

Chemometrics II: Regression and PLS (Building Predictive Models)

©Copyright 1996-2017
Eigenvector Research, Inc.
No part of this material may be
photocopied or reproduced in any form
without prior written consent from
Eigenvector Research, Inc.



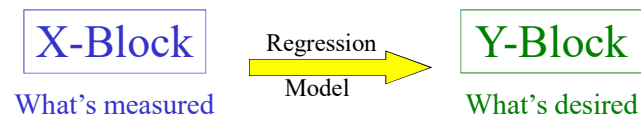
What Can Be Done with Regression?

- Analyte concentrations from spectra or other sensors
 - CH₄, H₂O, CO₂ in natural gas (NIR)
 - H₂, NH₃ in waste tanks (FTIR)
 - Sugar content of fruit (NIR)
- Prediction of property values
 - Octane of gasoline (NIR)
 - Ozone forming potential of automobile exhaust (FTIR)
- Sample classification (e.g., PLS-DA)
 - Detection of cervical cancer (ETF)
 - Detection of atherosclerotic (vulnerable) plaques (NIR)



3

Regression



Regression analysis creates a mapping between two blocks of data.

In contrast, PCA was used to explore the correlation structure within a single data block.

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



2

Outline

- Introduction
- Classical Least Squares (CLS)
- Inverse Least Squares (ILS) Models
- Multiple Linear Regression (MLR)
- Ridge Regression (RR)
- Principal Components Regression (PCR)
 - Cross-validation
- Partial Least Squares Regression (PLS)
 - Model Quality Measures
 - Determining of the Number of factors
 - Outlier Detection and Model Diagnostics
- Comparison of Methods on NIR Styrene Copolymer data
- A Unifying Theme: Continuum Regression (CR)
- Model Updating, Missing Data
- Summary



4

Course Materials

- These slides
- PLS_Toolbox or Solo 8.1 or later
- Data sets
 - From DEMS folder (installed with software)
 - plsdata (SFCM)
 - From EVRIHW folder (additional data sets)
 - EigenU_nir_data, SBRdata_EU

5



Data Preprocessing

- Everything that was said about preprocessing for PCA goes double for regression
- Data should be linearized, if possible
- Data is often mean-centered
- Variance scaling used when variables are in different units or greatly different magnitudes
- Many preprocessing methods available!
 - Goal: reduce extraneous variance, emphasize relevant variance
- Outlier elimination is critical to regression models

7



Conventions & Notation

- Rows correspond to *samples*, columns correspond to *variables*
- Notation:
 - \mathbf{X} = matrix of predictor variables
 - \mathbf{Y} = matrix (or vector \mathbf{y}) of predicted variables
 - M = number of samples (observations)
 - N_x = number of \mathbf{X} variables, N_y = number of \mathbf{Y} variables
 - \mathbf{T} = \mathbf{X} -block scores matrix, $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_K$ score vectors
 - \mathbf{U} = \mathbf{Y} -block scores matrix, $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K$ score vectors
 - \mathbf{P} = \mathbf{X} -block loads matrix, $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_K$ loadings vectors
 - \mathbf{Q} = \mathbf{Y} -block loads matrix, $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_K$ loadings vectors
 - \mathbf{W} = \mathbf{X} -block weights matrix, $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$ loadings vectors
 - Θ = ridge parameter

6



Classical Least Squares

- CLS can be used to develop calibration models
 - often used in spectroscopy
- The CLS model assumes the data follows:

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

where \mathbf{X} ($M \times N_x$) is the measured response, \mathbf{S} ($N_x \times K$) is a matrix of pure component responses, \mathbf{C} ($M \times K$) is a matrix of weights (e.g., concentrations) and \mathbf{E} ($M \times N_x$) is noise or an error matrix.

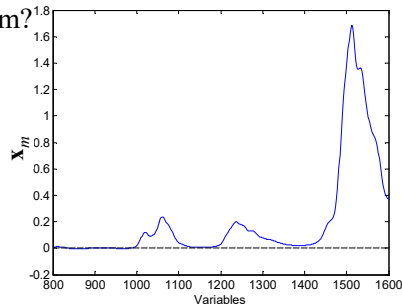
8



The CLS Model

- Given known pure component spectra, how much of each does it take to make up the observed m^{th} spectrum?

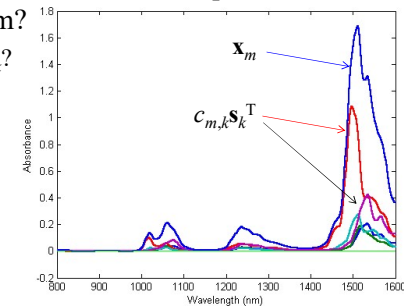
- $\mathbf{x}_m = \mathbf{c}_m \mathbf{S}^T + \mathbf{e}_m$
- $m = 1, \dots, M$
- $\mathbf{c}_m = [c_{m,1}, c_{m,2}, \dots, c_{m,K}]$
- $k = 1, \dots, K$



The CLS Model

- Given known pure component spectra \mathbf{S} , how much of each does it take to make up the observed spectrum?

- i.e.*, what are the $c_{m,k}$?



CLS (cont.)

- Once \mathbf{S} (the spectral “basis”) is known, \mathbf{c} , the degree to which each component contributes to a new sample \mathbf{x} , can be determined from

$$\mathbf{c} = \mathbf{x} \mathbf{S}^+$$

where \mathbf{S}^+ is the pseudo-inverse of \mathbf{S} , defined in CLS as

$$\mathbf{S}^+ = \mathbf{S}(\mathbf{S}^T \mathbf{S})^{-1}$$

- Problem: How to get \mathbf{S} ?
 - library, estimate from calibration measurements

Classical Least Squares

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

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

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

$$\mathbf{X} \mathbf{S} (\mathbf{S}^T \mathbf{S})^{-1} = \mathbf{C}$$

$$\mathbf{S}^+ = \mathbf{S} (\mathbf{S}^T \mathbf{S})^{-1}$$

- Note that $\mathbf{S}^T \mathbf{S}$ is $K \times K$ (analytes by analytes) and square

Estimating S

- Sometimes, **S** can be compiled *a priori* from a data base/spectral library, or from direct measurements of pure components
 - Problem: must account for all components that can contribute to **X**!
- S** can also be estimated from mixtures, provided all **C** are known and enough samples are available:

$$S^T = (C^T C)^{-1} C^T X$$

- Problem: The concentration of *every* analyte that contributes to **X** must be known!*

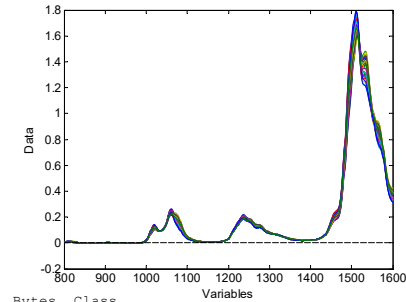
*Interferences and unknowns can be handled with GLS or ELS type models, but their basis must be estimated



13

CLS Example

- NIR data of pseudo-gasoline samples
 - absorbance at 401 channels
 - 30 samples
 - 5 analytes
- EigenU_nir_data.mat
- Data broken into
 - 25 calibration samples and
 - 5 test samples

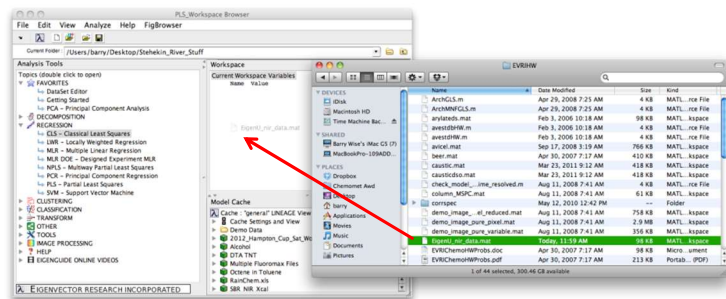


```
>> load EigenU_nir_data
>> whos
Name                Size      Bytes  Class
cal_conc             25x5       11002  dataset
cal_spec            25x401    96466  dataset
test_conc             5x5        10042  dataset
test_spec            5x401     32146  dataset
```



