- Line scanning
  - push broom
- Focal plane array
  - images can be acquired very quickly

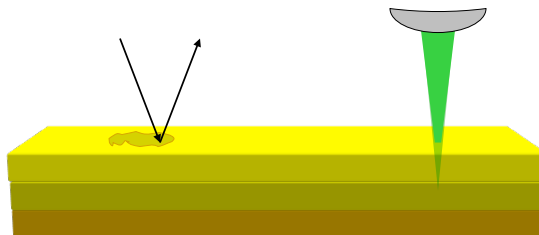


11



## Volumetric Analysis Techniques

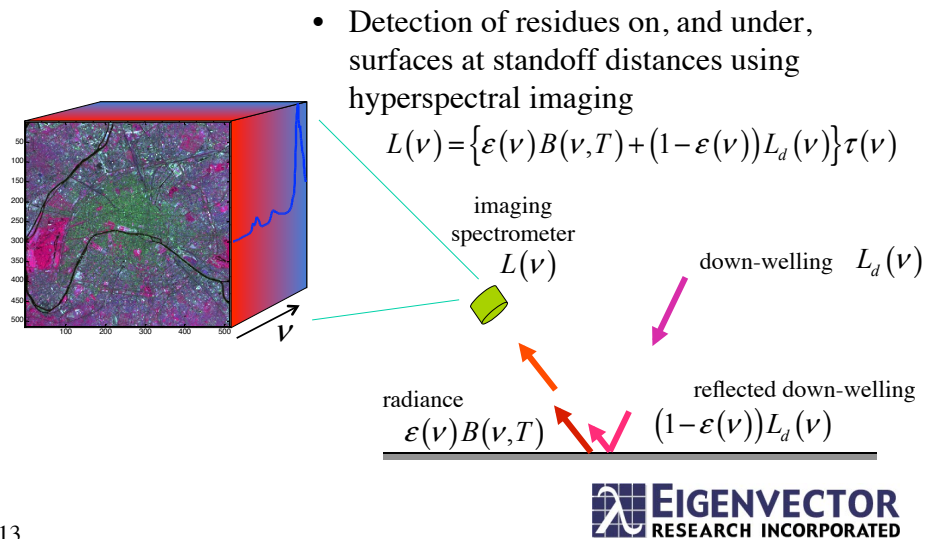
- Confocal Wavelength Resolved Imaging
- Surface Ablation Techniques
- Produces multivariate data in *3-dimensional space*



12



## Standoff and Remote Sensing



13

## Simple Image Analysis Tools

- TrendTool – Univariate Data Investigation
  - Analyze multivariate data using simple univariate measurements
- Image Manager – Data Manipulation and Analysis
  - Concatenating / Manipulating (e.g. rotation) Images
  - Particle Analysis
- Image Exploration Tools
  - Cross-section, Drill, and Magnification
- Preprocessing

14

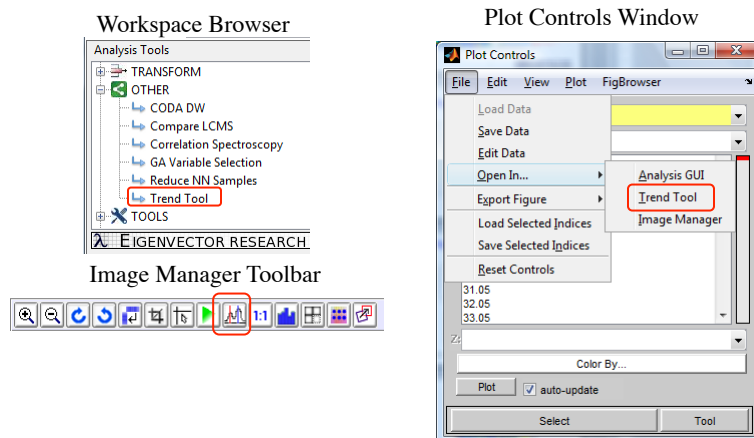
## TrendTool

- Display results of univariate calculations on multivariate data
  - Signal at given variable
  - Integrated signal across range of variables
  - Peak position
  - Peak width
- With or without baselines
- Ratio of measurements

15



## Opening TrendTool



16



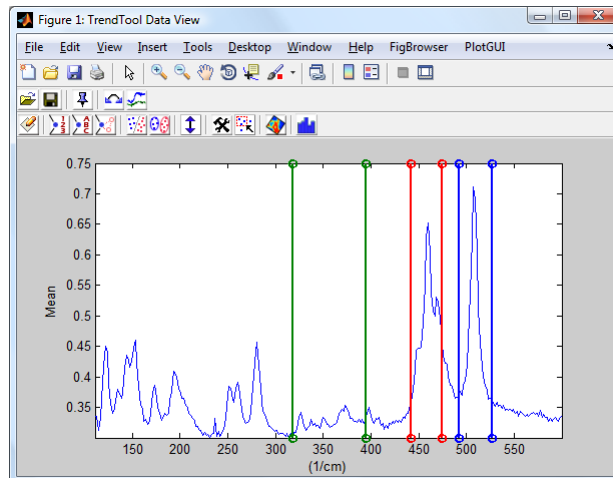
## TrendTool Windows: Data View

Use Data View to:

- Set analysis markers
- Choose analysis mode
- Select references and baseline points

Hints:

- Right-click white space to set marker or use toolbar button
- Drag markers to move
- Right-click markers to change types
- Use toolbar to save or load marker sets



17

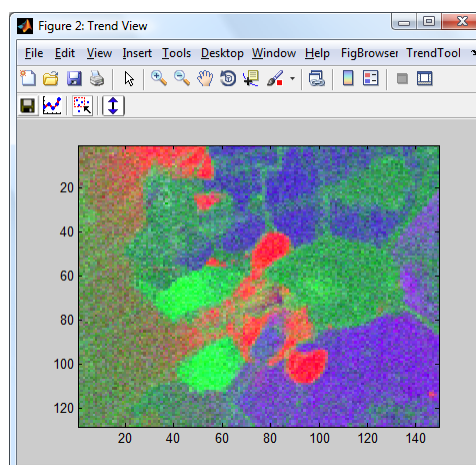
## TrendTool Windows: Trend View

Results displayed in Trend View

- Single marker displays with false-color
- Multiple markers display in RGB

Toolbar Buttons:

- [Autoscale icon] autoscale image
- [Select pixels icon] select pixels to display in Data View
- [Save/Spawn icon] save or spawn plot of results (respectively)



18

## TrendTool Analysis Modes

- **Height** – gives response at position (single marker)
- **Area** – gives integrated response between markers
- **Position** – gives position of peak response between markers
- **Width** – gives full width at half height between markers

"Add Reference" to subtract a single point baseline.  
 Convert reference to baseline (via right-click) to do two-point linear baseline.

"Normalize to Region" to normalize all regions to the response of the selected region.

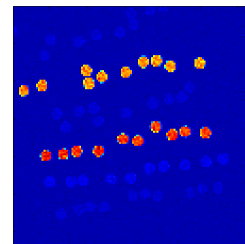
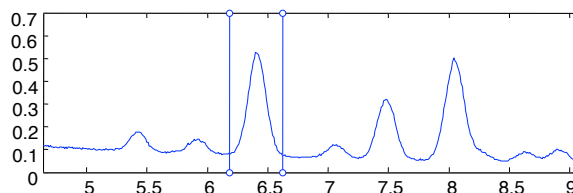
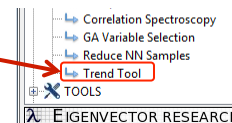
19



## TrendTool Example

**Example:** "wires" dataset  
 Energy Dispersive X-Ray Spectroscopy (EDS)  
 Image of wires composed of different alloys.

- Workspace Browser: Model Cache > Demo Data
- Drag "Wire Alloy Image" to TrendTool in Other Analysis Tools
- Use TrendTool to look at various peaks (right-click peak to change to peak type)



20



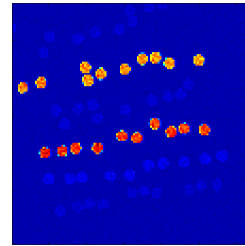
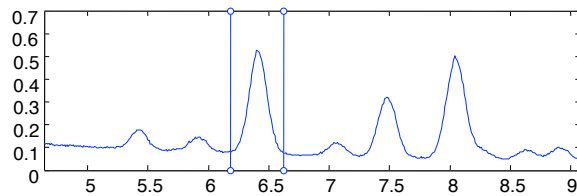


## Image Exploration

- Cross-section Tool – Transect of spatial dimension
- Drill Tool – Profile through variables of image
- Magnification Tool – Enhance spatial visibility

**Example:** "wires" dataset using TrendTool to look at one or more peaks...

Use "Spawn Results Plot" button on Trend View

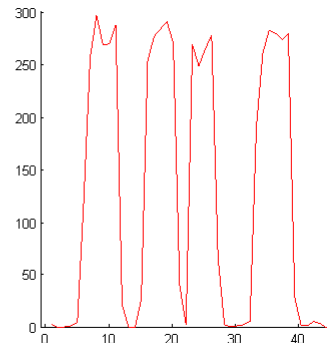
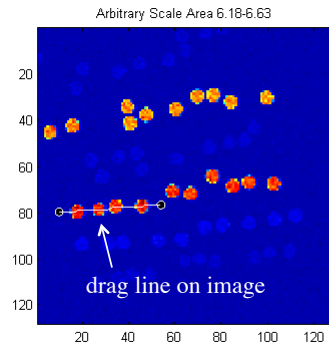
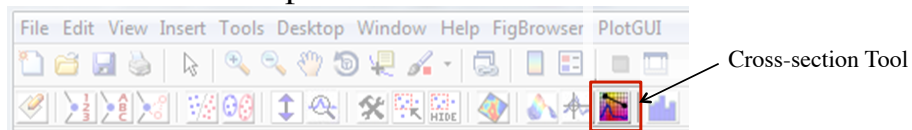


21



## Cross-Section Tool

- Transect of spatial dimensions



22



## Drill Tool

- Display of data under a given point

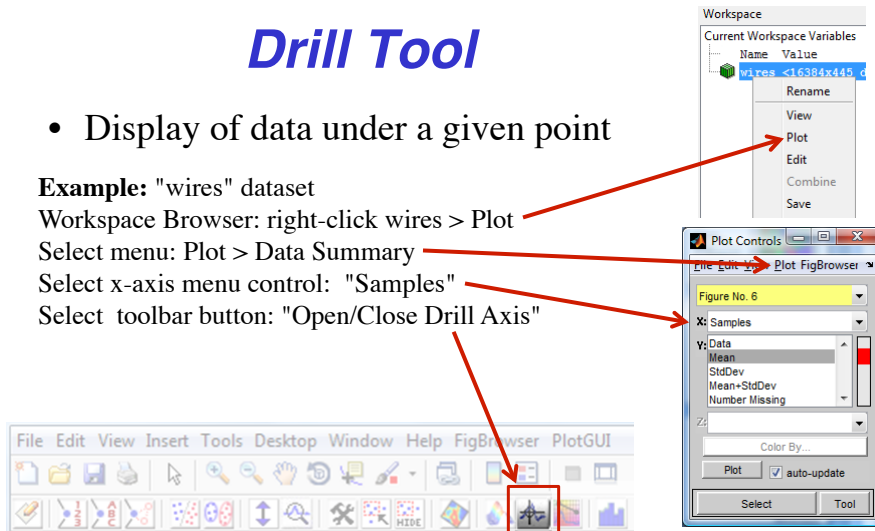
**Example:** "wires" dataset

Workspace Browser: right-click wires > Plot

Select menu: Plot > Data Summary

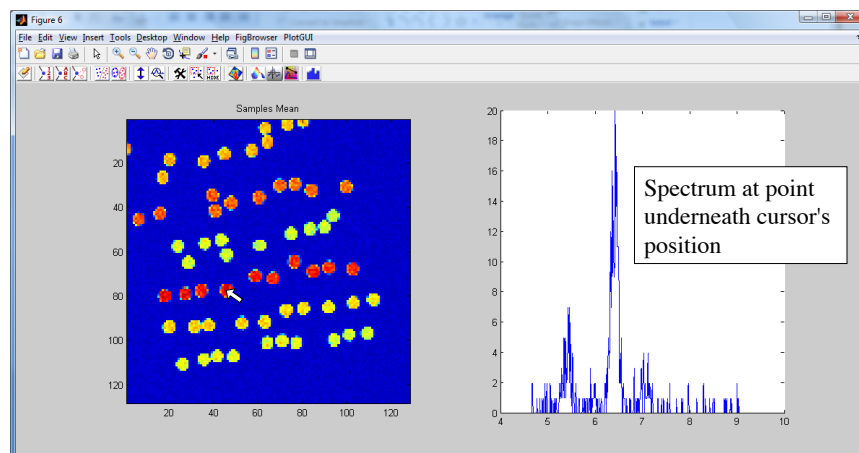
Select x-axis menu control: "Samples"

Select toolbar button: "Open/Close Drill Axis"



23

## Drill Tool

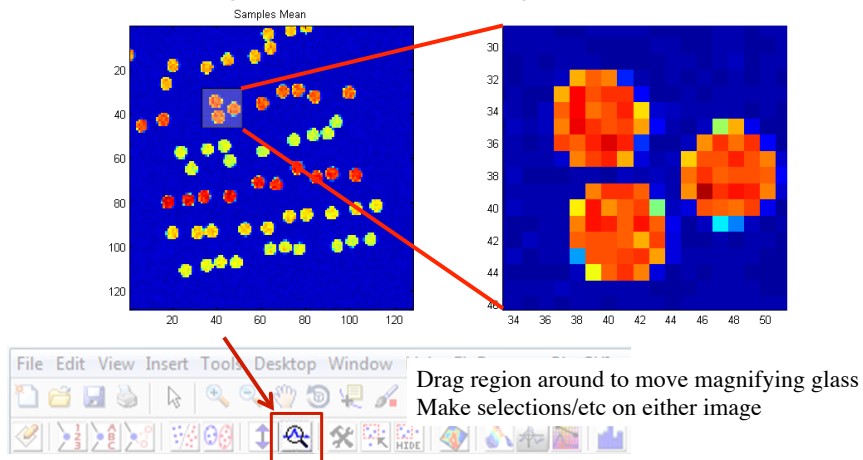


Double-click to view multiple spectra

24

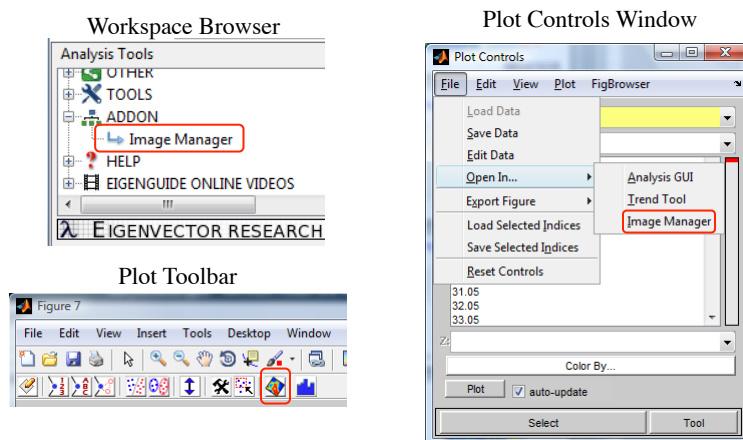
## Magnification Tool

- Show magnified view of image

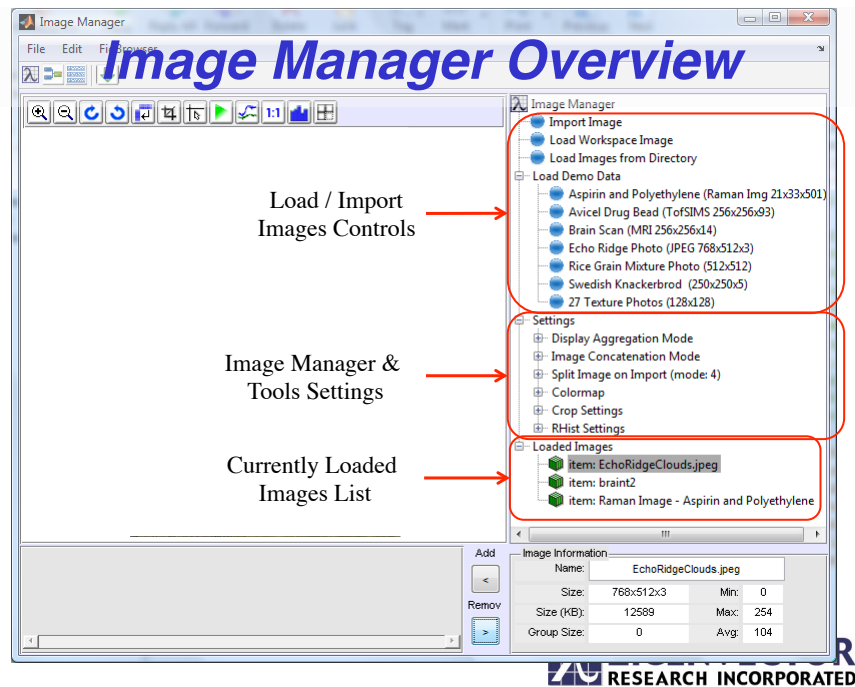


25

## Opening Image Manager

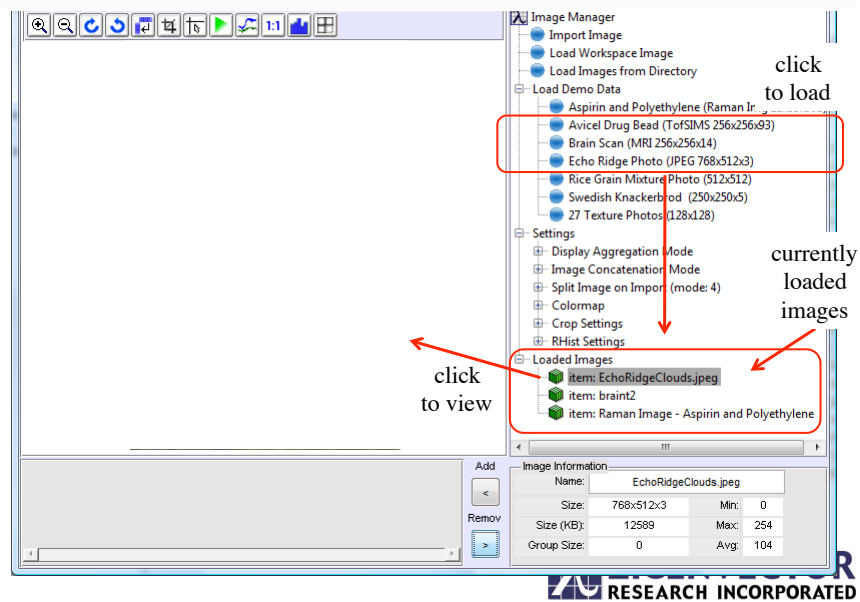


26



27

## Image Manager Overview

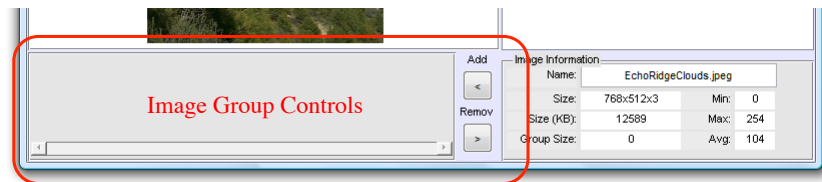


28

## Image Groups

Grouping allows you to:

- Combine images into a single DataSet for analysis
- Apply a univariate operation (rotate, crop, etc) to all images

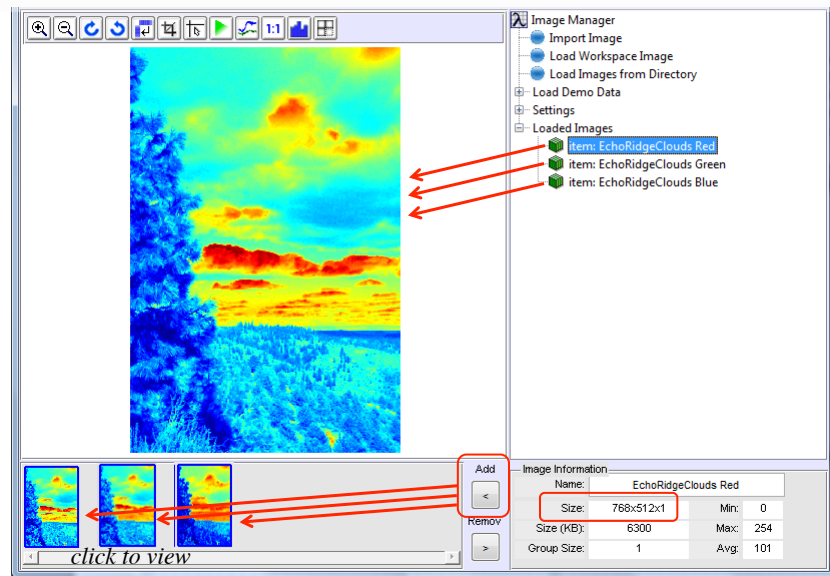


Example: combining three slabs of RGB image

29



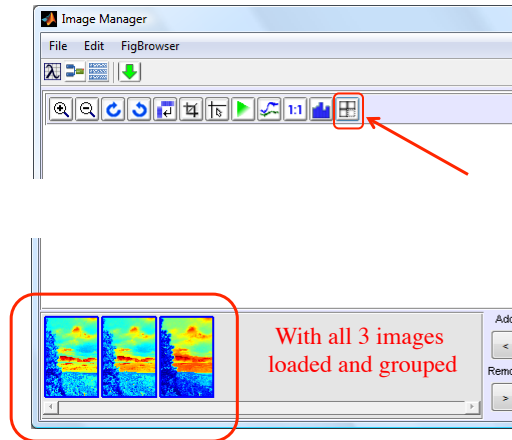
## Image Groups



30

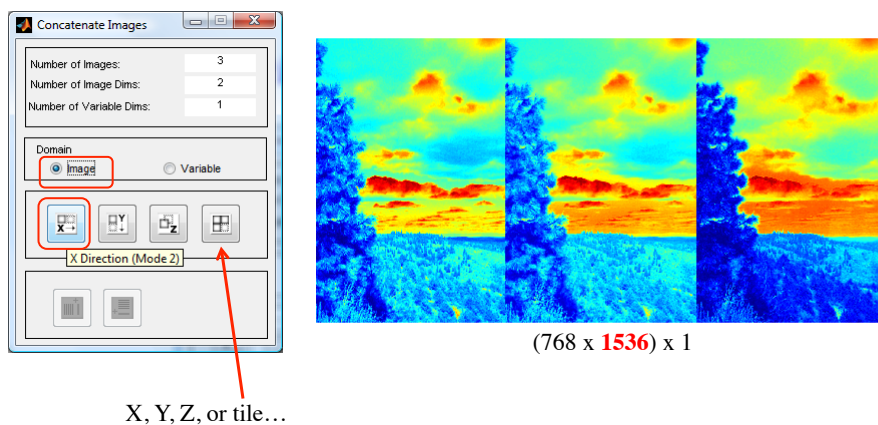


## Concatenating Images



31

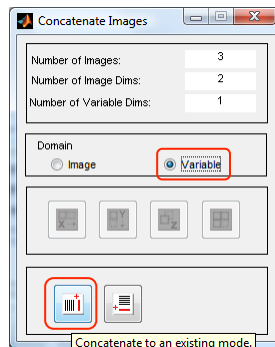
## Concatenating Images: Spatial Domain



32



## Concatenating Images: Variable Domain

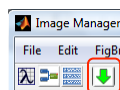
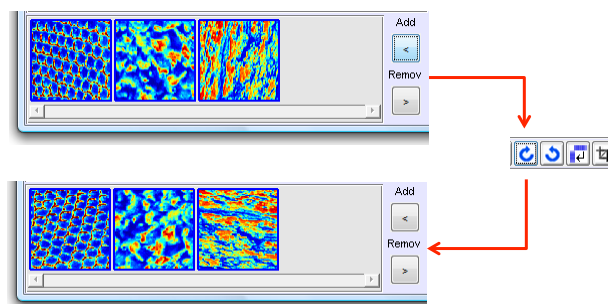


(768 x 512) x 3

33



## Group Manipulation Example: Rotation



Hint: to apply an action to only ONE image, click the "Apply Changes to Image Group" button until only one thumbnail is outlined in the image group pane.

34



## Particle Analysis

- Identify isolated regions (particles) in an image and give statistics on individual particles.
- Screen out particles and/or background.
- Create models based on particle statistics.
  - Particle outlier models (e.g. identify unusual particles)
  - Inferential models (e.g. drug activity based on particle statistics)
- Based on long-established ImageJ platform.

35



## Particle Analysis Example

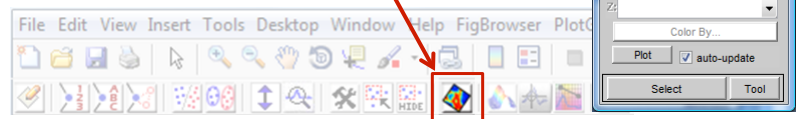
**Example:** "wires" dataset

Create Plot of data using Workspace Browser

Select menu: Plot > Data Summary

Select x-axis control: Samples

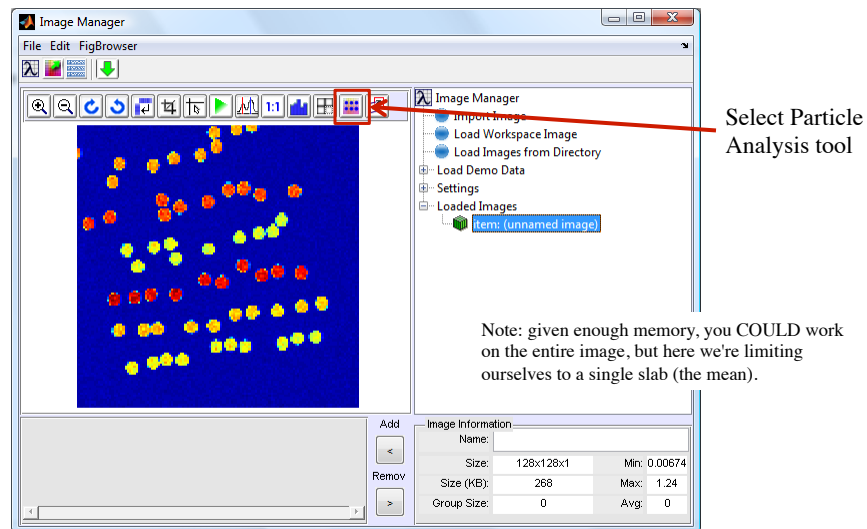
Select toolbar button: "Open in Image Manager"



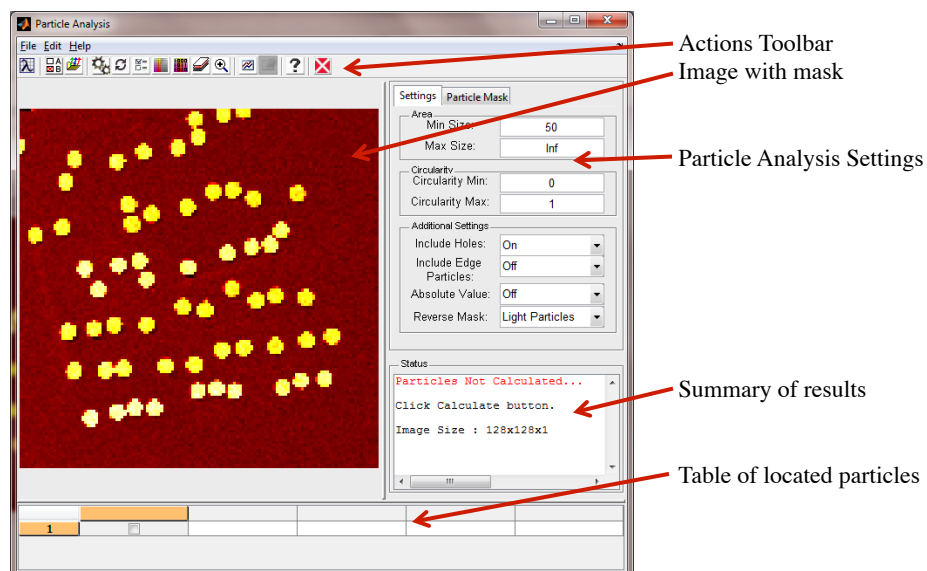
36



## Particle Analysis Example



37



38



## Particle Analysis Settings

- **Area Min/Max:** Ignore particles with area outside this range.
- **Circularity Min/Max:** Ignore particles outside this range.
- **Include Holes:** On = Include centers of particles even if below threshold.
- **Include Edge Particles:** On = Include particles which touch the edge of the image.
- **Absolute Value:** On = Consider positive and negative deviation from zero as "on" when making mask.
- **Reverse Mask:** Light Particles = Low signal is considered "off" (dark = not particle). Dark Particles = Low signal is considered "on" ("dark image" mode).

39



## Particle Mask Settings



Adjusts which pixels are considered particles

- **Threshold Slab:** For multi-slab images, which image slab is used to mask.
- **Threshold:** Signal level separating particles from background (slider adjusts or "Auto" checkbox does automatic threshold detection.)
- **Preprocessing:** Allows various operations on the binary image mask:
  - **Dilate:** Decrease mask around unmasked regions
  - **Erode:** Increase mask around unmasked regions.
  - **Smooth:** Smooth out noise in mask.

40



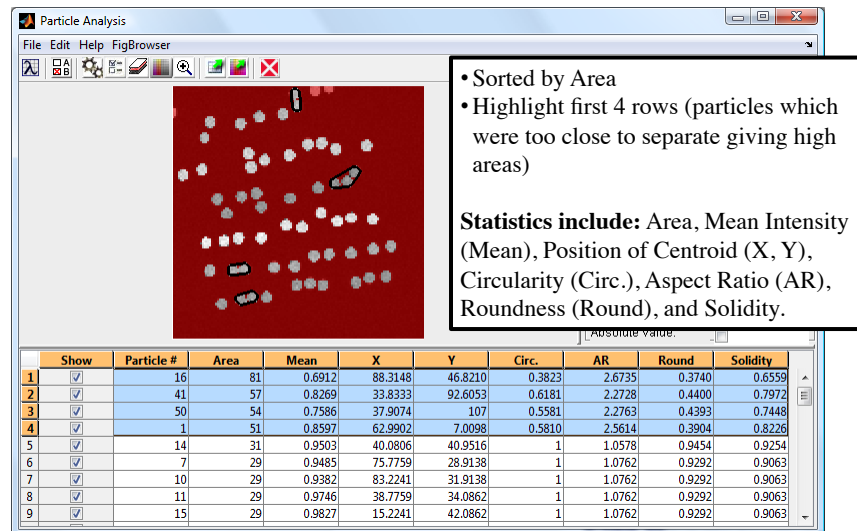
## Particle Analysis Example

- On "settings" tab, set Min Size to "2"
- On "Particle Mask" tab, set threshold to "0.4"
- Click "Recalc" button (next to threshold)
- Use Background Color and Grayscale settings to adjust display. 
- Select row of table to highlight corresponding particle.
- Select particle in image to highlight corresponding row of table.
- Sort by column using right-click menu.
- Use Export toolbar buttons to send table or image to Analysis. 

41



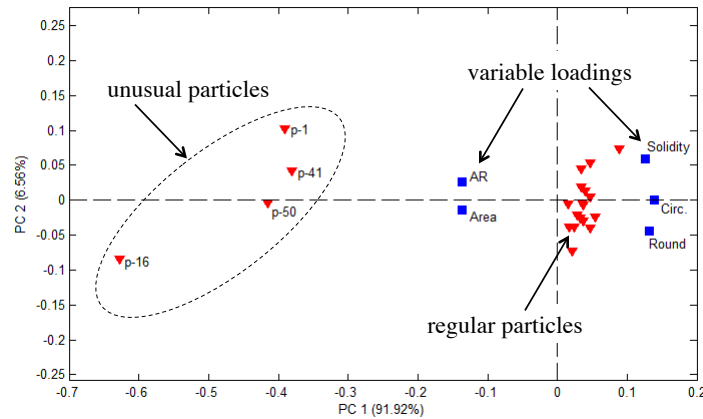
## Particle Analysis Example



42



## PCA of Particle Statistics Biplot of PCs 1 and 2



Autoscaled PCA model with mean intensity  
(Mean) and centroid (X, Y) variables excluded



43

## Using Preprocessing

	Show	Particle #	Area	Mean	X	Y	Circ.	AR	Round	Solidity
1	<input checked="" type="checkbox"/>	45	56	0.7529	33.71	92.7	0.6959	2.367	0.4224	0.855
2	<input checked="" type="checkbox"/>	17	46	0.7778	91	46	0.5318	2.408	0.4153	0.7541
3	<input checked="" type="checkbox"/>	15	32	0.6851	40	41	1	1	1	0.9412
4	<input checked="" type="checkbox"/>	37	29	0.7	111.1	80.81	1	1.064	0.9395	0.9063

•Add preprocessing:

- Erode (window = 3)
- Dilate (window = 3)

•Recalculate...