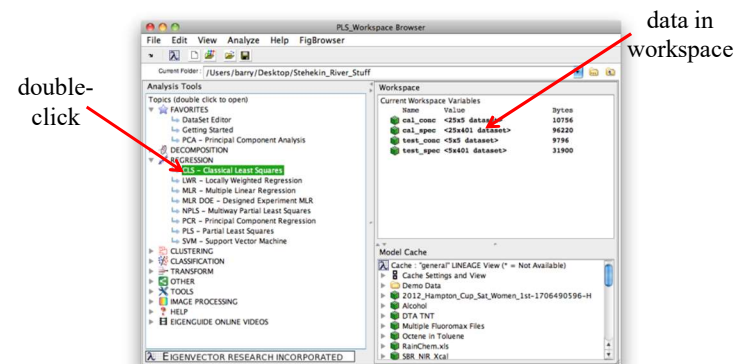
14

Load Data Into Browser



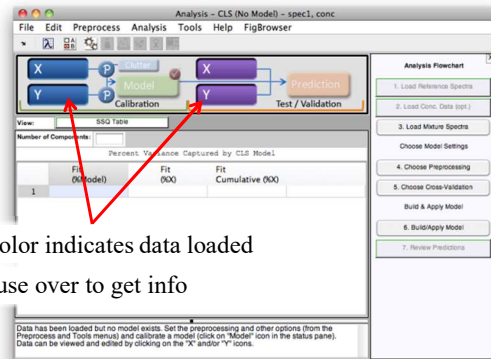
15

Start CLS Interface



16

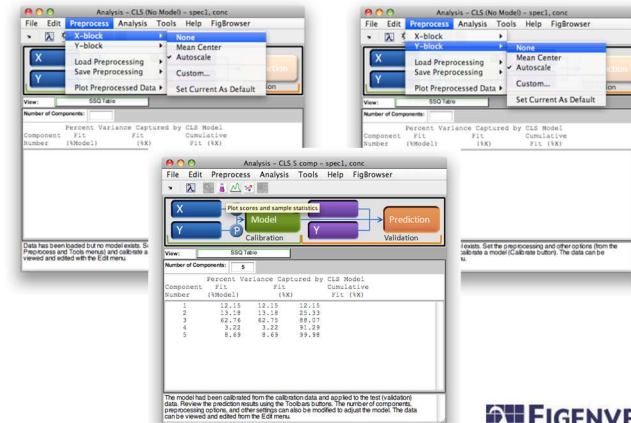
Data Loaded



17

EIGENVECTOR
RESEARCH INCORPORATED

Set Preprocessing to "none," calculate model

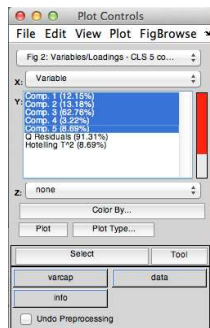


18

EIGENVECTOR
RESEARCH INCORPORATED

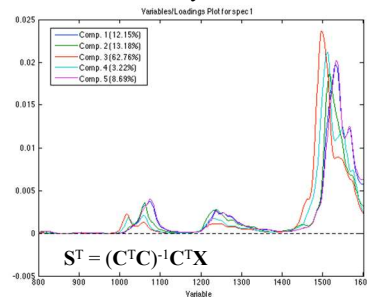
Pure Component Spectra

Click loadings "spectrum" icon, select all 5 components



19

S, estimated from mixtures, using known concentrations of all 5 analytes



Fit to Calibration and Estimate for Validation Samples

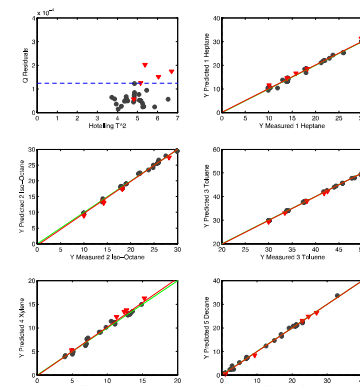
Click scores "flask" icon to get fits and predictions (test set).

Check "Show Cal Data with Test".

Calibration data (black)

Predicted test (red).

All analytes fit and predicted well.



20

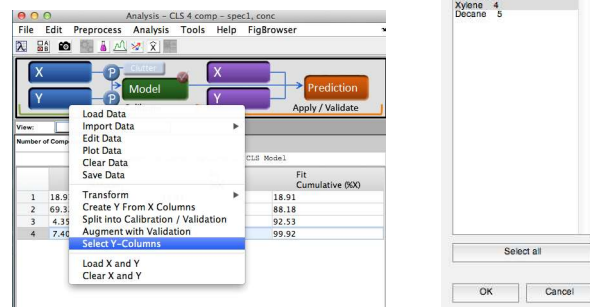
EIGENVECTOR
RESEARCH INCORPORATED

CLS Problem

- What if the concentration of 1 analyte was unknown?
- Repeat the CLS procedure using only the first 4 (of 5) analytes
- Attempt to predict concentrations of unused (test) samples

21

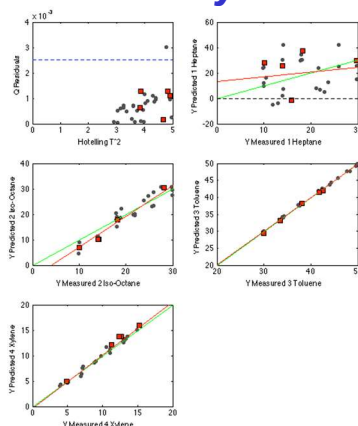
Select only the first four analytes and repeat



click 'cal Y: select Y-columns'

22

CLS Solution with One Analyte "Missing"



Click scores "flask" icon to get fits

Some analytes not fit (black) and not predicted (red) well, especially heptane

23

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

- Unlike CLS, ILS methods associate the noise with the predicted property, not the measured response

24

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 ILS regression model to be robust against them

Interferent: Any substance whose presence interferes with an analytical procedure and generates incorrect results (wiktionary)

25



Estimation of \mathbf{b} : MLR

- It is possible to estimate \mathbf{b} from

$$\mathbf{b} = \mathbf{X}^+ \mathbf{y}$$

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

- There are many ways to obtain a pseudo-inverse; the most obvious is multiple linear regression (MLR), a.k.a., Ordinary Least Squares (OLS)
- In this case, \mathbf{X}^+ is estimated from

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

26



Multiple Linear Regression

$$\mathbf{X}\mathbf{b} = \mathbf{y} + \mathbf{e}$$

$$\mathbf{X}\mathbf{b} = \mathbf{y}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

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

- Note that $\mathbf{X}^T \mathbf{X}$ is $N_x \times N_x$ and square

27



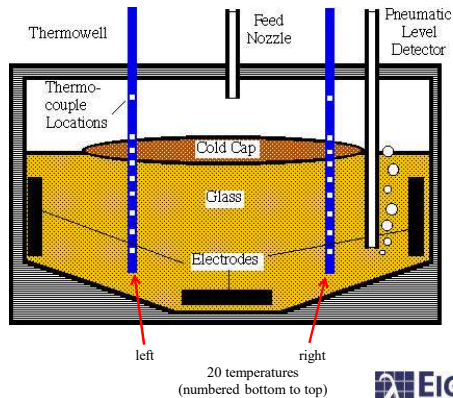
Problem with MLR

- Inverse of $\mathbf{X}^T \mathbf{X}$ only exists if ...
 - $\text{Rank}(\mathbf{X}) = N_x$, however $\text{rank}(\mathbf{X}) \leq \min(M, N_x)$
 - \mathbf{X} has more samples than variables *i.e.*, if $M > N_x$, and
 - problem with spectra
 - Columns of \mathbf{X} are not co-linear.
- Inverse may exist but be highly unstable if \mathbf{X} is nearly rank deficient (a.k.a., ill-conditioned).
- In these cases, small perturbations in the data (possibly due to noise) can produce very different results.

28



Slurry Fed Ceramic Melter: SFCM

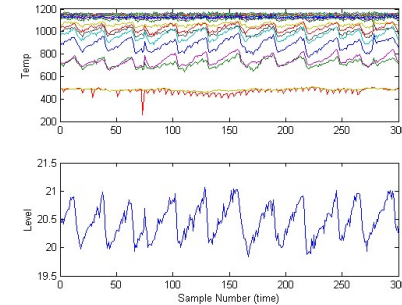


29

EIGENVECTOR
RESEARCH INCORPORATED

MLR Example

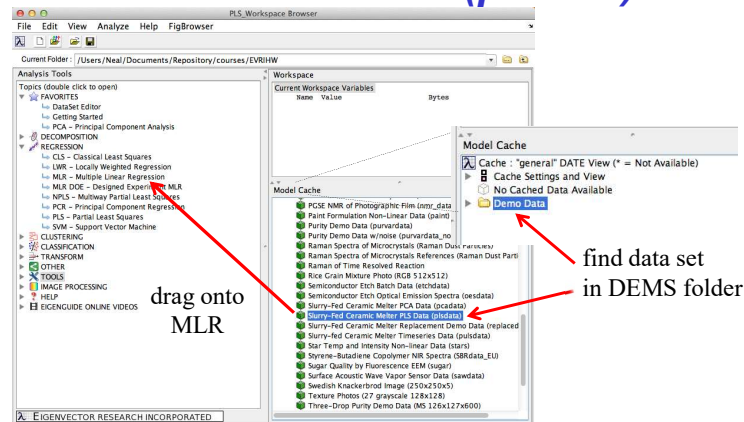
- Use MLR to obtain a relationship between temperature and level in a SFCM



30

EIGENVECTOR
RESEARCH INCORPORATED

Load SFCM Data (plsdata)