Only TWO joined particles



44



## Image-Oriented Preprocessing

- Image-specific preprocessing operates in pixel-space and are either Intensity or Binary based
- Intensity-Based Image Correction:
  - *Background Subtraction (Flatfield)*: Rolling-ball background subtraction for images.
  - *Min*: Min value over neighboring pixels. (filter out high-value pixels)
  - *Max*: Max value over neighboring pixels. (filter out low-value pixels)
  - *Mean*: Mean value over neighboring pixels. (filter out low/high pixels)
  - *Median*: Median value over neighboring pixels. (robust filter of low/high pixels)
  - *Trimmed Mean*: Trimmed mean value over neighboring pixels.
  - *Trimmed Median*: Trimmed median value over neighboring pixels.
  - *Smooth*: Spatial smoothing for images. (a weighted mean)

45



## Image-Oriented Preprocessing

- Binary-Based Image Correction
  - *Dilate*: Perform dilation on a binary image.
  - *Erode*: Perform erosion on a binary image.
  - *Close (Dilate+Erode)*: Perform dilation followed by erosion on a binary image.
  - *Open (Erode+Dilate)*: Perform erosion followed by dilation on a binary image.
- NOTE: Image-Oriented methods may break covariance (add multivariate rank) because variable slabs handled separately
- Standard variable-space preprocessing can be used too, but are spatially insensitive

46

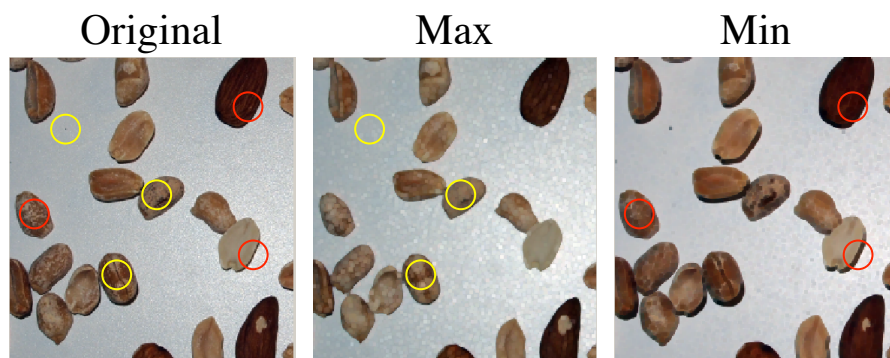


## ***Background Subtraction (Flat-field)***



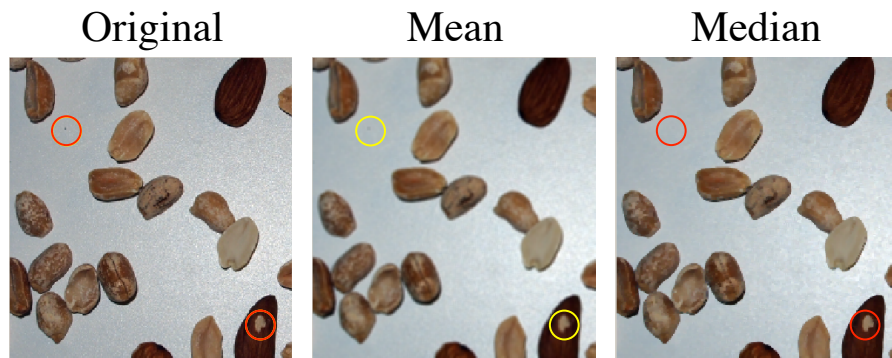
47

## ***Max & Min Preprocessing***



48

## *Mean & Median Preprocessing*



49

## *Displaying a Multivariate Image (4-10 Variables)*

- How to choose the 3 variables?
  - In which order should they be displayed?
- Doesn't choosing ignore potential information in the remaining variables?
- How could information be extract from the image?
- What happens when we go to more variables? ...
- .... Factor-based techniques
  - use the correlation structure to enhance S/N
  - really good for hyperspectral

50

## ***MIA: PCA-Based Methods***

- Many methods are based on the spectroscopic information in an image
  - although spatial information is ignored mathematically
  - images are examined for spatial structure
- PCA (Principal Components Analysis)
  - Exploratory analysis
- SIMCA (Soft Independent Method Class Analogy)
  - Classification

51



## ***Image PCA***

- Matricizing
- PCA: scores, scores images, loadings
  - unusual samples Q and  $T^2$
  - score-score plots, density plots
  - linking scores and image plane(s)
  - contrast enhancement

52



## PCA Math Summary

- For a data matrix  $\mathbf{X}$  with  $M$  samples and  $N$  variables (generally assumed to be mean centered and properly scaled), the PCA decomposition is

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

Where  $R = \lfloor \min\{M, N\} \rfloor$ , and the  $\mathbf{t}_k \mathbf{p}_k^T$  pairs are ordered by the amount of variance captured.

- Generally, the model is truncated to  $K$  PCs, leaving some small amount of variance in a residual matrix  $\mathbf{E}$ :

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

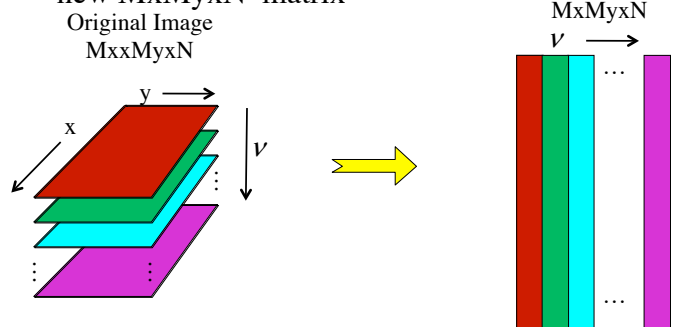
- where  $\mathbf{T}$  is  $M \times K$  and  $\mathbf{P}$  is  $N \times K$ .

53



## Matricizing (a.k.a. Unfolding)

- PCA works on  $\mathbf{X}$  ( $M \times N$ ) but the image is  $M \times x \times y \times N$ 
  - reshape by matricizing such that each pixel is a row in a new  $M \times My \times N$  matrix



54



## Properties of PCA

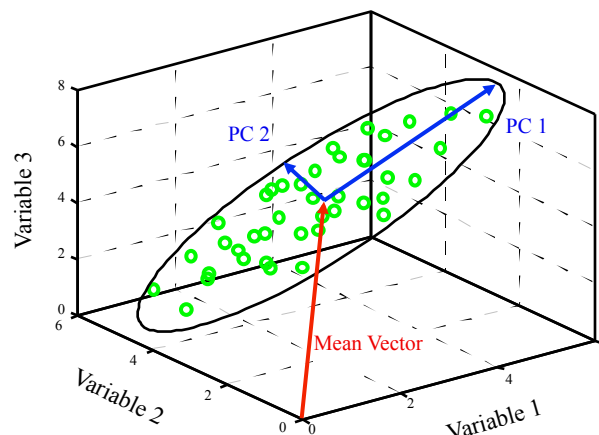
$$\mathbf{X} = \begin{bmatrix} \mathbf{t}_1 & \mathbf{t}_2 & \dots & \mathbf{t}_K \end{bmatrix} \begin{bmatrix} \mathbf{p}_1^T \\ \mathbf{p}_2^T \\ \dots \\ \mathbf{p}_K^T \end{bmatrix} + \mathbf{E}$$

- $\mathbf{t}_k, \mathbf{p}_k$  ordered by amount of *variance captured*
  - $\lambda_k$  are the eigenvalues of  $\mathbf{X}^T \mathbf{X} \rightarrow \mathbf{X}^T \mathbf{X} \mathbf{p}_k = \lambda_k \mathbf{p}_k$
  - $\lambda_k$  are  $\propto$  variance captured
- $\mathbf{t}_k$  (*scores*) form an orthogonal set  $\mathbf{T}_K (M \times K)$ 
  - describe relationship between *samples*  $\rightarrow$  *pixels* ( $M = M_x M_y$ )
- $\mathbf{p}_k$  (*loadings*) form an orthonormal set  $\mathbf{P}_K (N \times K)$ 
  - describe relationship between *variables*

55



## PCA Graphically



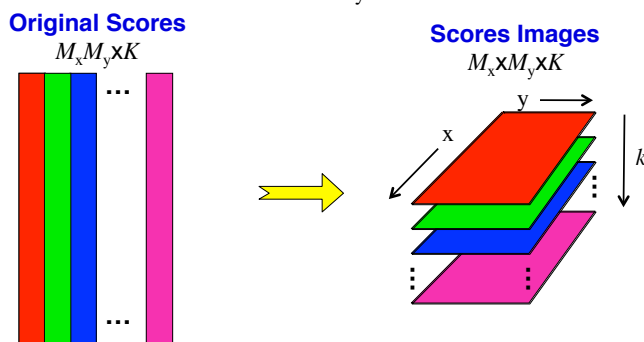
56





## Reshape Scores To Images

- PCA gives scores  $\mathbf{T}$  ( $M \times K$ ) which is reshaped to scores images ( $M_x \times M_y \times K$ )
  - each score vector is a  $M_x \times M_y$  scores image

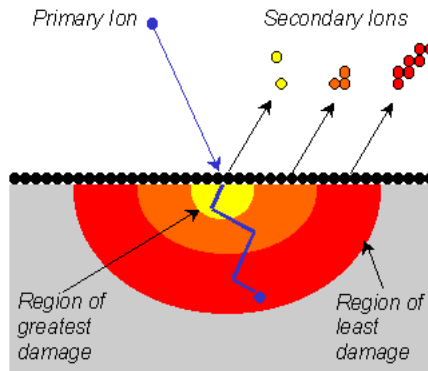


## Plots / Images for PCA

- scores and loadings plots are interpreted in pairs
  - plot  $\mathbf{t}_k$  vs sample number
    - find relationship between *samples*  $\rightarrow$  *pixels*
    - each  $M_x \times M_y \times 1$  score vector is reshaped to a  $M_x \times M_y$  matrix that can be visualized as a "*scores image*" showing spatial relationships between pixels
  - $\mathbf{p}_k$  vs variable number
    - relationship between *variables* responsible for observations in samples
- it is useful to plot  $\mathbf{t}_{k+1}$  vs.  $\mathbf{t}_k$  and  $\mathbf{p}_{k+1}$  vs.  $\mathbf{p}_k$ 
  - examine image and score / score plots

## TOF-SIMS of PMMA and Deuterated Polystyrene

- Time of flight secondary ion mass spectroscopy used for surface analysis
- Mass spectrum for each pixel
- Thanks to Physical Electronics for the data

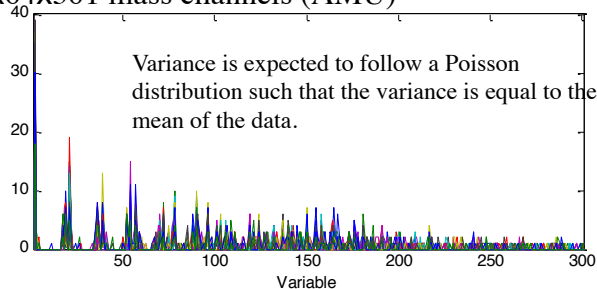


59



## Example Data

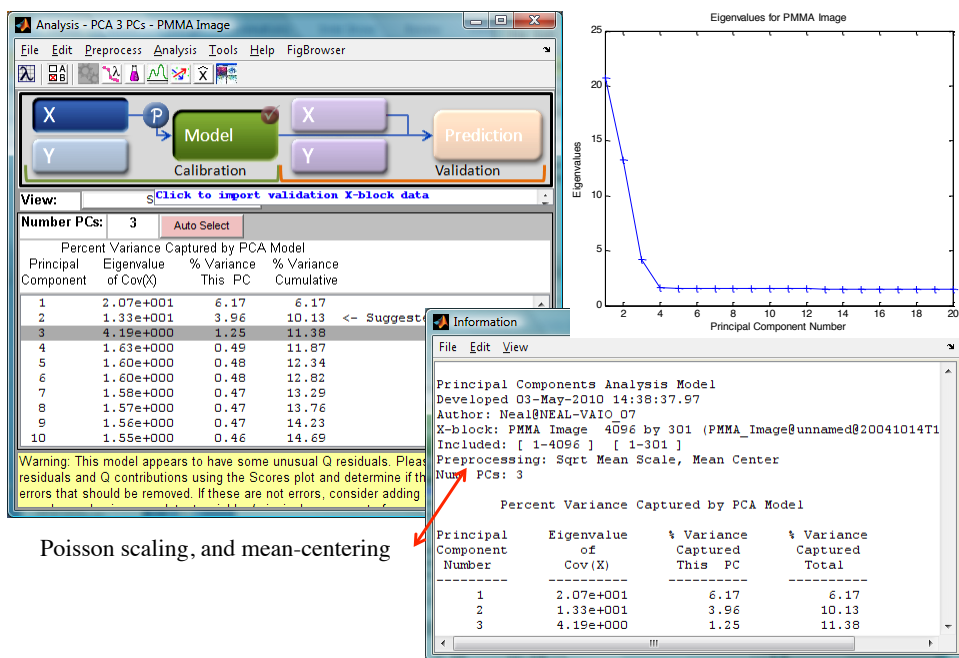
- Data is positive SIMS spectrum at each pixel (point) on a 64x64 grid
- 64x64x301 mass channels (AMU)



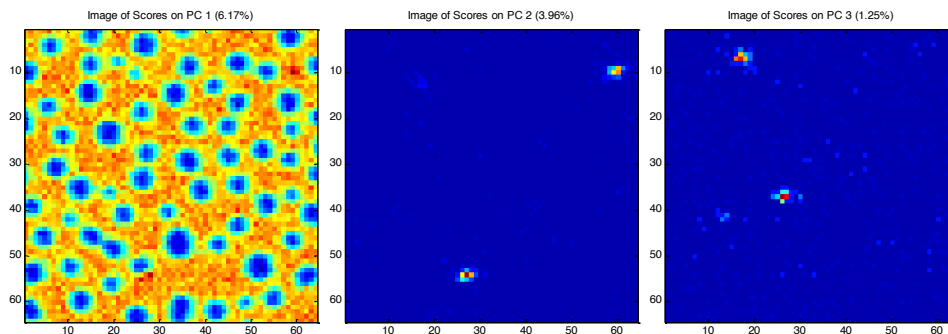
M.R. Keenan, "Multivariate Analysis of Spectral Images Composed of Count Data," in *Techniques and Applications of Hyperspectral Image Analysis*, H. F. Grahn and P. Geladi, eds. (John Wiley & Sons, West Sussex, England), 89-126, 2007.

60





61



Scores images show islands of polystyrene in PMMA and two sources of unusual variance



62

## PCA Statistics

- Limits can be set for
  - Q residual: lack of fit statistic
    - for a row of  $\mathbf{E}$ ,  $\mathbf{e}_m$ , and a row of  $\mathbf{X}$ ,  $\mathbf{x}_m$ ,  $m = 1, \dots, M$ 

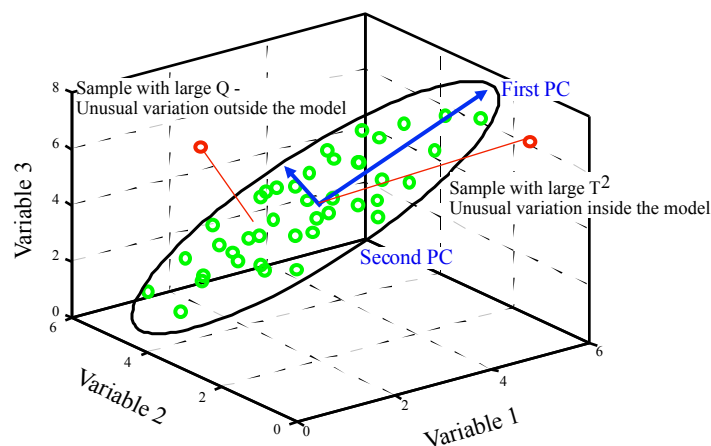
$$Q_m = \mathbf{e}_m \mathbf{e}_m^T = \mathbf{x}_m (\mathbf{I} - \mathbf{P}_K \mathbf{P}_K^T) \mathbf{x}_m^T$$
  - Hotelling's  $T^2$  statistic
    - for a row of  $\mathbf{T}_K$ ,  $\mathbf{t}_m$ , and  $K \times K$  diagonal matrix  $\boldsymbol{\Lambda}$ 

$$T_m^2 = \mathbf{t}_m \boldsymbol{\Lambda}^{-1} \mathbf{t}_m^T = \mathbf{x}_m \mathbf{P}_K \boldsymbol{\Lambda}^{-1} \mathbf{P}_K^T \mathbf{x}_m^T$$
  - and also for individual columns:
    - scores,  $\mathbf{t}_{mk}$
    - residuals  $\mathbf{e}_{mk}$

63

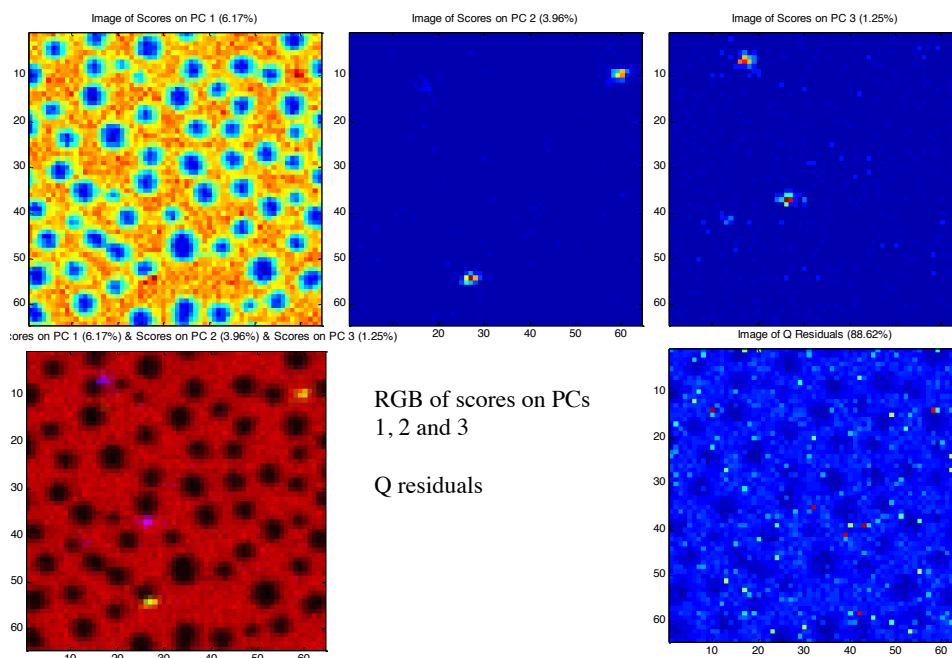


## Geometry of Q and $T^2$



64





Analysis - PCA (No Model) - New Image DataSet

File Edit Preprocess Analysis Tools Help FigBrowser

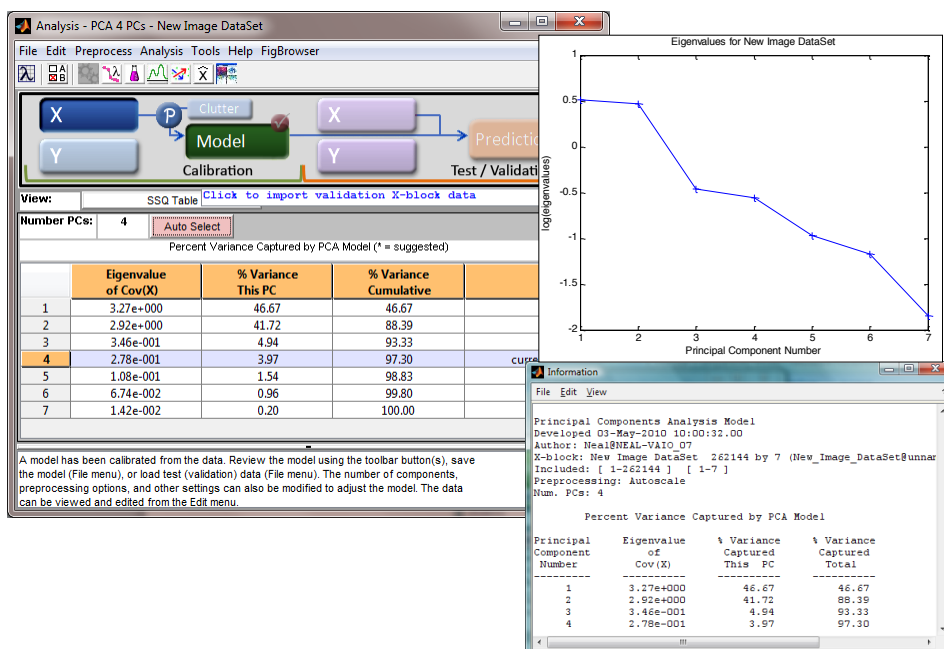
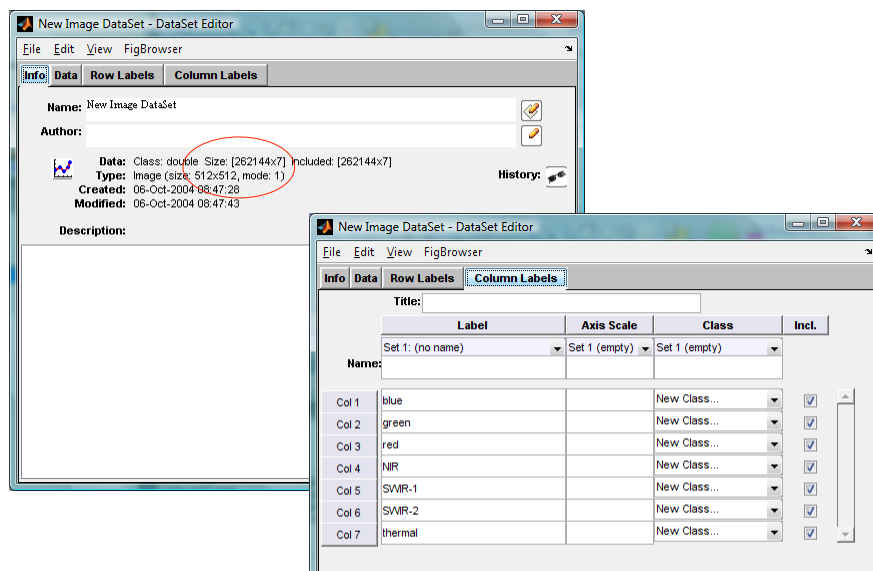
View: SSQ Table

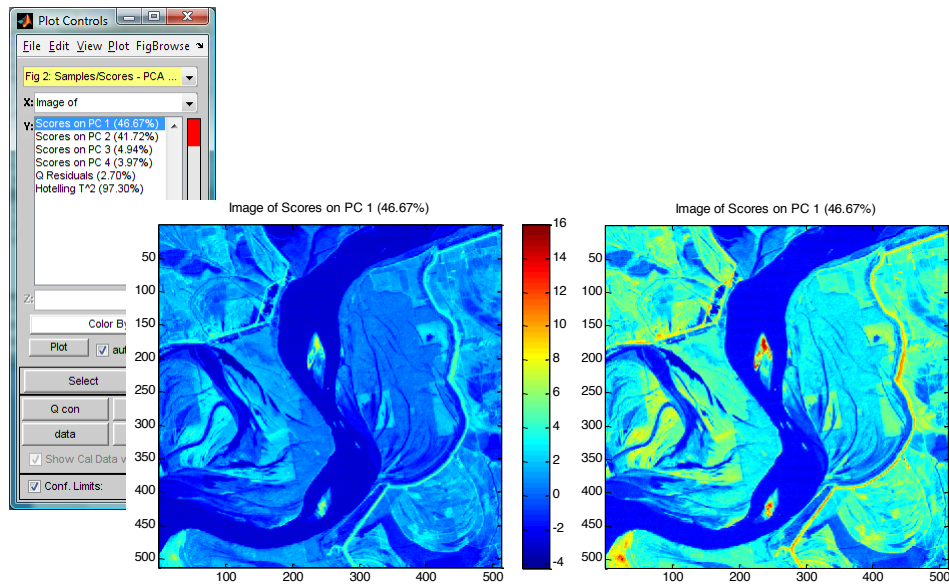
Number PCs: Auto Select

Percent Variance Captured by Principal Component Eigenvalue of Cov(X) % Variance Captured by This PC

Data has been loaded but no model exists. Set the preprocessing and other options (from the Preprocess and Tools menus) and calibrate a model (Calibrate button). The data can be viewed and edited with the Edit menu.

...\\MIADData

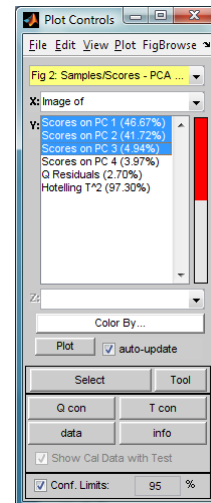




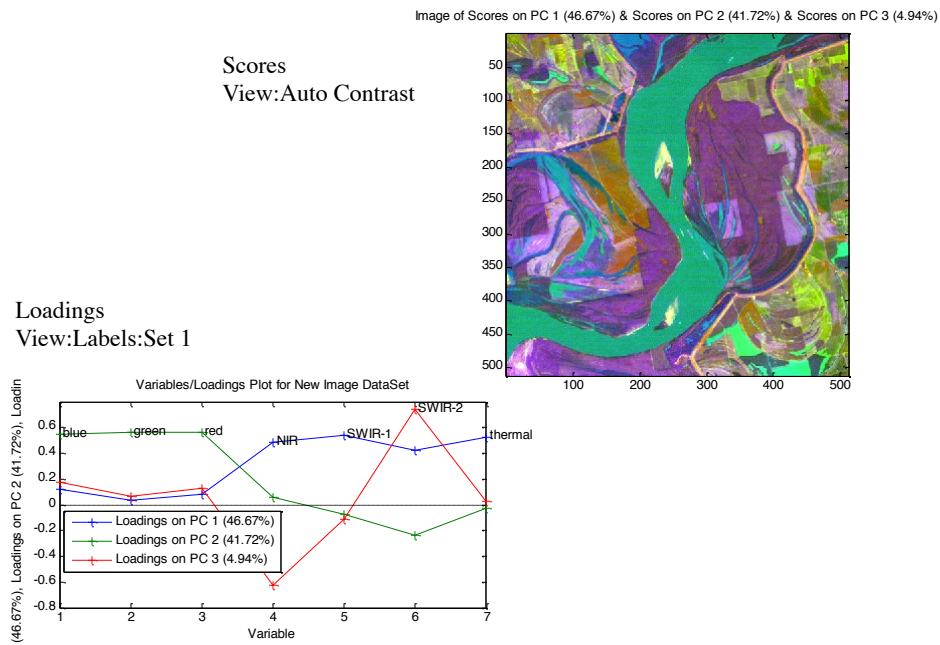
69

## Creating Color Images

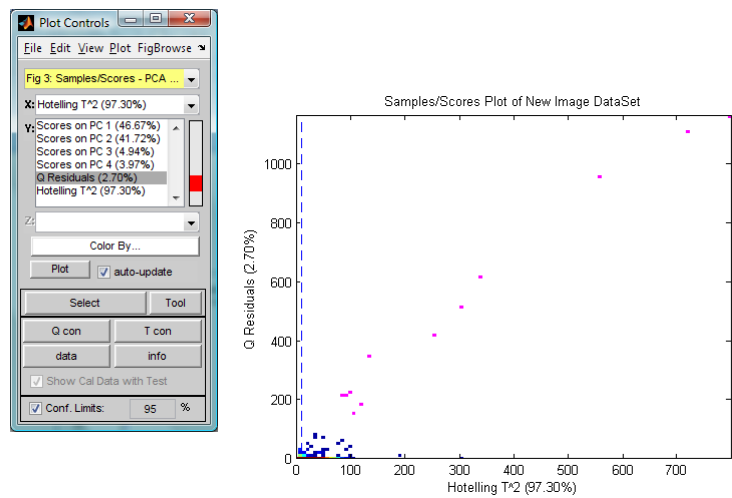
- Images are made of three colors:  
red, green and blue
  - e.g., scaled to integers 0-255 for 8-bit color
- Scores can be used to define the colors
  - PC 1 = red, PC 2 = green, PC 3 = blue



70



71



pixels with high Q and T<sup>2</sup> have been selected

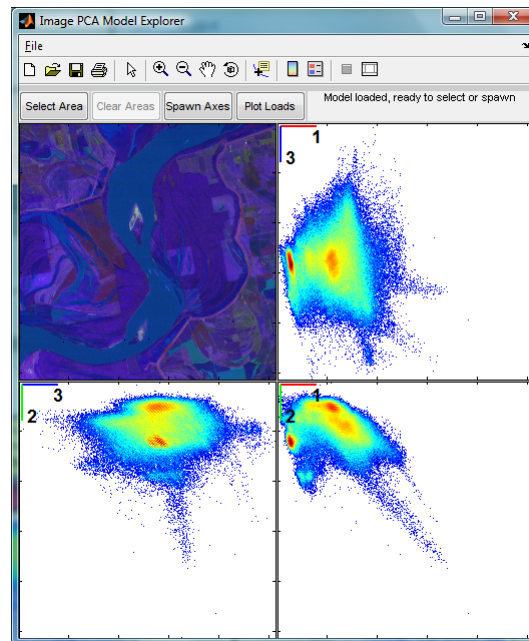
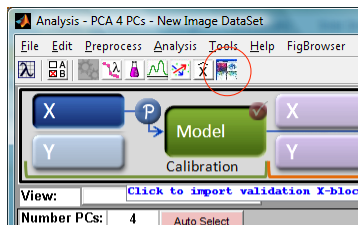
72



## Bivariate Scores Plots

- Plotting  $\mathbf{t}_{k+1}$  vs.  $\mathbf{t}_k$  (score / score plots)
- Problem: lot's of points
  - $512 \times 512 = 262144$  points with lot's of them falling on top of each other (big blobs)
- Density plots
  - count the number of points that lie on top of each other (have same score / score value)
  - color code according to density
  - use log to allow easy comparison between large and small number densities

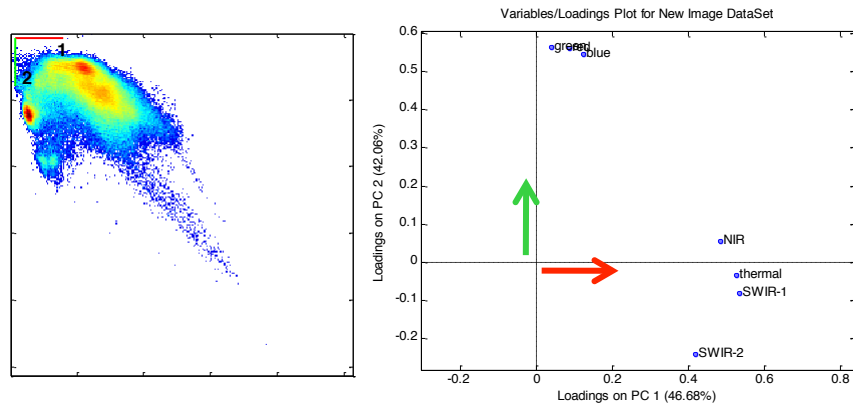
73



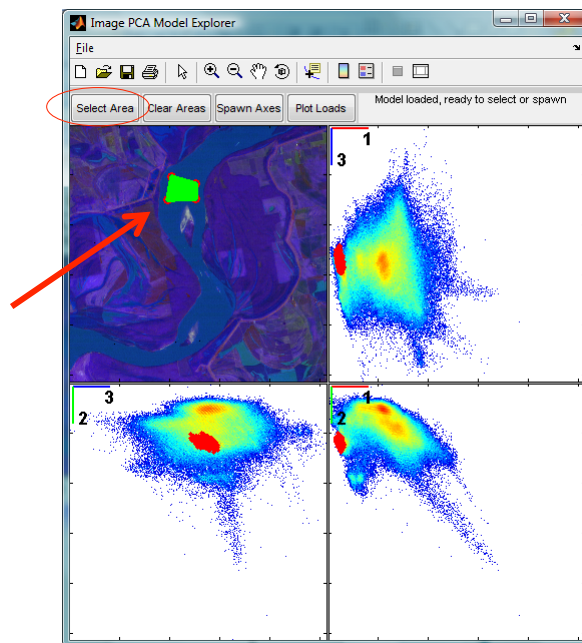
74



## Scores and Loadings

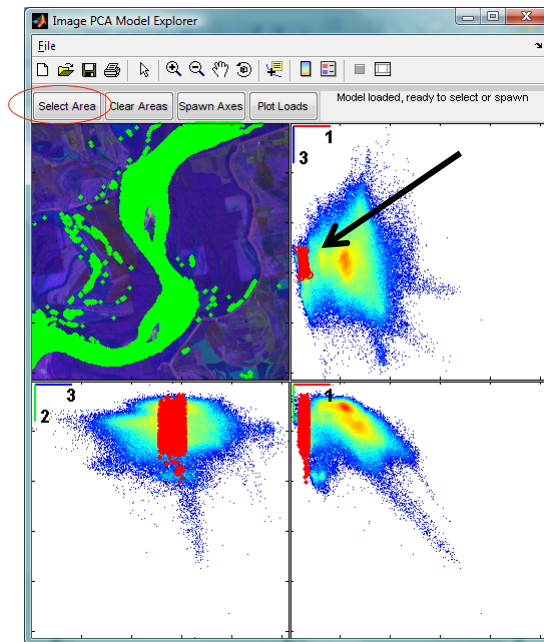


75



selecting an area w/in  
the image plane  
shows where it lies in  
the scores space

76



selecting an area w/in the scores space shows where it lies in the image plane

images can be explored to find similarities and differences w/in an image

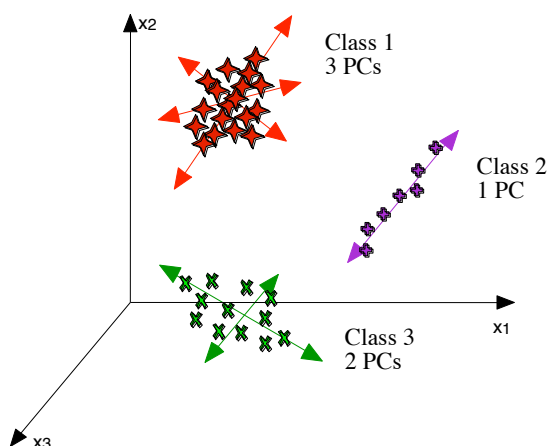
77

## ***SIMCA***

- Supervised pattern recognition / classification technique
  - the model is a collection of PCA models
  - each "class" is a separate PCA model
  - new samples are compared to all of the PCA models and scores,  $T^2$  and  $Q$  are compared to statistical limits on each model
  - samples can belong to one, none or more than one class

78

## A SIMCA Model

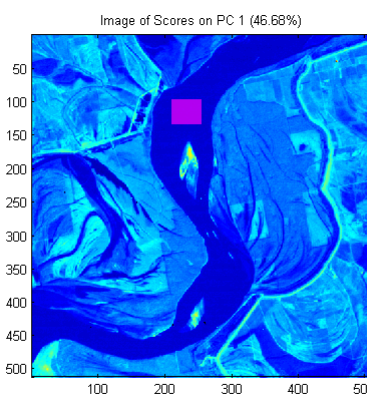
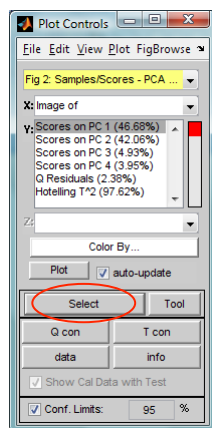


79



## SIMCA Example

For SIMCA, classes need to be defined.  
Use the selection tool to select regions in the image that are expected to be similar and to be modeled as a single class.