31

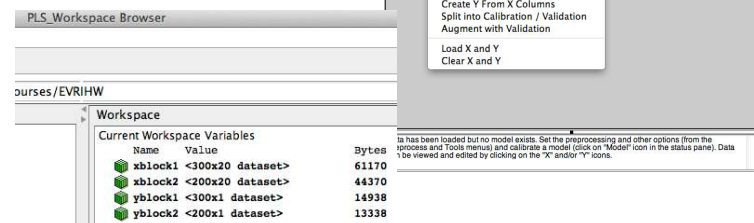
EIGENVECTOR
RESEARCH INCORPORATED

Edit Data

xblock1 and yblock1
will load as calibration

xblock2 and yblock2
will load as validation set

Edit calibration X-block data



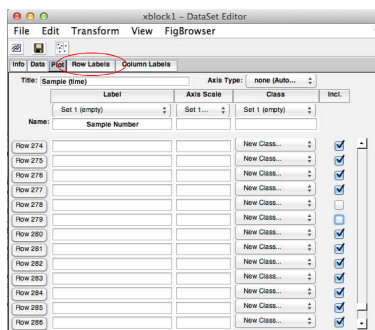
32

EIGENVECTOR
RESEARCH INCORPORATED

Remove Outliers

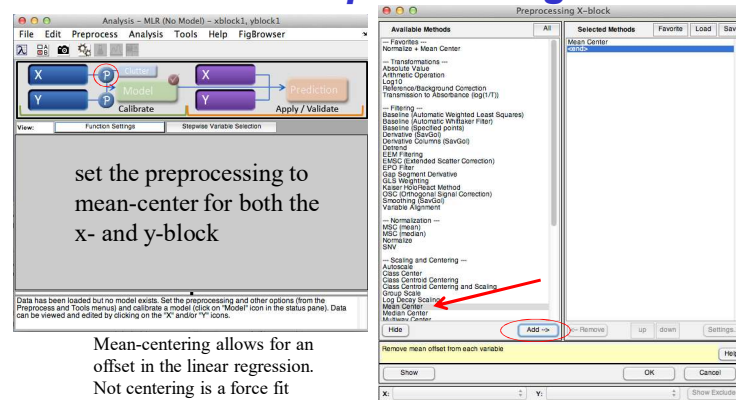
Select “Row Labels”
tab in DataSet Editor

Exclude samples 73,
167, 278 and 279 from
xblock1



33

Set Preprocessing

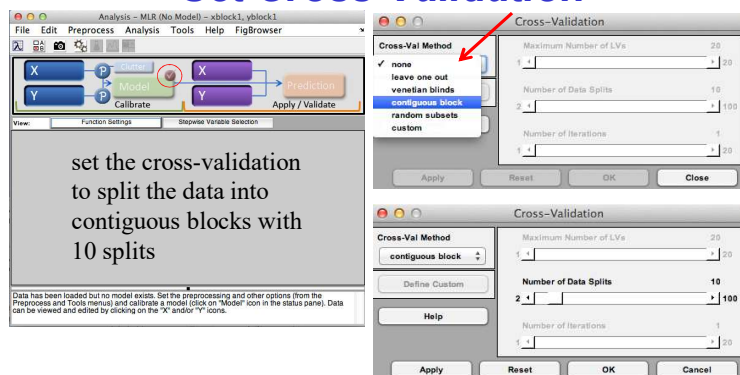


set the preprocessing to
mean-center for both the
x- and y-block

Mean-centering allows for an
offset in the linear regression.
Not centering is a force fit
through zero.

34

Set Cross-Validation



set the cross-validation
to split the data into
contiguous blocks with
10 splits

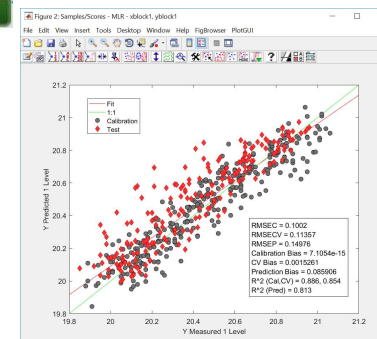
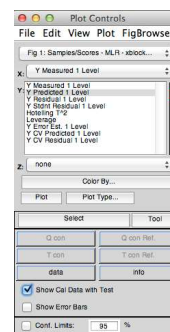
35

MLR Fit to Calibration and Prediction on Validation Data

Click calculate model



Click scores “flask”



36

Ridge Regression

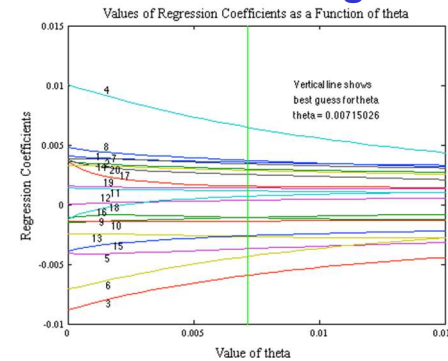
- Ridge Regression (RR) is one way to deal with ill-conditioned problems
- RR gets its name because a constant is added to the “ridge” of the covariance matrix in the formation of the pseudo-inverse:

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X} + \mathbf{I}\theta)^{-1} \mathbf{X}^T$$

- The addition of the ridge ($\mathbf{I}\theta$ term) stabilizes the inverse and shrinks the values of the coefficients
 - this “ridging” is known as matrix regularization

37

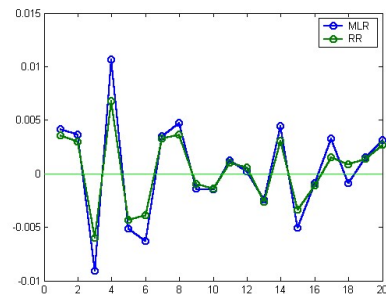
RR Shrinkage



38

`[brr, theta] = ridge(mx, my, 0.015, 31);`

RR and MLR Regression Vectors



39

Problem with MLR and RR

- RR helps stabilize the inverse
 - ridging biases the regression
 - how to determine the ridge parameter θ ?
- MLR does not work when $M < N_x$
- Possible solution: eliminate variables
 - how to choose which variables to keep?
 - stepwise regression or other variable selection
 - lose multivariate advantage - signal averaging
- Another solution: use PCA to reduce original variables to some smaller number of factors
 - retains multivariate advantage
 - noise reduction aspects of PCA

40

Principal Components Regression

- Principal Components Regression (PCR) is one way to deal with ill-conditioned problems
- Property of interest \mathbf{y} is regressed on PCA scores:

$$\mathbf{X}^+ = \mathbf{P}_K \left(\mathbf{T}_K^T \mathbf{T}_K \right)^{-1} \mathbf{T}_K^T$$

- Problem is to determine K the number of factors to retain in the formation of the model

41



Principal Components Regression

$$\mathbf{T}_K \mathbf{b}_{pc} = \mathbf{y} + \mathbf{e} = \mathbf{X} \mathbf{P}_K \mathbf{b}_{pc} \iff \mathbf{b} = \mathbf{P}_K \mathbf{b}_{pc}$$

$$\mathbf{T}_K \mathbf{b}_{pc} = \mathbf{y}$$

$$\mathbf{T}_K^T \mathbf{T}_K \mathbf{b}_{pc} = \mathbf{T}_K^T \mathbf{y}$$

$$\mathbf{b}_{pc} = \left(\mathbf{T}_K^T \mathbf{T}_K \right)^{-1} \mathbf{T}_K^T \mathbf{y}$$

$$\mathbf{b} = \mathbf{P}_K \left(\mathbf{T}_K^T \mathbf{T}_K \right)^{-1} \mathbf{T}_K^T \mathbf{y}$$

$$\mathbf{X}^+ = \mathbf{P}_K \left(\mathbf{T}_K^T \mathbf{T}_K \right)^{-1} \mathbf{T}_K^T$$

- Note that $\mathbf{T}_K^T \mathbf{T}_K$ is $K \times K$ and square

42



Cross-Validation

- Divide data set into J sample subsets
- For **each subset** ($j=1, \dots, J$):
 - Build PCA model using samples in the **remaining** subsets
 - Apply the model to subset j samples
 - Calculate PRESS (Predictive Residual Sum of Squares) for the subset samples:

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

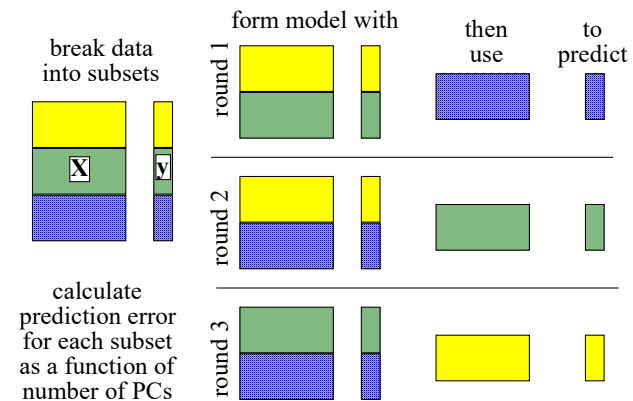
- Look for minimum or “knee” in CumPRESS curve

$$\text{RMSECV} = \left(\frac{1}{M} \sum_{j=1}^J \mathbf{e}_j^2 \right)^{1/2}$$

43



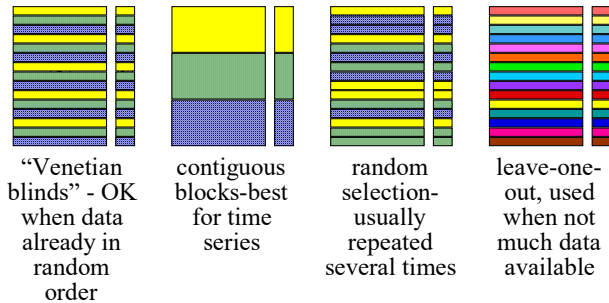
Cross-Validation Graphically



44



Formation of Test Sets



What else?

Custom selection, based on prior knowledge!



45

Cross-Validation Considerations

- Cross-validation method selection criteria
 - *Number* of objects in dataset, M
 - *Order* of objects in dataset
 - *Objective* of cross-validation (specific type of error?)
 - Presence/absence of *replicates*
 - Remember the objective is to mimic future performance
- “Traps” to avoid
 - “Replicate sample trap”
 - Different replicates in both model and test set
 - “External subset selection trap” - extrapolation
 - Test set “space” outside of model set “space”



46

Cross-validation Usage Matrix 1/2

DATASET TYPES	CROSS-VALIDATION METHODS				
	Venetian Blinds	Contiguous Blocks	Random Subsets	Leave-One Out	Custom
GENERAL CV Method Properties	<ul style="list-style-type: none"> • Easy • Relatively quick 	<ul style="list-style-type: none"> • Easy • Relatively quick 	<ul style="list-style-type: none"> • Easy • Can be slow, if m or number of iterations large • Selection of subsets unknown 	<ul style="list-style-type: none"> • Easiest! (Only one parameter) • Avoid using if $m > 20$ 	<ul style="list-style-type: none"> • Flexible • Requires time to define splits
Small data sets (<20 objects)	•	•	<ul style="list-style-type: none"> • OK, if many iterations done 	<ul style="list-style-type: none"> • Good choice.... • ...unless DOE data • OK, but.... 	<ul style="list-style-type: none"> • often needed to avoid the external subset selection trap
randomly-distributed objects	<ul style="list-style-type: none"> • Good choice 	<ul style="list-style-type: none"> • Good choice 	<ul style="list-style-type: none"> • Good choice • Can take a while with large m, many iterations 	<ul style="list-style-type: none"> • Can take a while with large m 	•



47

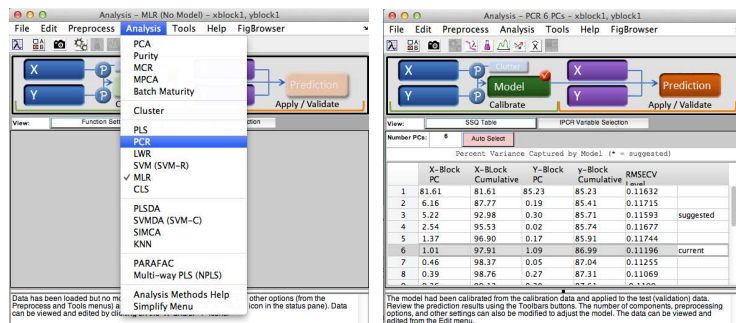
Cross-Validation Usage Matrix 2/2

DATASET TYPES	CROSS-VALIDATION METHODS				
	Venetian Blinds	Contiguous Blocks	Random Subsets	Leave-One Out	Custom
time-series data	<ul style="list-style-type: none"> • Useful for assessing NON-temporal model errors • Can be optimistic with low number of data splits 	<ul style="list-style-type: none"> • Useful for assessing temporal stability of model 	•	•	•
Batch data	<ul style="list-style-type: none"> • Useful for assessing predictability within batches/parts of batches 	<ul style="list-style-type: none"> • Useful for assessing predictability between batches/parts of batches 	•	•	<ul style="list-style-type: none"> • Can manually select “batch-wise” test sets
Blocked data (replicates)	<ul style="list-style-type: none"> • Beware the replicate sample trap (optimistic results) 	<ul style="list-style-type: none"> • Good way to avoid replicate sample trap • Beware the external subset selection trap! 	<ul style="list-style-type: none"> • Can use to avoid replicate sample trap (high number of splits, high iterations preferable) 	<ul style="list-style-type: none"> • overly optimistic results, due to replicate sample trap 	•
Designed Experiment (DOE) data	<ul style="list-style-type: none"> • Dangerous, unless object order is randomized 	<ul style="list-style-type: none"> • Dangerous, unless object order is randomized 	•	<ul style="list-style-type: none"> • Not recommended (external subset selection trap) 	<ul style="list-style-type: none"> • often needed to avoid the external subset selection trap



48

Switch Analysis from MLR to PCR, Calculate Model

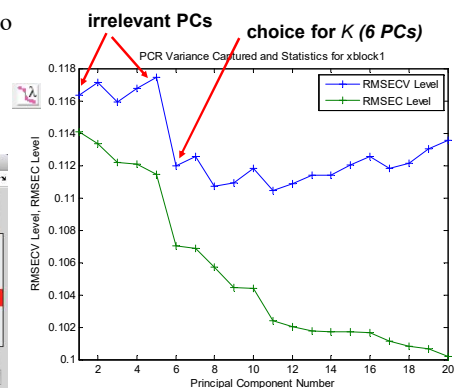


49

PCR Cross-Validation Example

Recall: time series data so contiguous block CV

Click variance “lambda” icon to get CV results



50

PCR Variance Captured

Percent Variance Captured by PCR Model

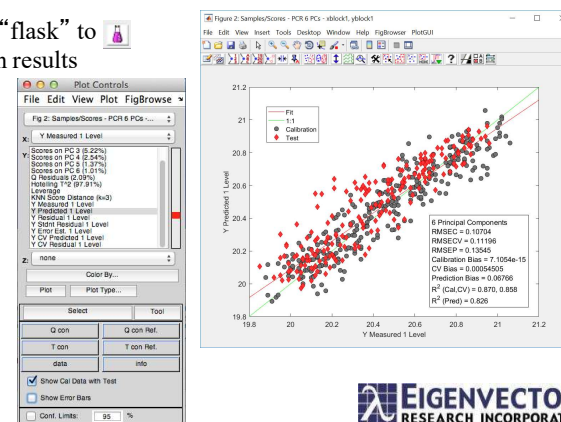
PC #	-----X-Block-----		-----Y-Block-----	
	This PC	Total	This PC	Total
1	81.61	81.61	85.23	85.23
2	6.16	87.77	0.19	85.41
3	5.22	92.98	0.30	85.71
4	2.54	95.53	0.02	85.74
5	1.37	96.90	0.17	85.91
6	1.01	97.91	1.09	86.99
7	0.46	98.37	0.05	87.04
8	0.39	98.76	0.27	87.31
9	0.36	99.12	0.30	87.61
10	0.24	99.37	0.02	87.63

51

PCR Model Fit to Calibration Data and Validation Predictions

Click scores “flask” to get prediction results

Adjust Plot Controls to get desired plot



52

Problems with PCR

- Some PCs not relevant for prediction, but are only relevant for describing variance in **X**
 - leads to local minima and increase in PRESS
- This is a result of PCs determined without using information about property to be predicted **y**
- A solution is to find factors using information from **y** and **X**

53



Partial Least Squares

- PLS is related to PCR and MLR
 - PCR captures maximum variance in **X**
 - MLR achieves maximum correlation between **X** and **Y**
 - PLS tries to do both by maximizing covariance between **X** and **Y**
- Requires addition of weights **W** to maintain orthogonal scores
- Factors calculated sequentially by projecting **Y** through **X**

$$\mathbf{X}^+ = \mathbf{R}_K (\mathbf{T}_K^T \mathbf{T}_K)^{-1} \mathbf{T}_K^T = \mathbf{W}_K (\mathbf{P}_K^T \mathbf{W}_K)^{-1} (\mathbf{T}_K^T \mathbf{T}_K)^{-1} \mathbf{T}_K^T$$

54

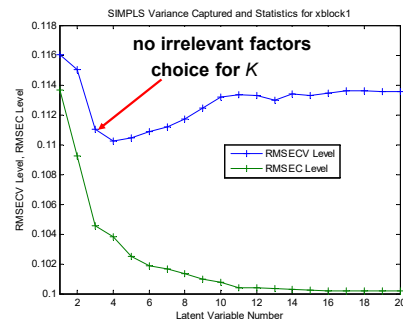


PLS Cross-Validation Example

Set Analysis to PLS

Calculate model

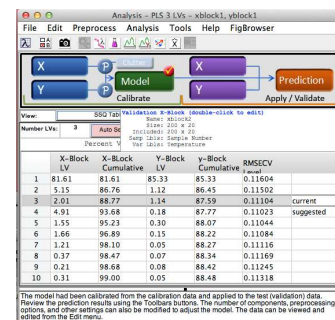
Click variance “lambda”
to get CV results



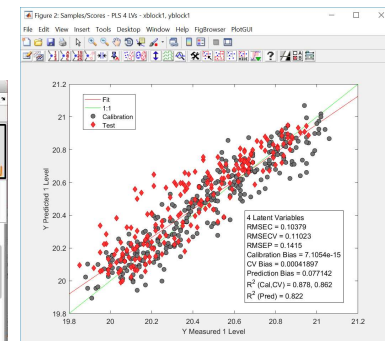
55



Set number of LVs to 4



56



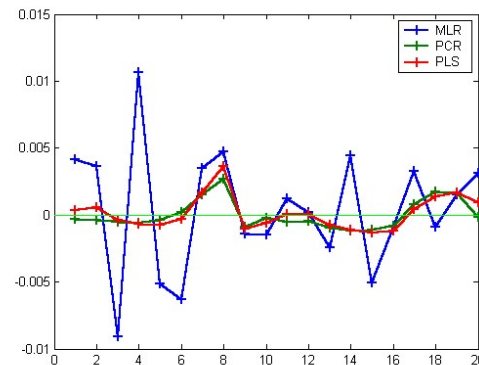
PLS Variance Captured

Percent Variance Captured by PLS Model				
LV #	---X-Block---		---Y-Block---	
	This LV	Total	This LV	Total
1	81.61	81.61	85.33	85.33
2	5.15	86.76	1.12	86.45
3	2.01	88.77	1.14	87.59
4	4.91	93.68	0.18	87.77
5	1.55	95.23	0.30	88.07
6	1.66	96.89	0.15	88.22
7	1.21	98.10	0.05	88.27
8	0.37	98.47	0.07	88.34
9	0.21	98.68	0.08	88.42
10	0.31	99.00	0.05	88.48