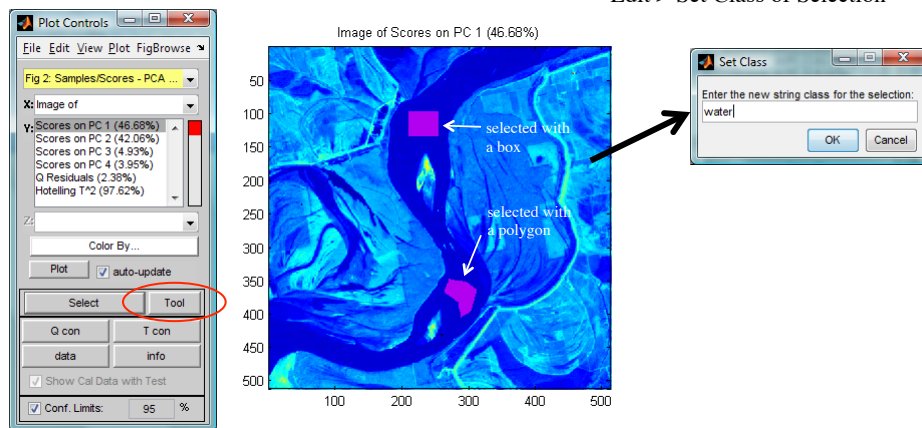
80



## SIMCA Example

- Use the Tool to change the selection tool.
- Hold shift to select multiple regions.

Edit > Set Class of Selection

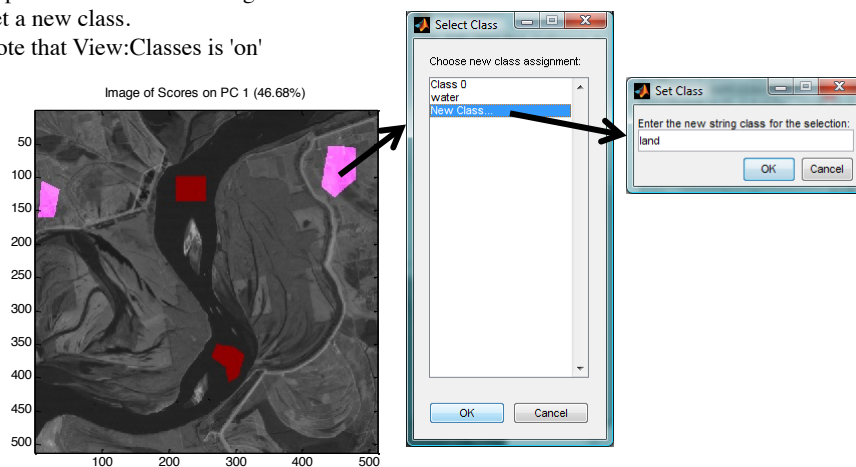


81

**EIGENVECTOR**  
RESEARCH INCORPORATED

## SIMCA Example

- Repeat to select different regions.
- Set a new class.
- Note that View:Classes is 'on'

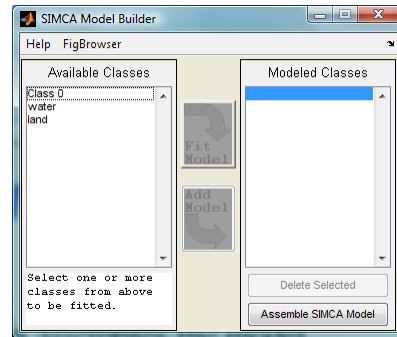
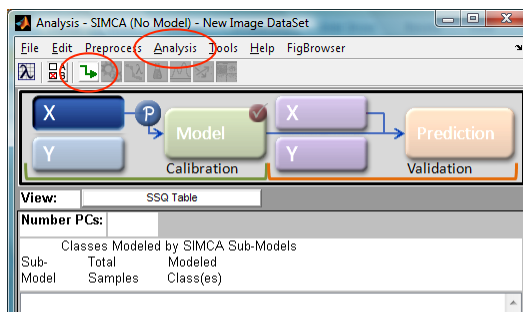


82

**EIGENVECTOR**  
RESEARCH INCORPORATED

## SIMCA Model Builder

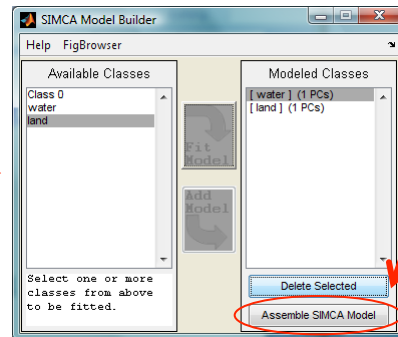
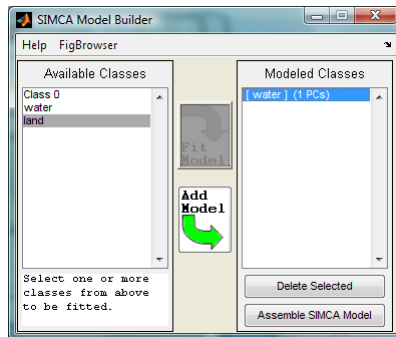
- SIMCA requires a selection of classes to be modeled and then assembles the model
  - Analysis:SIMCA



83

## Model of Each Class

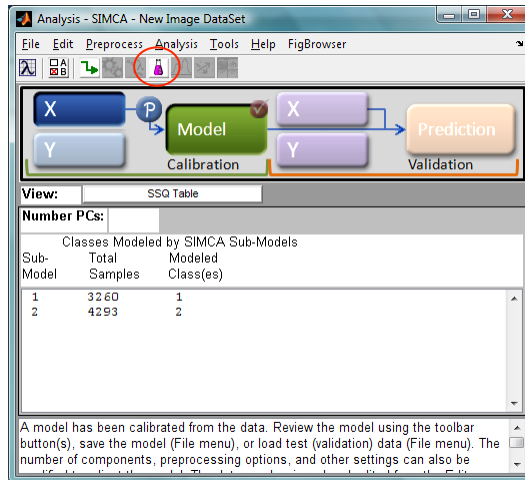
- Each class is modeled using PCA
  - highlight a class and then "fit model"
  - select the number of PCs, etc., then "add model"



84

## SIMCA Example

- The SIMCA model consists of two PCA models

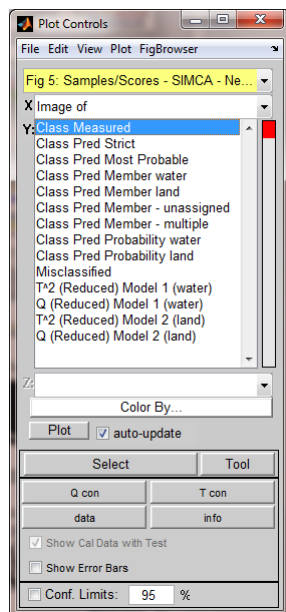


- Data from the entire image will be projected onto each PCA model.
- Scores, Q and  $T^2$  are calculated for each model and it is determined which model the data is closest to.
- Click the scores button to examine the images.

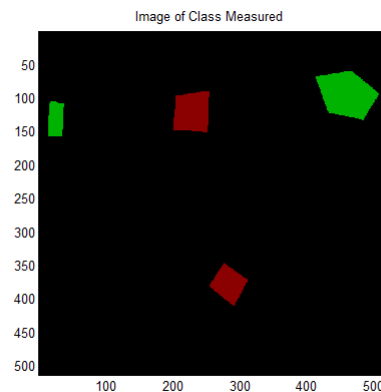
85



## SIMCA Model Predictions



- "Class Measured" = where the classes were selected.
- "Reduced" means that the statistic was normalized by the limit of the corresponding statistic (e.g., to the 95% CL).

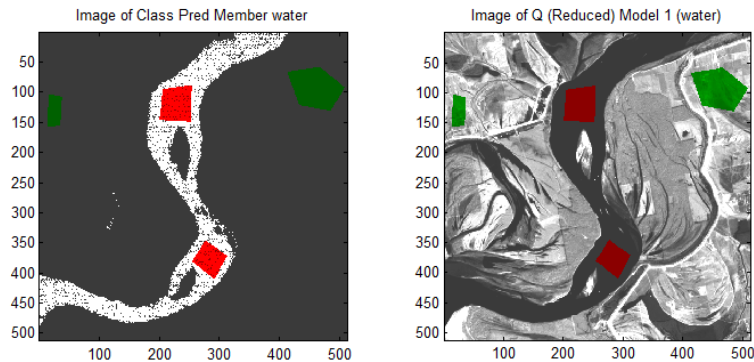


86



## Model 1 Predictions

- Model 1 (w/in set limits for both Q and  $T^2$ )
- Reduced Q on Model 1 (dark is low)

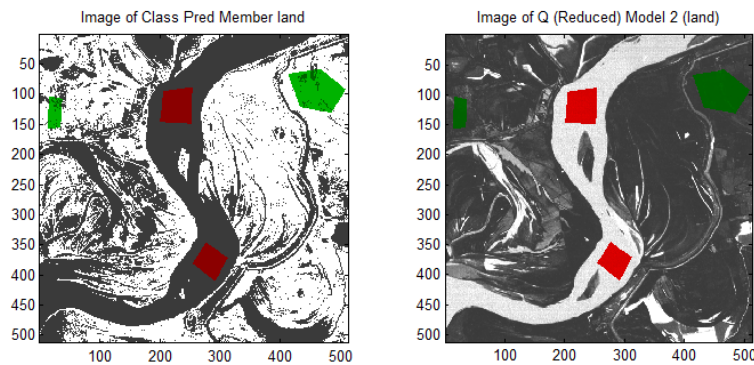


87



## Model 2 Predictions

- Model 2 (w/in set limits for both Q and  $T^2$ )
- Reduced Q on Model 2 (dark is low)



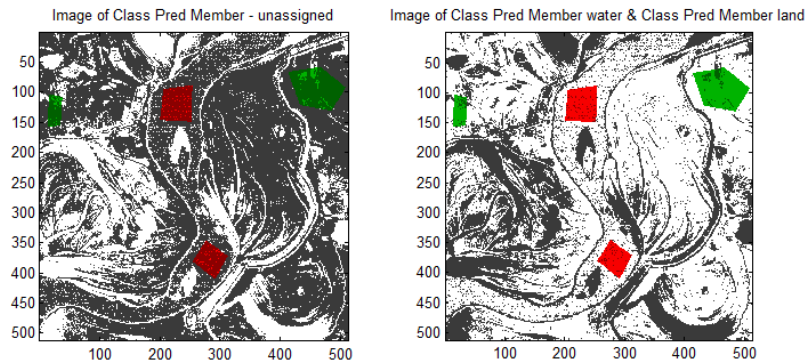
88





## *In Model and Not-In-Any Model*

- Outside of both models (left)
- Inside either model (right)

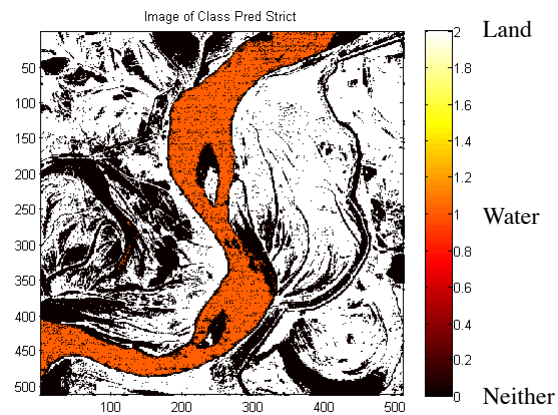


89



## *"Strict" Class Predictions*

- Strict predictions require probability of 50% or greater for one class only
- (Note: turn off classes to view)



90



## Image PCA Conclusions

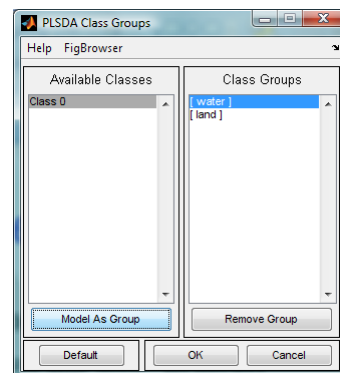
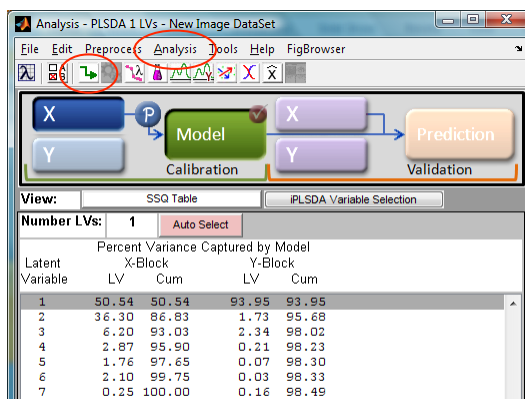
- Image PCA is a useful unsupervised pattern recognition technique for exploring images
  - scores and loadings are useful for determining what original variables are responsible for differences observed in an image
    - score-score plots and linked score plots
    - contrast enhancement might be needed to see small changes
- Image SIMCA is a useful supervised pattern recognition technique
  - find similar / dissimilar portions of an image very quickly

91



## PLSDA Model Builder

- PLS discriminant analysis requires a selection of classes to be modeled
  - Analysis:PLSDA

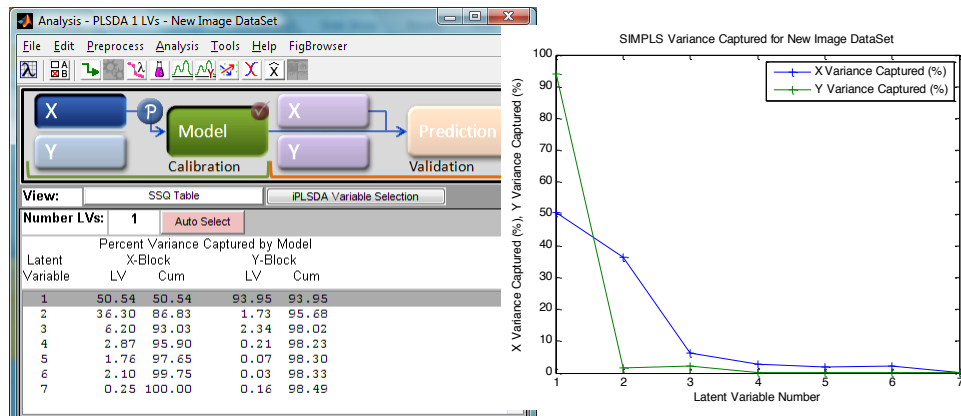


92

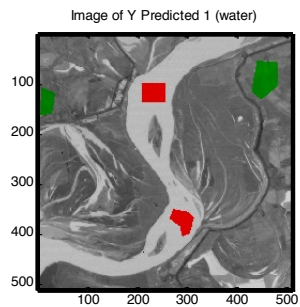


# PLSDA Maximizes Class Separation on a PLS Model

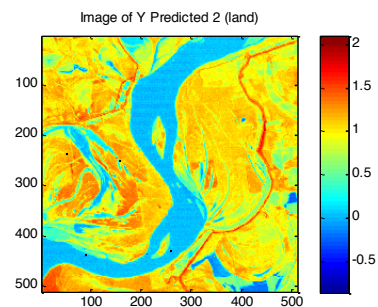
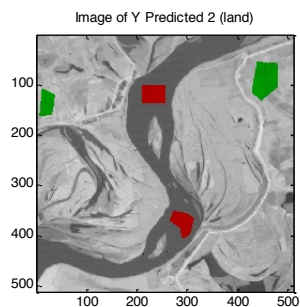
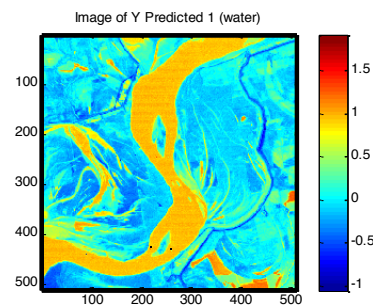
- PLS (selection of factors, cross-validation, etc.)