57



Regression Vectors



58



NIPALS: PLS Algorithm

Choose $\mathbf{u}_1 = \mathbf{y}$ or one column of \mathbf{Y}

$$\mathbf{w}_1 = \frac{\mathbf{X}^T \mathbf{u}_1}{\|\mathbf{X}^T \mathbf{u}_1\|} \quad (1)$$

$$\mathbf{t}_1 = \mathbf{X} \mathbf{w}_1 \quad (2)$$

$$\mathbf{q}_1 = \frac{\mathbf{u}_1^T \mathbf{t}_1}{\|\mathbf{u}_1^T \mathbf{t}_1\|} \quad (3)$$

$$\mathbf{u}_1 = \mathbf{Y} \mathbf{q}_1 \quad (4)$$

Check for convergence by comparing \mathbf{t}_1 to previous \mathbf{t}_1 . If $\mathbf{Y} = \mathbf{y}$ skip (3) and (4) and continue

$$\mathbf{p}_1 = \frac{\mathbf{X}^T \mathbf{t}_1}{\|\mathbf{t}_1^T \mathbf{t}_1\|} \quad (5)$$

$$\mathbf{p}_{1\text{new}} = \frac{\mathbf{p}_{1\text{old}}}{\|\mathbf{p}_{1\text{old}}\|} \quad (6)$$

$$\mathbf{t}_{1\text{new}} = \mathbf{t}_{1\text{old}} \|\mathbf{p}_{1\text{old}}\| \quad (7)$$

$$\mathbf{w}_{1\text{new}} = \mathbf{w}_{1\text{old}} \|\mathbf{p}_{1\text{old}}\| \quad (8)$$

Find the regression coefficient for the inner relation:

$$b_1 = \frac{\mathbf{u}_1^T \mathbf{t}_1}{\mathbf{t}_1^T \mathbf{t}_1} \quad (9)$$

After calculating scores and loadings for first Latent Variable, the X and Y-block residuals are calculated:

$$\mathbf{E}_1 = \mathbf{X} - \mathbf{t}_1 \mathbf{p}_1^T \quad (10)$$

$$\mathbf{F}_1 = \mathbf{Y} - \mathbf{u}_1 \mathbf{q}_1^T \quad (11)$$

Repeat entire procedure replacing \mathbf{X} and \mathbf{Y} with their residuals

Non-linear iterative partial least squares (NIPALS).
Geladi, Paul; Kowalski, Bruce (1986), "Partial Least Squares Regression: A Tutorial",
Analytica Chimica Acta **185**: 1–17, doi:10.1016/0003-2670(86)80028-9

59



Other PLS Algorithms

- It can be shown that \mathbf{w}_1 is given by
$$\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}_1 = \lambda \mathbf{w}_1$$
- The SIMPLS algorithm uses an orthogonalization of a Krylov sequence (faster than NIPALS algorithm)
- The important thing to remember is:

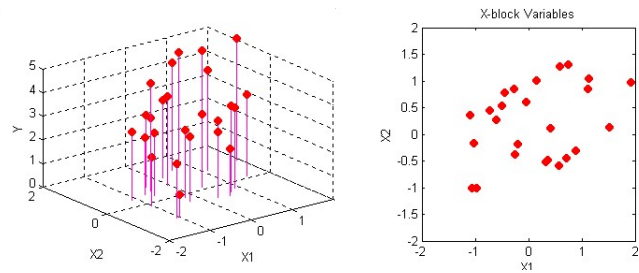
PLS finds factors in \mathbf{X} which are correlated with \mathbf{Y} while describing large amounts of variance in \mathbf{X}

Sijmen de Jong, "SIMPLS: an alternative approach squares regression,"
Chemometrics and Intelligent Laboratory Systems, **18** (1993) 251-263

60

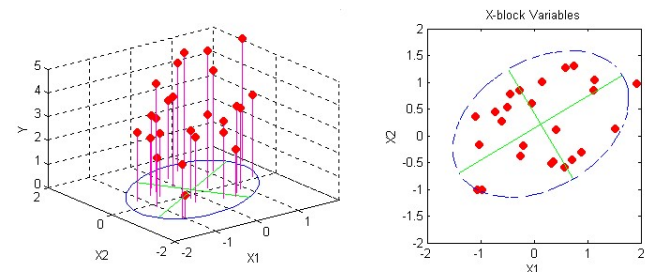


Y Projected onto X Plane



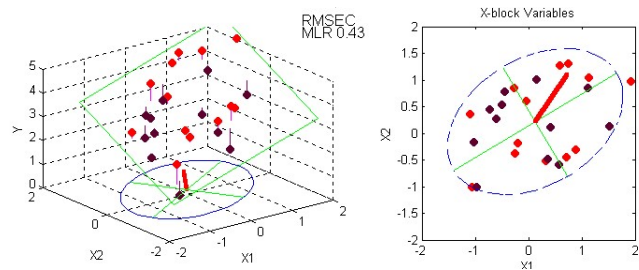
61

PCA of X-Block



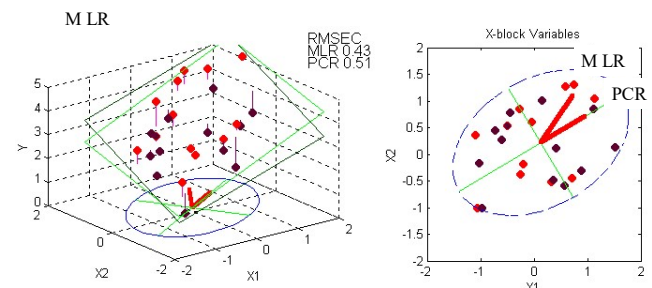
62

MLR Regression Vector and Surface



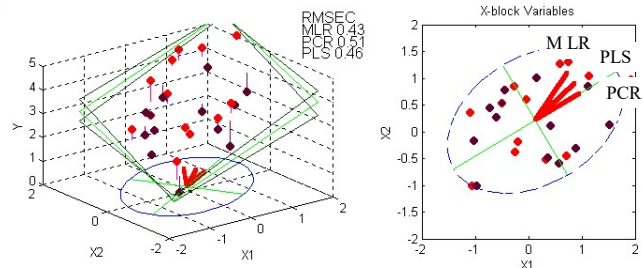
63

PCR Regression Vector and Surface



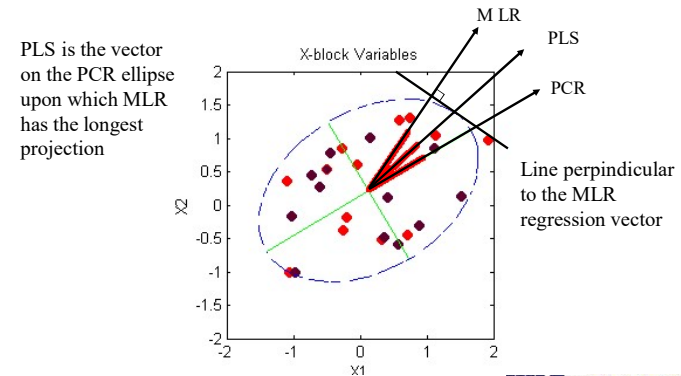
64

PLS Regression Vector and Surface



65

Geometric Relationship of MLR, PCR, and PLS



66

PLS for Multivariate Y

- PLS can be used to relate multivariate **X** to multivariate **Y** (*a.k.a.*, PLS2)
 - outer relationships

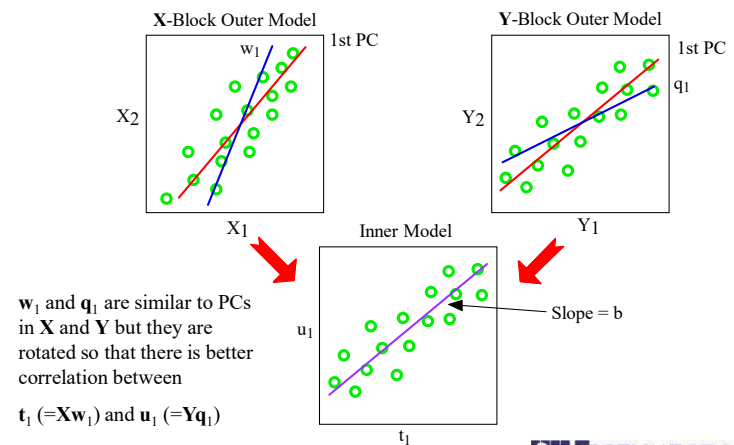
$$\mathbf{X} = \mathbf{T}_K \mathbf{P}_K^T + \mathbf{E}$$

$$\mathbf{Y} = \mathbf{U}_K \mathbf{Q}_K^T + \mathbf{F}$$
 - inner relationship

$$\mathbf{U}_K = \mathbf{T}_K \mathbf{B}_K$$
 - i.e.*, the scores in **Y** are linear combinations of the scores in **X**

67

PLS2



68

Model Quality Measures

- Root Mean Square Error (RMSE) Metrics

- RMSEC
- RMSECV
- RMSEP

- In units of the Y variable!

- Correlation Coefficient (r)

- Unit-less
- Considers the range of Y

$$\sqrt{\frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{m}}$$

$$\frac{\sum_{i=1}^m (\hat{y}_i - \bar{\hat{y}})(y_i - \bar{y})}{\sqrt{\left(\sum_{i=1}^m (\hat{y}_i - \bar{\hat{y}})^2\right) \left(\sum_{i=1}^m (y_i - \bar{y})^2\right)}}$$



69

Root Mean Square Error (RMSE) Metrics

- These are used to assess a model's *fit to the data* and *predictive ability on new data*
- Measures “average” deviation of model estimates from the measured data
- Measure of *fit* - root mean squared error of *calibration* (RMSEC)

$$\text{RMSEC} = \sqrt{\frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{m}}$$

i's refer to all samples used to build the model



70

Cross-Validation Error

- RMSEC measures *fit to the model data*. RMSECV (root mean squared error of cross-validation) is an estimate of *predictive power on new data*.
- RMSECV is a function of the number of factors k and how the test sets were selected

$$\text{RMSECV} = \sqrt{\frac{\sum_{j=1}^J \sum_{i=1}^{m_j} (y_i - \hat{y}_i)^2}{m_j}} = \sqrt{\frac{\text{PRESS}}{m_j}}$$

J's refer to different CV subsets

i's refer to CV subset samples- not used to build CV models



71

Prediction Error

- Prediction error is often used to *validate* a model and is a true measure of the *predictive power on new data*
- Measure of *prediction error* - root mean squared error prediction (RMSEP)

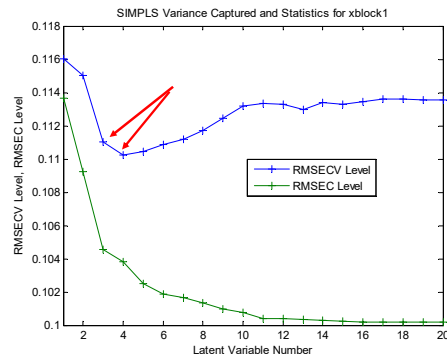
$$\text{RMSEP} = \sqrt{\frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{m}}$$

i's refer to samples NOT used to build the model



72

RMSE metrics, as a function of factor (PC, LV)



RMSEC and RMSECV can also be used to determine the optimal number of factors (LVs, PCs) to be used in a model

Comparison of Models

- MLR, PCR, and PLS models were constructed using SFCM data: Calibration used (xblock1) and test used (xblock2).

	MLR	PCR	PLS
RMSEC	0.1002	0.1070	0.1038
RMSECV	0.1136	0.1120	0.1102
RMSEP	0.1498	0.1355	0.1415

- Fit** and **prediction** are two entirely different aspects of a model's performance

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
 - if data is time series use block CV due to correlated noise
 - be conservative, models are more often over-fit than under-fit
 - 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?

Model Diagnostics

- Diagnostics useful for finding outliers/uniques
- X**-block Q residual and T^2
- X**-block leverage and studentized **Y**-block residuals
- Try SFCM example **without removing outliers**

Build PLS Model on SFCM Data

- Construct a linear regression for yblock1 from xblock1 (time series data)
 - predict level of slurry fed ceramic melter (Y-block)
 - using melter temperatures (X-block)
- Test the model on xblock2 and yblock2

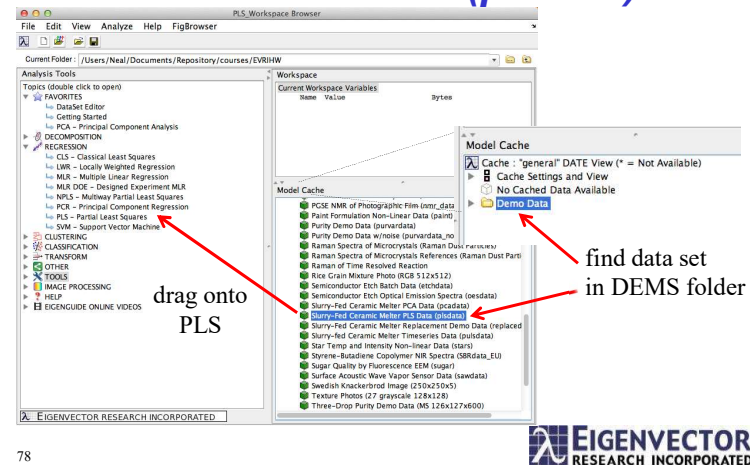
IF data still loaded, can do:

- Edit/Calibration/X-block Data
- Row Labels Tab
- right-click on "Incl" : "Clear/Reset" (checks all rows)
- calculate model



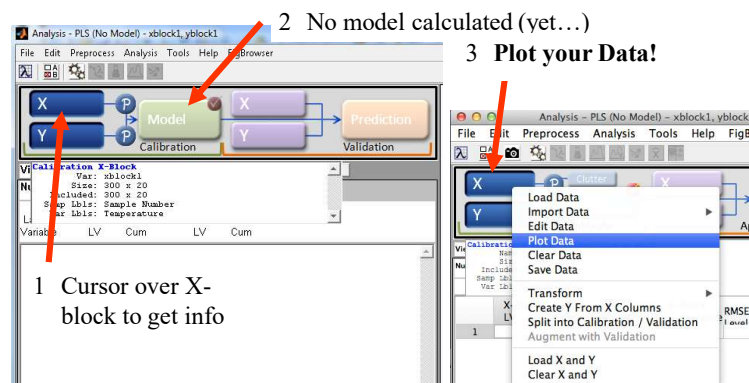
77

Load SFCM Data (plsdata)



78

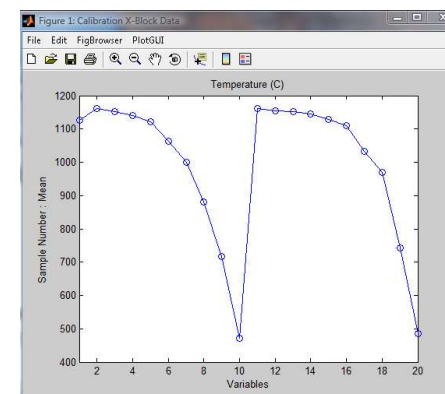
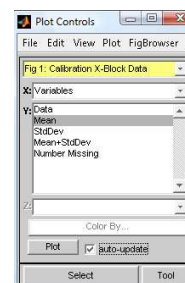
Data: loaded but not analyzed



79

Plot Your Data

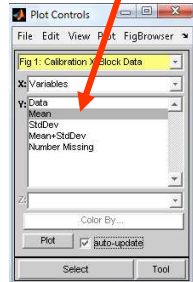
Default plot
Column / Variable
mean



80

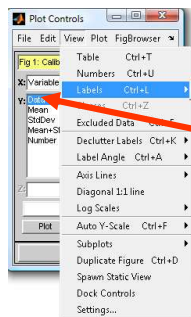
Plot Your Data

- 1 Plot control default
can look at summary stats



81

the Plot control generates plots in
MATLAB figure windows



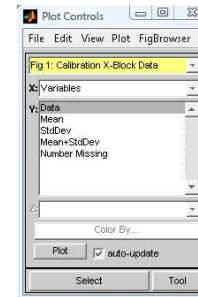
- 2 under view menu
check
labels: Temperature

- 3 under Y: menu
highlight Data

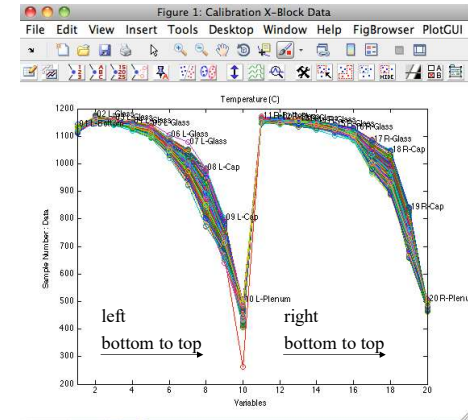


Plot Your Data

all samples vs.
variables



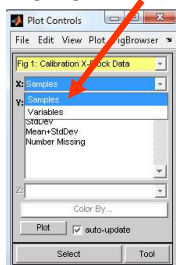
82



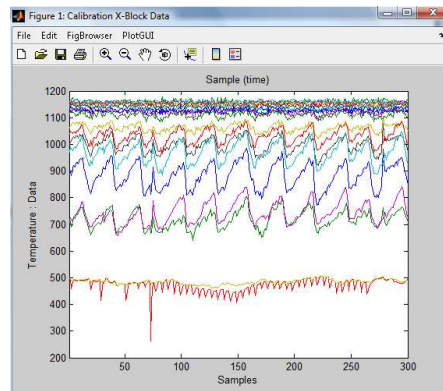
Plot Your Data

all variables vs.
samples (time)