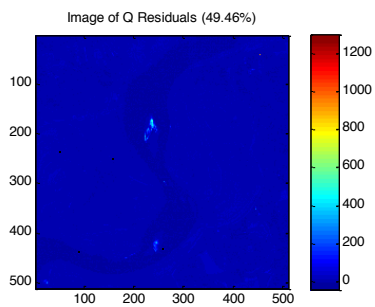
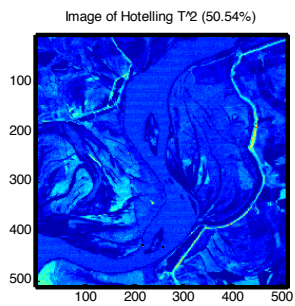
93



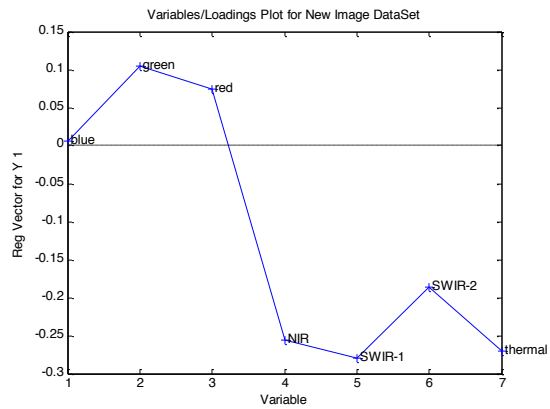
- Data from the entire image are projected onto the PLSDA model.
- Light shows high predictions on each class.
- Click the scores button to examine the images.
- View:Classes (uncheck Set 1)



94



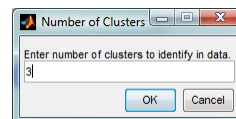
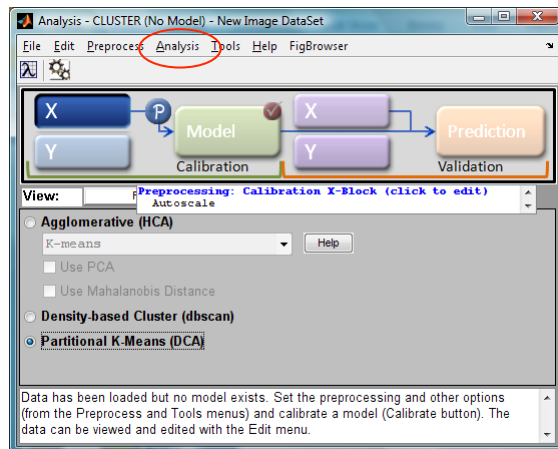
- Inspect  $T^2$  and Q
- Regression vector suggests that green and red increase relative to IR channels for water relative to land



95

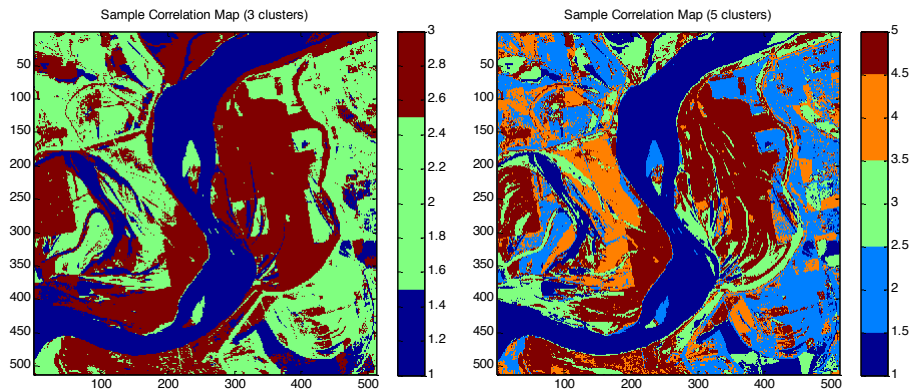
## Cluster Analysis

- Analysis:Cluster



96

## *Results for 3 and 5 Clusters*



97



## *Image PLSDA and Clustering Conclusions*

- If classes (regions) are known, PLSDA is a useful supervised pattern recognition technique for exploring images
  - can often bring out more contrast than PCA
- Image clustering is a useful unsupervised pattern recognition technique (guess number of clusters)
  - find similar / dissimilar portions of an image very quickly
- Results of all analysis methods must be consistent

98



## ***Comments on Presenting Images***

- Images are representations of spatial and chemical information, ...
- but they can be mis-used.
  - users can control colors and contrasting and select channels or PCs (or rotations thereof)
  - as a result some things can be highlighted while others can be hidden
- It is important to report how images were constructed
  - the work must be reproducible

99



## ***Other Ways of Focusing on Variance of Interest***

- Maximum Autocorrelation Factors – find variance with spatial correlation
- Maximum Difference Factors – find variance with spatial transitions (multivariate edge detection)
- Generalized Least Squares Weighting – ignore variance from specified regions

100



## Maximum Autocorrelation Factors for Multivariate Images

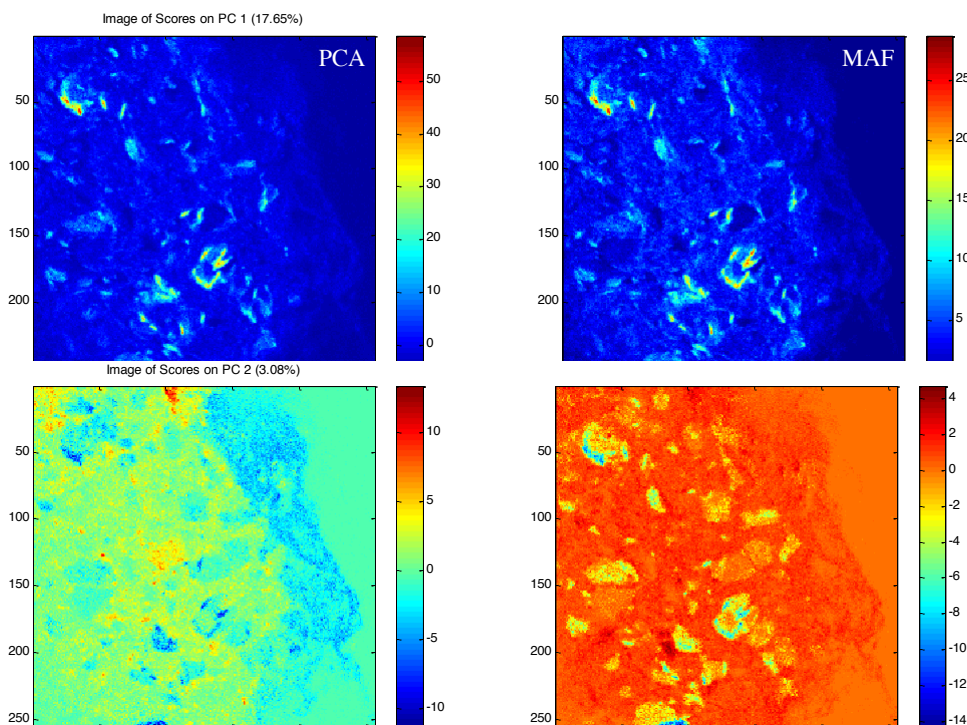
- For MNF, the clutter was intra-class variance
- For MAF, the clutter is the first spatial difference
  - the first difference should be high on edges and just noise w/in clusters
  - the result is the same generalized eigenvector problem as MNF with different clutter  $\Sigma_C$

T.A. Blake, J.F. Kelly, N.B. Gallagher, P.L. Gassman and T.J. Johnson, "Passive detection of solid explosives in Mid-IR hyperspectral images," *Anal Bioanal Chem*, **395**, 337-348, 2009.

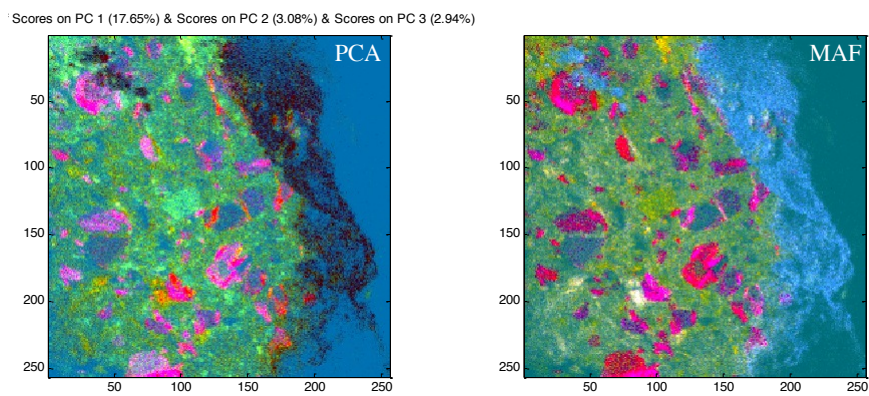
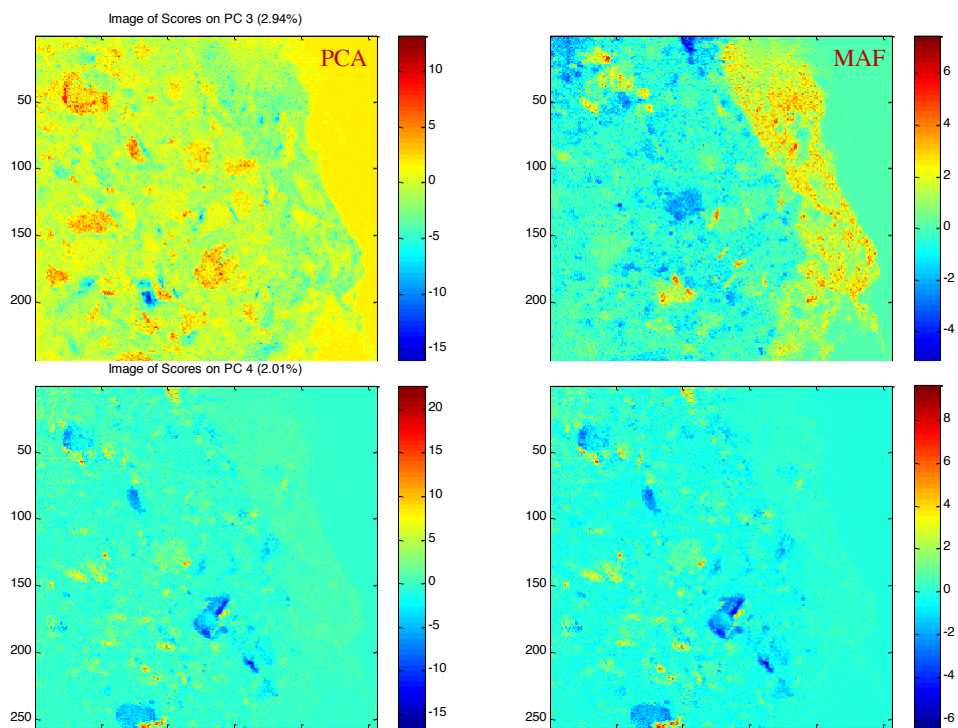
N.B. Gallagher, J.F. Kelly, T.A. Blake, "Passive infrared hyperspectral imaging for standoff detection of tetryl explosive residue on a steel surface," *Whispers* 2010, June 14-16, Reykjavik, Iceland



101







RGB images after auto-contrasting



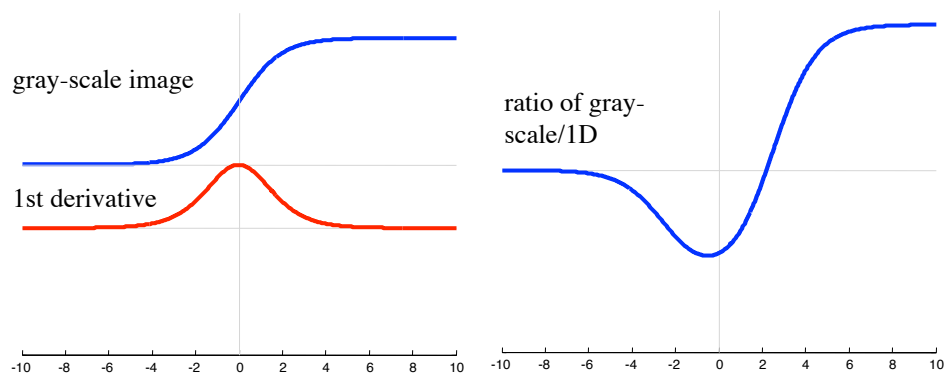
## Maximum Difference Factors MDF

- For MNF and MAF  $\Sigma_X$  was the covariance of the image
- For MAF  $\Sigma_C$  was the covariance of the first spatial difference and in MNF it was estimated from intra-class variance
- For MDF  $\Sigma_X$  is the covariance of the first spatial derivative of image, and  $\Sigma_C$  is the covariance of the second spatial derivative
  - the result is multivariate edge detection
  - often show magnitudes  $\sqrt{dx^2+dy^2}$

105



## MAF

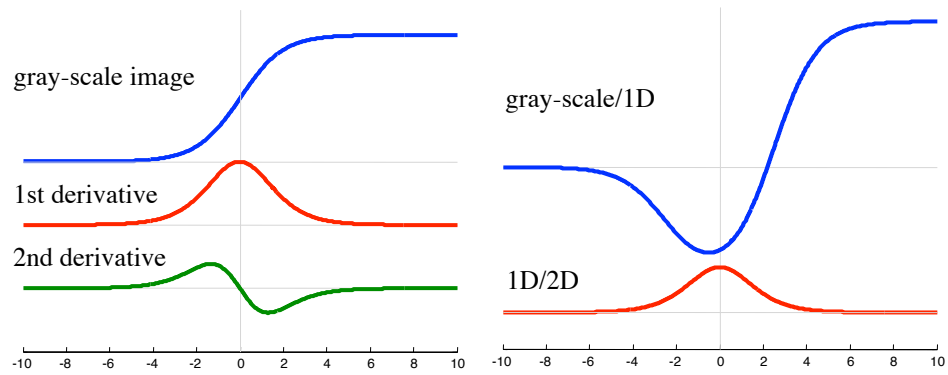


MAF finds locations in the image where the ratio of gray-scale to first derivative is a maximum

106

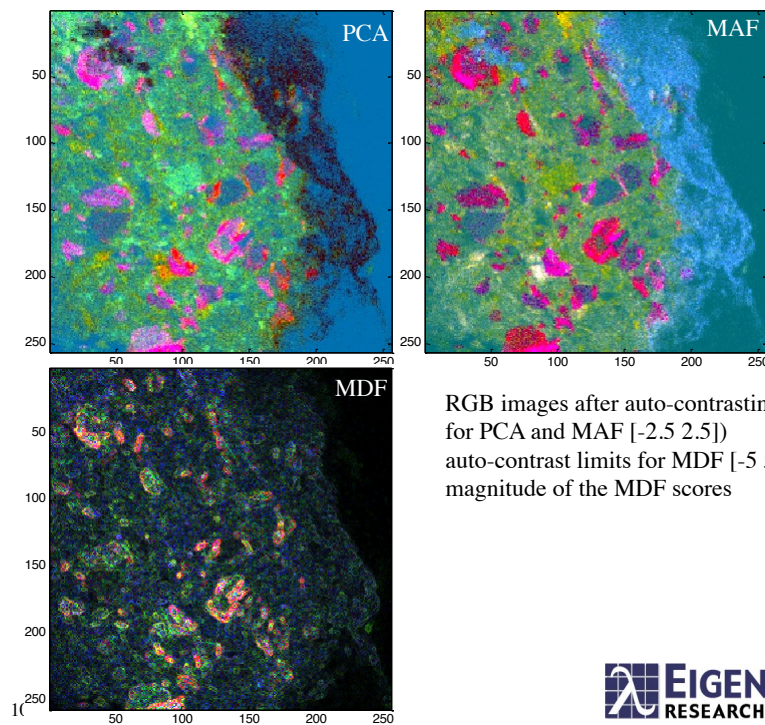


## MDF



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

107



## Measured Signal Includes Clutter

- Clutter is present in all measurements
  - X-block, Y-block



- Use knowledge of physics and chemistry to create a linear relationship
  - non-linearity w/in X-block adds factors in X
  - non-linearity between X- and Y-blocks adds error

109



## Why is Clutter Bad?

- Attempt to maximize S/C via pre-processing or the model e.g., MAF
- Methods that don't remove net analyte signal (NAS) are preferable
  - NAS is the portion of spectrum  $s_i$  unique to analyte  $i$  and orthogonal to all other factors in  $S_{-i}$ , and  $S/C \sim |NAS|$
- Adding clutter tends to add something parallel to  $s_i$  thus lowering NAS
  - Increases estimation error

110



## GLS

- GLS can be used for target detection, classification and quantification
  - need a model of the clutter and a spectrum of pure component(s)
  - no need for buckets of calibration samples
    - in some cases these can't be acquired
  - a.k.a. matched filter and Aitken estimator
    - Turin, George L., "An Introduction to Matched Filters." *IRE Transactions on Information Theory*, 6(3) 1960: 311-329. (this is in a special issue on matched filters) and is used extensively in the remote sensing community [e.g., Burr T, Hengartner N (2006) *Sensors* 6:1721-1750] and has also been referred to as an "adaptive matched filter" to highlight the fact that the clutter covariance can be easily modified resulting in a new filter.
    - Aitken, A., "On Least Squares and Linear Combinations of Observations", *Proceedings of the Royal Society of Edinburgh*, 1935, **55**, 42-48
    - E.g., T.A. Blake, J.F. Kelly, N.B. Gallagher, P.L. Gassman and T.J. Johnson, "Passive detection of solid explosives in Mid-IR hyperspectral images," *Anal Bioanal Chem*, **395**, 337-348, 2009.