under X: menu
highlight Samples

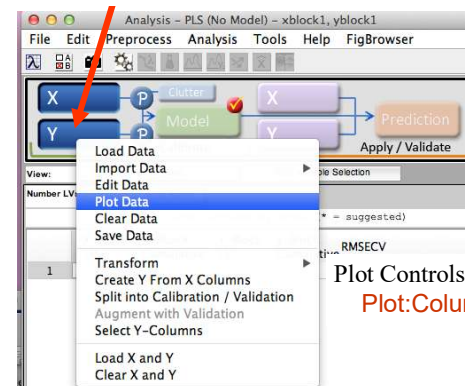


83

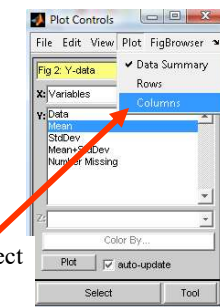


Plot Y data

Click



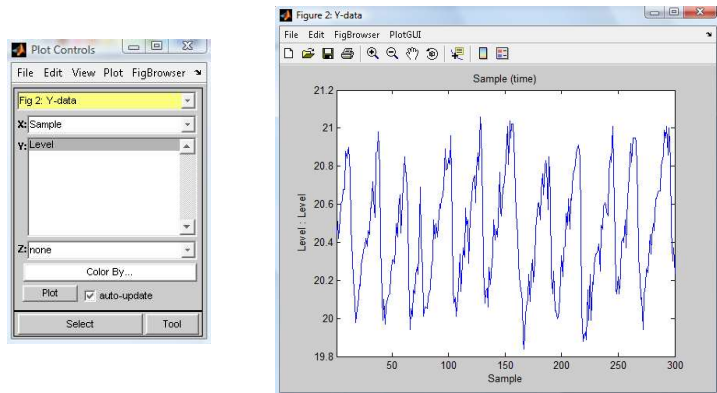
Plot Controls: Select
Plot: Columns



84



Plot Y Data

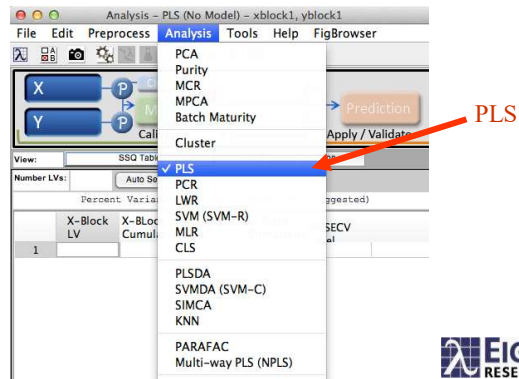


Plot Your Data: Summary

- Bottom temperatures higher than top temperatures
 - surface and plenum space is cooler than the bottom
- Trend in time
 - “saw-tooth” pattern showing correlation between some temperatures and level

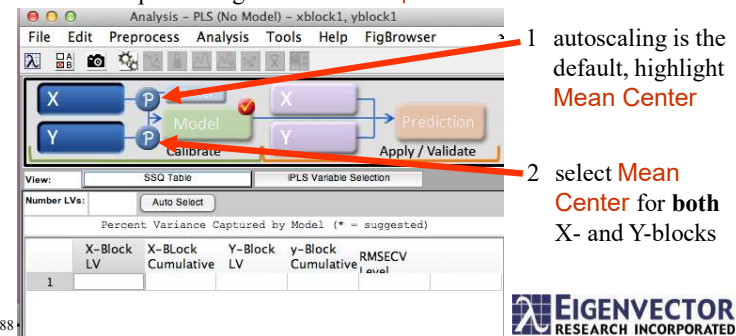
Which Regression?

- BACK to Analysis Window, then Analysis menu
- Choose PLS...



How Should We Scale the Data?

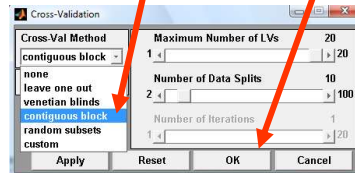
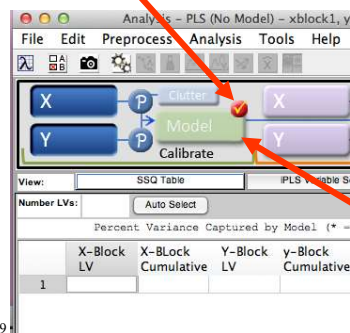
- Variables are in same units and there's a reason to believe that variance is associated with signal: Suggests mean-centering.
- X Preprocessing is set under **Preprocess:X-block**



How to Cross-Validate?

- Time series data suggests contiguous block cross-validation

- click **check** for Cross-Validation
- contiguous block, then OK



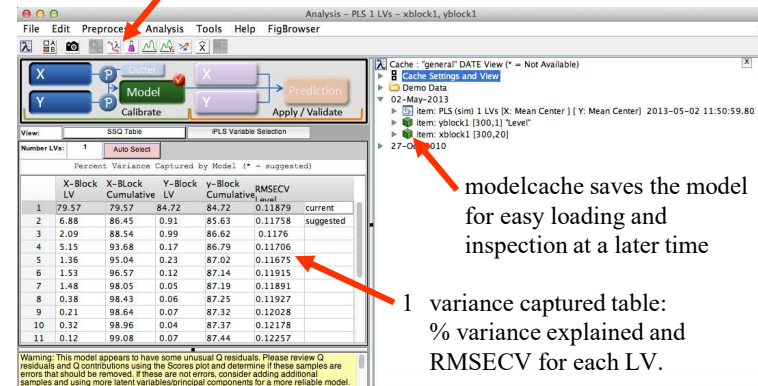
- click **Model** button to perform cross-validation and build the regression model



89

Regression Results

- Click **Eigenvalue** button to plot RMSEC and RMSECV

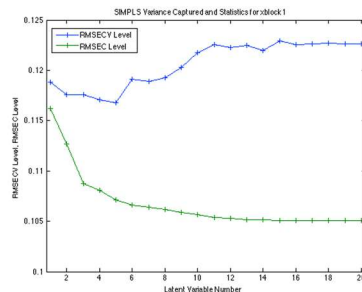
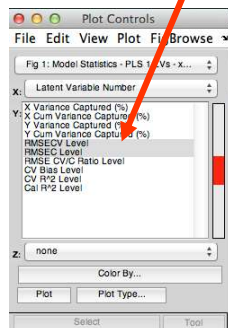


modelcache saves the model for easy loading and inspection at a later time

- variance captured table: % variance explained and RMSECV for each LV.

RMSECV Plot

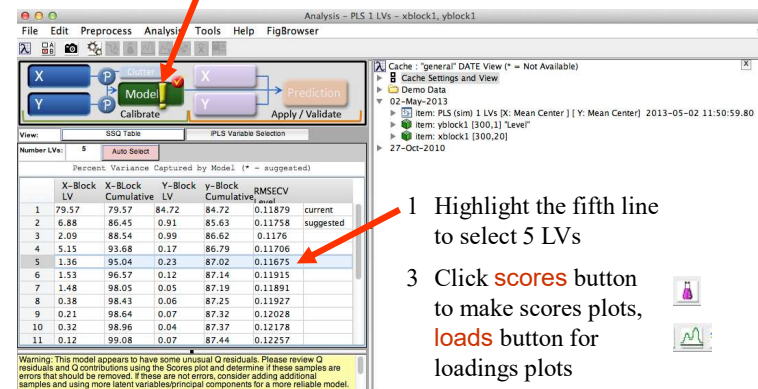
Plot RMSEC and RMSECV vs. LV.



91

Choose Number of LVs

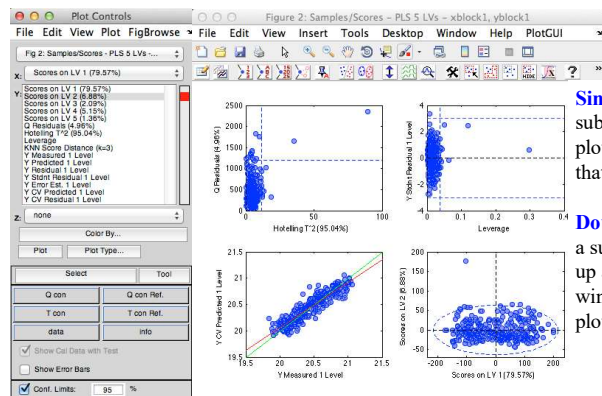
- Click **Model** button to reconstruct the model with 5 LVs



- Highlight the fifth line to select 5 LVs

- Click **scores** button to make scores plots, **loads** button for loadings plots

Scores Summary Plot

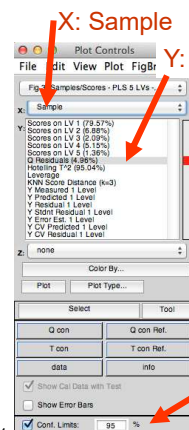


Single-click on a subplot: brings up plot controls for that plot

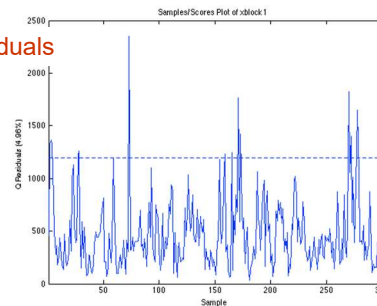
Double-click on a subplot: brings up a new figure window with that plot only