111



## GLS Weighting for ILS and PCA

- GLS weighting can be applied to inverse least squares models (e.g., PLS) and PCA

$$\mathbf{X}\mathbf{b} = \mathbf{c} \quad \text{inverse least squares model}$$

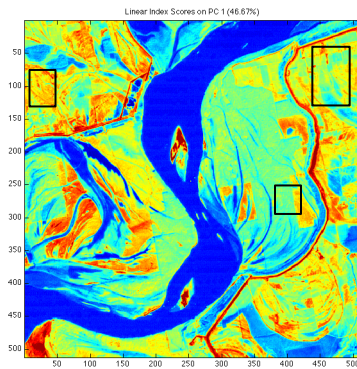
$$\mathbf{X}_w = \mathbf{X}\mathbf{W}_c^{-1/2} \quad \begin{array}{l} \text{weighting of } \mathbf{X} \text{ can be considered a generalization of} \\ \text{autoscaling and is a pre-whitening step} \\ \\ \text{can also be applied to standardization where the clutter} \\ \text{covariance is the difference matrix between instruments} \end{array}$$

H. Martens, M. Høy, B.M. Wise, R. Bro and P.B. Brockhoff, "Pre-whitening of data by covariance-weighted pre-processing," *J. Chemo.*, **17**(3), 153-165 (2003).

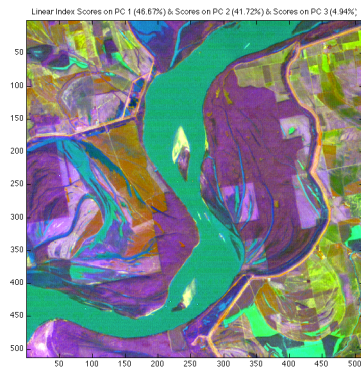
112



## ***Landsat Image of Mississippi***



Scores on PC 1

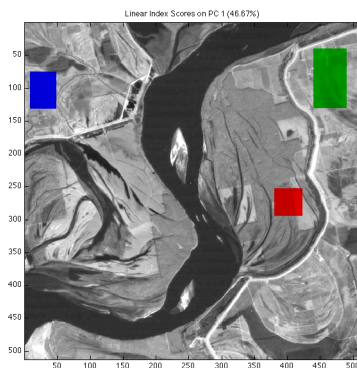


Scores on First 3 PCs

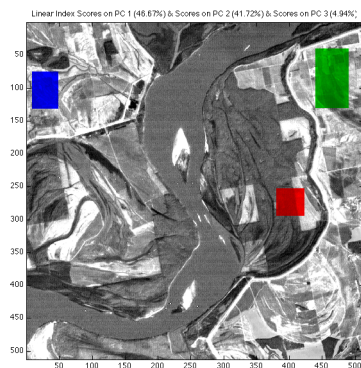
113



## ***Select Classes with Clutter to Down-weight***



Scores on PC 1

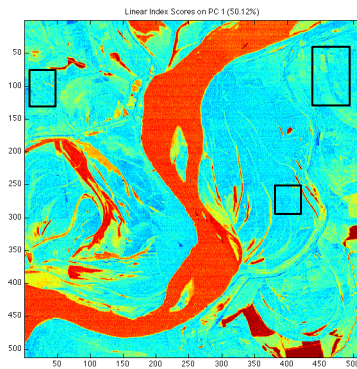


Scores on First 3 PCs

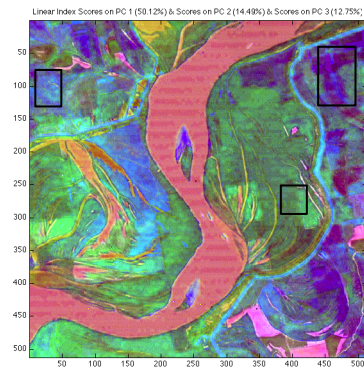
114



## PCA after GLS-Weighting



Scores on PC 1

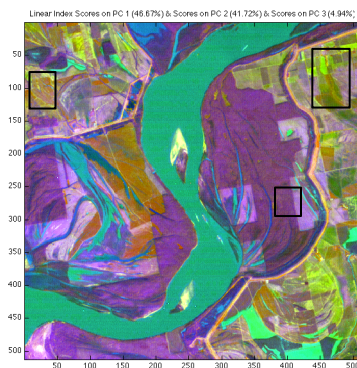


Scores on First 3 PCs

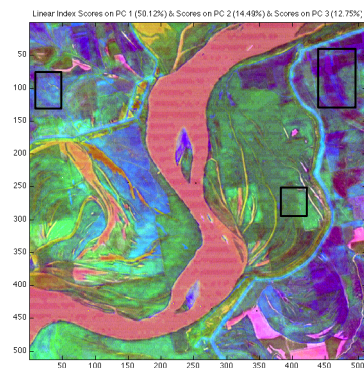
115



## PCA With and Without GLS-Weighting



Without GLS-Weighting



With GLS-Weighting

116





## ***Comments on Filters***

- Savitsky-Golay
  - For derivatives OR smoothing (noise reduction)
- Fourier
  - Remove high-frequency (noise) or low-frequency (baseline) components
  - Typically- NOT “windowed”
    - Position (wavelength) information not considered
- Wavelets
  - Extracting information by BOTH frequency and position
    - Allows BOTH feature selection and pre-processing!
  - filters that are based on window-size (scale)
    - orthogonal and oblique basis functions can be used
- Statistics w/in windows
  - Mean, Median, Max, Min

`SPATIAL_FILTER, LINE_FILTER, BOX_FILTER`

117



## ***Multivariate Image Regression and Quantitative Analyses***

- Inverse Least Squares models (Partial Least Squares – PLS)
- Classical Least Squares (CLS)
- Multivariate Curve Resolution (MCR)

118



## Multivariate Image Regression

- Inverse least squares models
  - PCR, PLS
  - Similar to PCA for X-block
    - matricizing, scores, scores images, loadings, unusual samples Q and T<sup>2</sup>, score-score plots, density plots, linking scores and image plane(s), contrast enhancement
  - Add predictions of a y-block
    - $\mathbf{y} = \mathbf{X}\mathbf{b}$
    - predict a property
    - used for PLS-discriminant analysis

119



## Inverse Least Squares

- Inverse least squares (ILS) models assume that the model is of the form:  $\mathbf{X}\mathbf{b} = \mathbf{y} + \mathbf{e}$ 

where  $\mathbf{y}$  ( $M \times 1$ ) is a property to be predicted,  
 $\mathbf{X}$  ( $M \times N_x$ ) is the measured response,  
 $\mathbf{e}$  ( $M \times 1$ ) is an error vector, and  
 $\mathbf{b}$  ( $N_x \times 1$ ) is a vector of coefficients
- It is possible to estimate  $\mathbf{b}$  from  $\mathbf{b} = \mathbf{X}^+\mathbf{y}$ 

where  $\mathbf{X}^+$  is the pseudo-inverse of  $\mathbf{X}$

120





## Advantage of ILS Methods

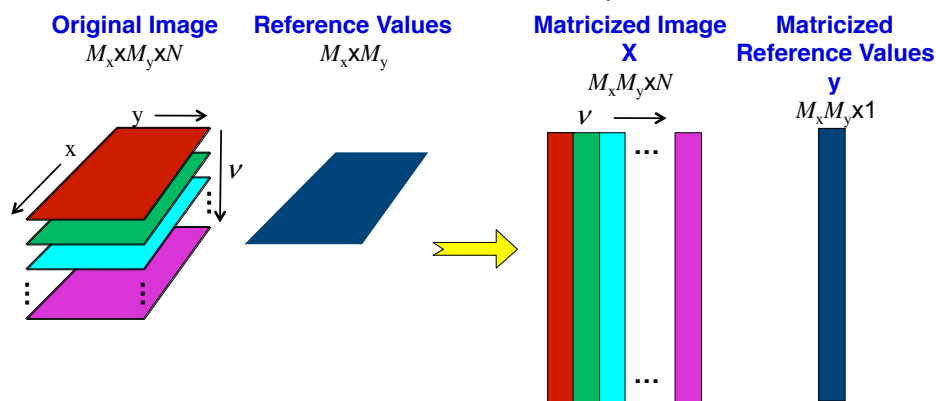
- ILS methods (including MLR, PCR, PLS, CR) don't require the concentration of all analytes, including interferences, be known ...
- ...however, interferences must vary in the calibration data set for the the ILS regression model to be robust against them
  - clutter factors must vary in the calibration data – it's best if they vary such that they are orthogonal to the target of interest
- Disadvantage is that reference values must be available in a representative number of pixels

121



## Unfolding ILS

- The image is  $M_x \times M_y \times N$  and it is reshaped by matricizing such that each pixel is a row in a  $M_x M_y \times N$  matrix



122



## Estimation of b

- There are many ways to obtain a pseudo-inverse
- Multiple linear regression (MLR)<sup>1</sup>  $\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ 
  - problems with rank deficiency and ill conditioning
- Principal components regression (PCR)<sup>2</sup>  $\mathbf{X}^+ = \mathbf{P}_K (\mathbf{T}_K^T \mathbf{T}_K)^{-1} \mathbf{T}_K^T$
- Partial least squares (PLS)<sup>2,3</sup>  $\mathbf{X}^+ = \mathbf{W}_K (\mathbf{P}_K^T \mathbf{W}_K)^{-1} (\mathbf{T}_K^T \mathbf{T}_K)^{-1} \mathbf{T}_K^T$ 
  - cross-validation used to select number of factors

<sup>1</sup>Draper, N. and Smith, H., "Applied Regression Analysis, Second Edition", John Wiley & Sons, New York, NY (1981).

<sup>2</sup>Martens, H. and Næs, T., "Multivariate Calibration", John Wiley & Sons, New York, NY (1989).

<sup>3</sup>M. Andersson, "A comparison of nine PLS1 algorithms," *J. Chemom.*, **23**(10), 518-529 (2009)

## Model Performance Measures

- Average measures
  - Root mean square error of calibration (RMSEC)  $\text{RMSEC} = \sqrt{\frac{\sum_{m=1}^M (y_m - \hat{y}_m)^2}{M}}$ 
    - measure of fit
  - RMSE cross-validation (RMSECV)  $\text{RMSECV}_K = \sqrt{\frac{\sum_{j=1}^J \sum_{m=1}^{M_j} (y_m - \hat{y}_m)^2}{M}} = \sqrt{\frac{\text{PRESS}_K}{M}}$ 
    - approximate measure of prediction
  - RMSE prediction (RMSEP)  $\text{RMSEP} = \sqrt{\frac{\sum_{m=1}^{M_p} (y_m - \hat{y}_m)^2}{M_p}}$ 
    - measure of prediction
- Estimation error includes leverage

Faber, N.M. and Bro, R., *Chemomem. and Intell. Syst.*, 61, 133-149 (2002).

## ***For PCR and PLS: Number of PCs or LVs***

- Choice is not always simple
- A few rules of thumb
  - $\sqrt{M}$  a good choice for number of splits
  - useful to do repeated CVs with different data ordering
    - want subsets to span the data space
  - be conservative, models are more often overfit than underfit
  - best choice is often not the global minimum PRESS
  - look for minimum of PRESS and work backwards if improvement is not at least 2%
  - $RMSEC < RMSECV$  by more than ~20% indicates overfit
  - look at variance captured in **X** and **Y**. Is it significant with respect to what you know about the data?

125



## ***Model Diagnostics***

- Diagnostics useful for finding outliers/uniques
- **X**-block Q residual and  $T^2$
- **X**-block leverage and studentized **Y**-block residuals

126



## Unfold ILS Comments

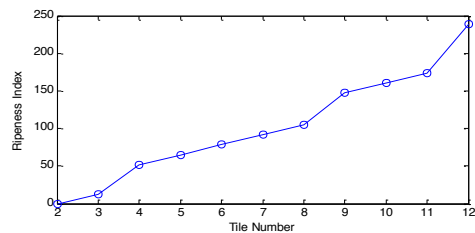
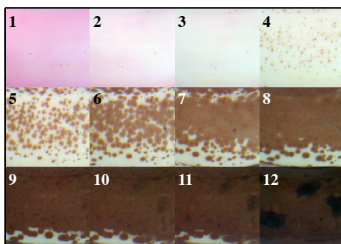
- Easy, can use existing code with rearranged data
- Statistics reasonably well defined
- No second order advantage to be lost!
  - images are multi-mode but the spatial mode is not bi-linear

127



## Banana Ripeness by PLS

- Goal: Develop an automated (objective) method to assess banana ripeness
- X-Block RGB Images of Bananas at various stages of ripeness (Tiled)
- Y-Block Ripeness index for each tile

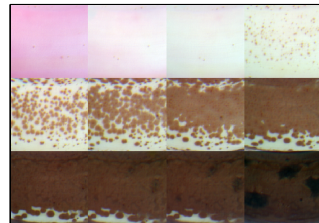


Data Courtesy Kim Esbensen  
University of Ålborg, Denmark

128



## Two-Dimensional Calibration Data

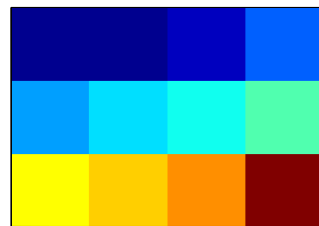


X-block

Image-based calibration takes advantage of high sampling rate of imaging (40 thousand samples for each tile!)

Y-block assumes a constant reference value for each image.

Unfold blocks before PLS



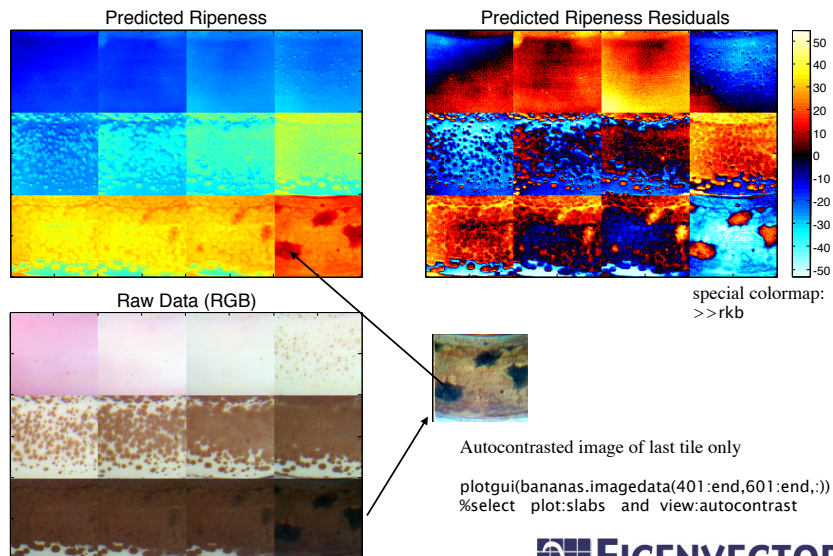
Y-block

Note: Does not inherently take spatial correlation into account.

129



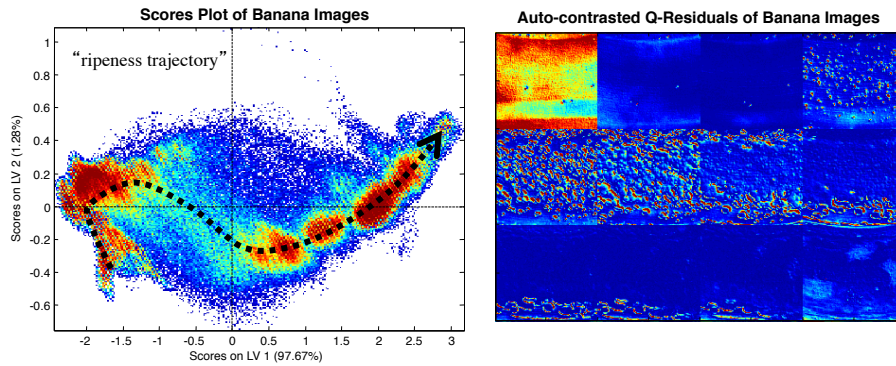
## Banana Predictions



130



## Banana Scores and Q Residuals



131

## Classical Least Squares Models

- Classical Least Squares (CLS)
  - alternative to ILS models often used in imaging where reference values are unknown but target spectrum is known
  - extended mixture model, generalized least squares
    - accounting for clutter [unknown (but characterizable) interferences] in CLS
  - often used in spectroscopic applications and remote sensing

132

## CLS Models

- Classical Least Squares (CLS)

- also uses 'unfolded image'

$$\mathbf{X} = \mathbf{C}\mathbf{S}^T + \mathbf{E} \quad \hat{\mathbf{C}} = \mathbf{X}\mathbf{S}(\mathbf{S}^T\mathbf{S})^{-1}$$

- requires spectrum of all chromophores
  - often cited as reason for ILS

- Extended Mixture Model (ELS)

$$\mathbf{X} = [\mathbf{C} \quad \mathbf{T}][\mathbf{S} \quad \mathbf{P}]^T + \mathbf{E} \quad \begin{bmatrix} \hat{\mathbf{C}} & \hat{\mathbf{T}} \end{bmatrix} = \mathbf{X}[\mathbf{S} \quad \mathbf{P}]\left([\mathbf{S} \quad \mathbf{P}]^T[\mathbf{S} \quad \mathbf{P}]\right)^{-1}$$

- where  $\mathbf{P}$  is a sub-space that span the systematic clutter variance

- Generalized Least Squares (GLS)

$$\mathbf{X} = \mathbf{C}\mathbf{S}^T + \mathbf{E} \quad \hat{\mathbf{C}} = \mathbf{X}\mathbf{W}^{-1}\mathbf{S}(\mathbf{S}^T\mathbf{W}^{-1}\mathbf{S})^{-1}$$

- where  $\mathbf{W}$  is the clutter covariance (might center also)
  - requires characterization of clutter
  - similar to requirement that interferences vary in ILS

133



## MCR

- Based on the classical least squares (CLS) model, attempt to estimate  $\mathbf{C}$  and  $\mathbf{S}$  given  $\mathbf{X}$ :

$$\mathbf{X} = \mathbf{C}\mathbf{S}^T + \mathbf{E}$$

where

$\mathbf{X}$  is a  $M \times N$  matrix of measured responses,

$\mathbf{C}$  is a  $M \times K$  matrix of pure analyte contributions,

$\mathbf{S}$  is a  $N \times K$  matrix of pure analyte spectra, and

$\mathbf{E}$  is a  $M \times N$  matrix of residuals.

134



## *MCR Objective*

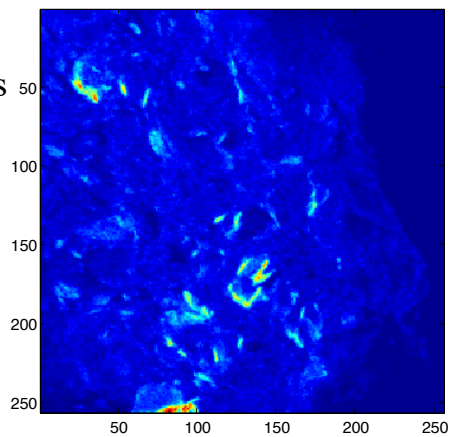
- Decompose a data matrix into chemically meaningful factors
  - pure analyte spectra
  - pure analyte concentrations
- Easy to interpret
  - provides chemically / physically meaningful information
  - caveats:
    - rotational and multiplicative ambiguity
    - use of constraints

135



## *Imaging Mass Spec*

- Image is 256x256x90
- The mass spectrum was 41945 mass channels selected and binned into 93 channels
- Image of total ion count
  - false color



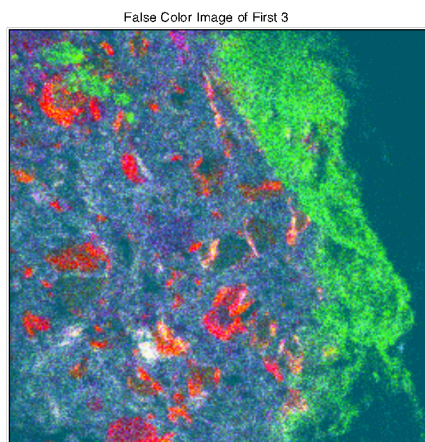
136





## PCA Score Image

Pretty picture,  
but loadings are  
very difficult to  
interpret!



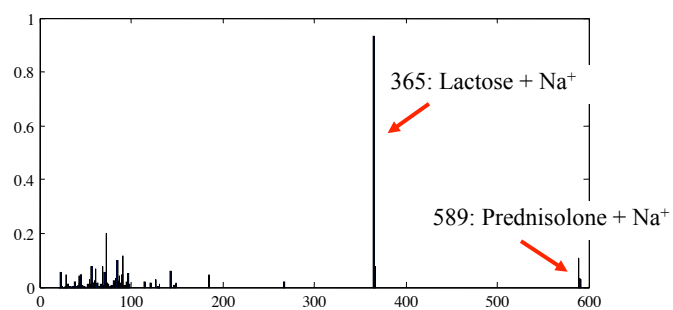
137

## MCR (ALS) on TOF-SIMS Image

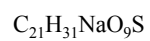
- Non-negative constraints on both C and S
- Initialize with pure/extreme samples (i.e. pixels)
- Recover 6 interpretable spectra and concentration profiles
- Showing Score Images – image was unfolded with each pixel as a separate sample then the scores are re-folded to form images

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).

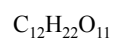
138



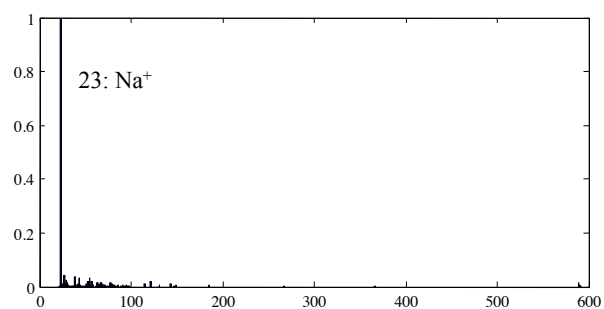
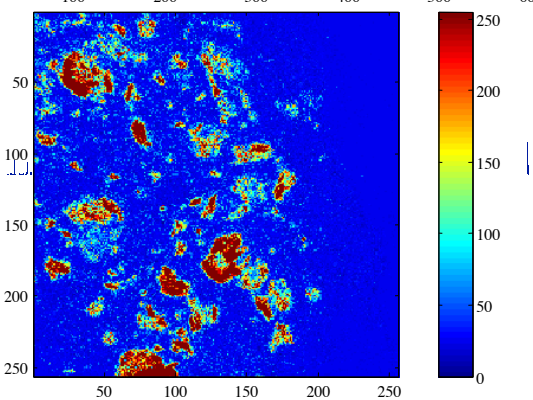
Prednisolone:



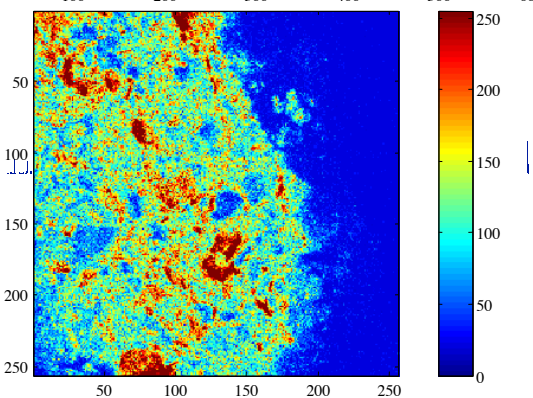
Lactose:

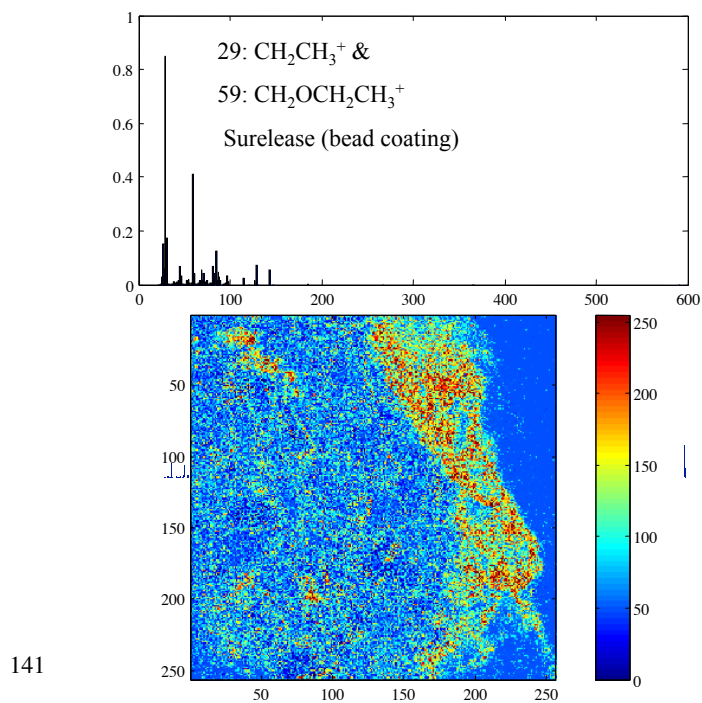


139



140

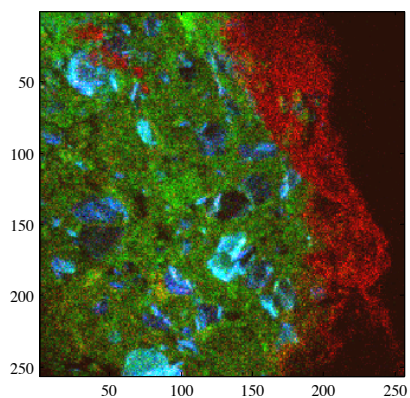




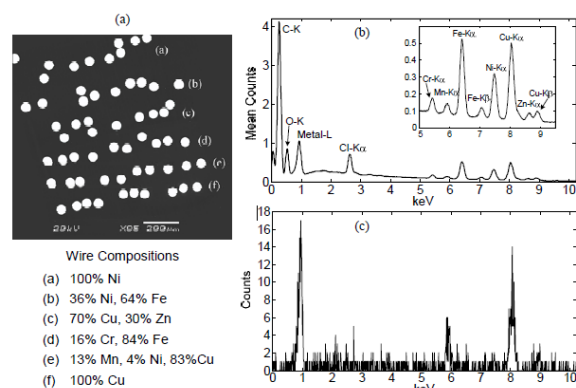
## *RGB “Chemical” Image*

Red: Surelease (bead coating)  
Green: Na  
Blue: Prednisolone (drug)

only 3 of 6 factors extracted are  
shown



## Energy Dispersive Spectrometry (EDS) Image of Wires

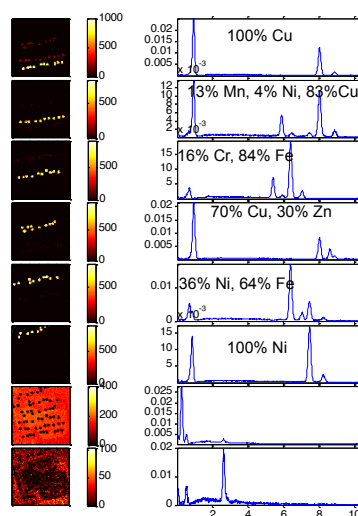


M.R. Keenan, Multivariate Analysis of Spectral Images Composed of Count Data, In: H. F. Grahn, P. Geladi (eds.), Techniques and Applications of Hyperspectral Image Analysis, pp. 89-126, Wiley & Sons, 2007

143



## MCR Results on Wires



144

Using Image Contrast (angle) constraint



## Example of Dealing w/ Clutter

- MIA Example: Multivariate Curve Resolution (MCR)
  - Perform EMSC magnitude and slope correction (more later ...)
    - reference is an estimate of the resin spectrum with robust fitting
    - allow glucose, lysine,  $\text{CaSO}_4$  spectra to pass the filter
    - Gallagher, Blake, Gassman, *J. Chemometr.*, **19**(5-7), 271-281 (2005).
  - Step 2: Account for scratches using spatial constraints:
    - Scores from a PCA of region 2778 to 1790  $\text{cm}^{-1}$  w/ 2<sup>nd</sup> derivative preprocessing capture variability due to scratch features
    - Equality constraints on **C**: components 4 to 11 → the scratches
      - Soft equality Constraints on **S**: components 1 to 3
        - » Factor 1: resin, Factor 2: lysine (w/~  $\text{CaSO}_4$ ), Factor 3: glucose
- Linear mixture model referred to as an extended mixture model

$$\mathbf{X} = [\mathbf{C} \quad \mathbf{T}] [\mathbf{S} \quad \mathbf{P}]^T + \mathbf{E}$$

desired factors
interferences

145

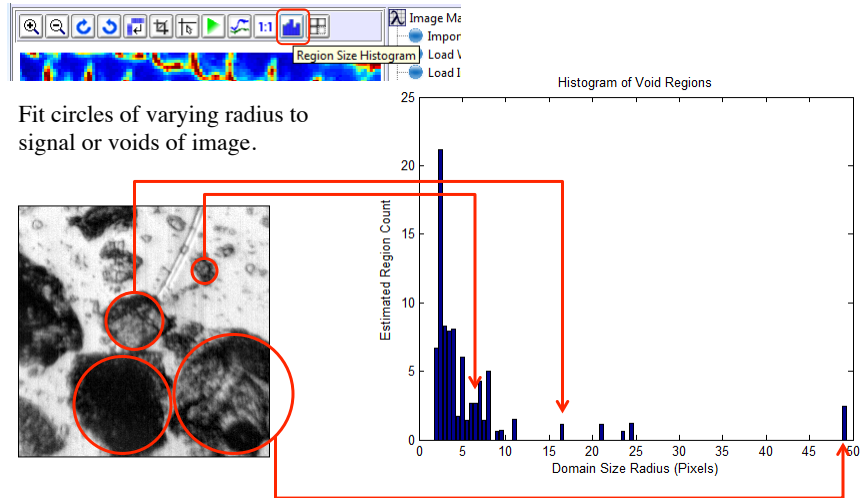


## OTHER TOOLS

146



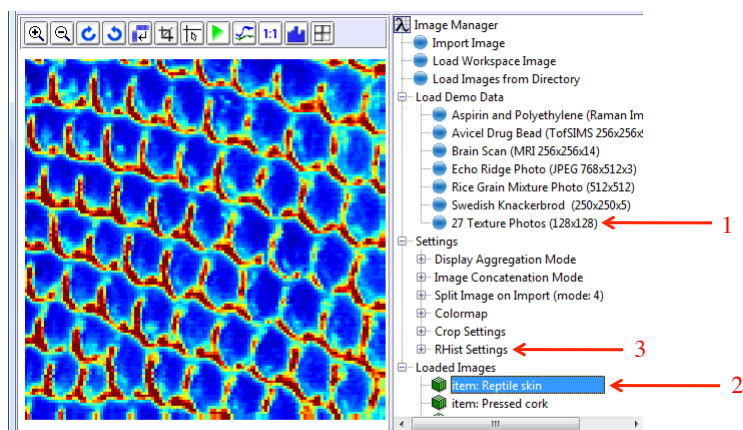
## Radial Region Histograms



147



## Region Histogram Example: Reptile Skin



148



## Region Histogram Settings

Option Name	Value
Display	
plots	final
Standard	
units	radius
space	void
minsize	2
stepsize	0.5
inthreshold	
dilatefill	on
buffer	10

### Plots –

- Final (only histogram)
- Detailed (includes filling map)

### Units –

- Radius, Diameter, Area

### Space –

- Signal: measure signal

- Void: measure lack of signal ←

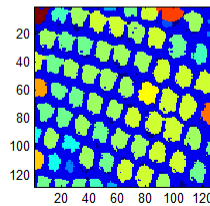
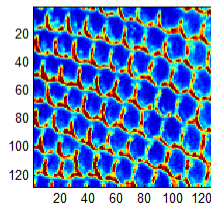
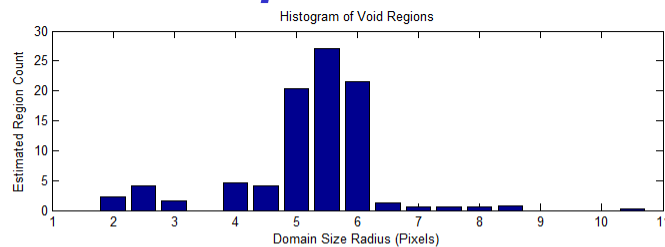
### Dilatefill –

- (On/Off) accommodates non-circular regions by adjusting circle to fill space

149



## Region Histogram Example: Reptile Skin



Hint: Right-click axes to spawn as separate figure.

150



## *Filled Region Map*