X-Block Q Residuals (by sample)

1 set plot controls

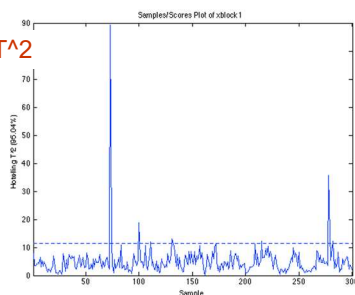
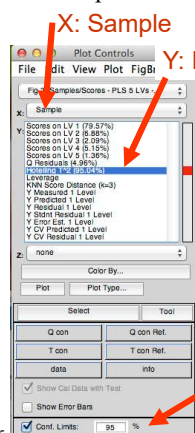


2 Check Conf. Limits



X-Block Hotelling T^2 (by sample)

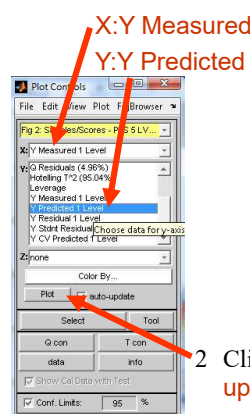
1 set plot controls



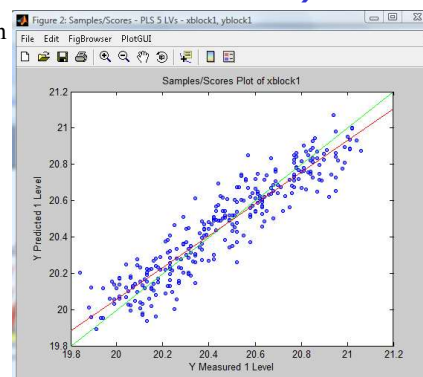
2 Check Conf. Limits

Calibration Curve (Predicted vs. Measured)

1 Click Scores Button



2 Click Plot if auto-update is not checked

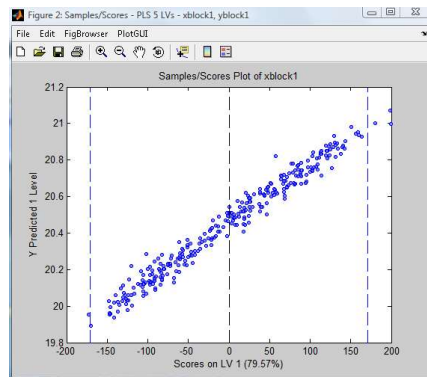
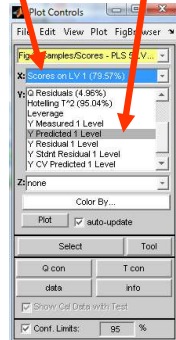


Predicted Y vs. LV 1

1 Click **Scores** Button

X:LV 1

Y:Y Predicted



EIGENVECTOR
RESEARCH INCORPORATED

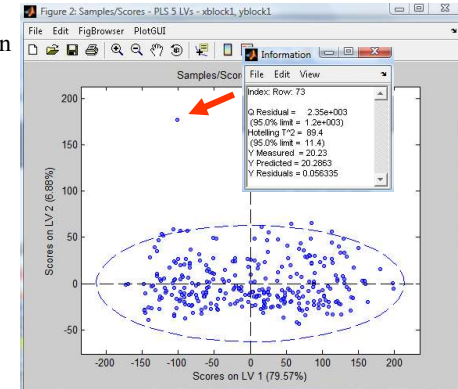
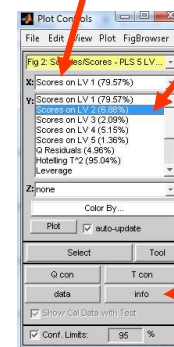
97

LV 2 vs. LV 1

1 Click **Scores** Button

X:Scores on LV 1

Y:Scores on LV 2



2 Click **info** Button
and select sample

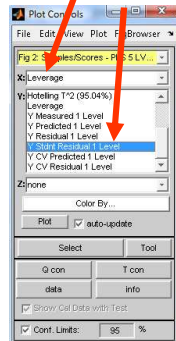
EIGENVECTOR
RESEARCH INCORPORATED

Studentized Y Residual vs. Leverage

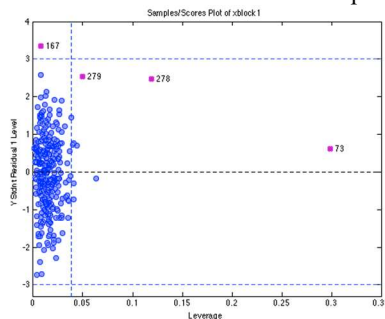
1 Click **Scores** Button

X:Leverage

Y:Y Stdnt Residual



2 select
outliers
3 view
numbers



EIGENVECTOR
RESEARCH INCORPORATED

Calculation of Studentized Residuals

- Given the pseudo-inverse \mathbf{X}^+ the leverage for a sample $\mathbf{x}_m = \mathbf{x}(m,:)$ and column $\mathbf{X}^+(\cdot, m)$ is

$$l_m = \mathbf{x}_m \mathbf{X}^+(\cdot, m)$$

- Studentized residuals for column m^{th} of \mathbf{y} , $t_{e,m}$

$$e_m = \hat{y}_m - y_m$$

$$\sigma = \left(\frac{1}{M-K} \mathbf{e}^T \mathbf{e} \right)^{1/2}$$

$$t_{e,m} = \frac{e_m}{\sigma(1-l_m)^{1/2}}$$

A studentized residual is the quotient resulting from the division of a residual by an estimate of its standard deviation and is a form of t-statistic.

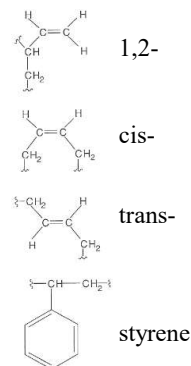
EIGENVECTOR
RESEARCH INCORPORATED

100

How Much Leverage is Too Much?

- In PLS and PCR a good rule of thumb is $3K/M$, where K is the number of LVs or PCs, and M is the number of samples
- In MLR, use $2N_x/M$, where N_x is the number of X-block variables

Regression Example

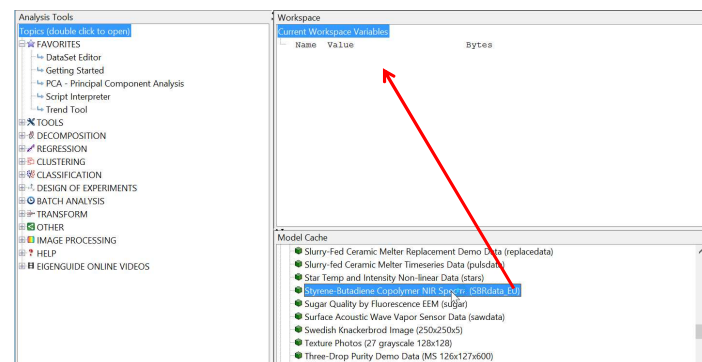


- NIR transmission spectra of styrene-butadiene copolymers
- Different amounts of 4 analytes
 - All 4 are known for all samples (by NMR)
- Data file: [SBRdata_EU.mat](#)
 - 60 calibration samples in arrays **Xcal**, **Ycal**
 - 10 test samples in arrays **Xtest**, **Ytest**

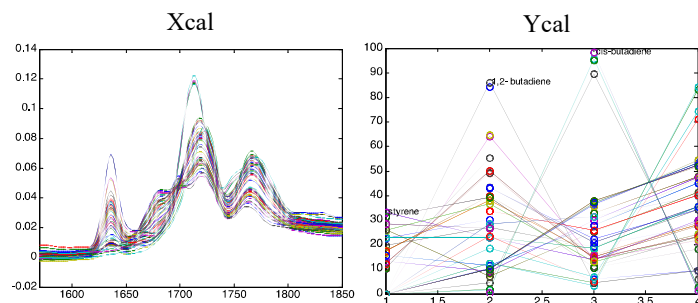
Regression Methods

- Compare the styrene predictions using the four different regression approaches below.
- CLS – no centering
- MLR (stepwise) – no centering
- PCR – no centering
- PLS – no centering
- Additionally: show results with mean centering

Load Into Workspace

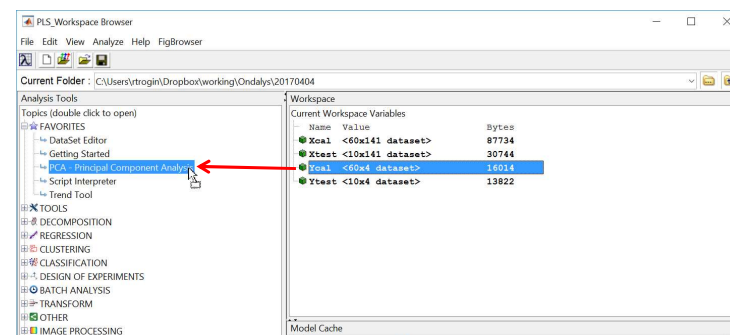


Load and Plot Calibration Data



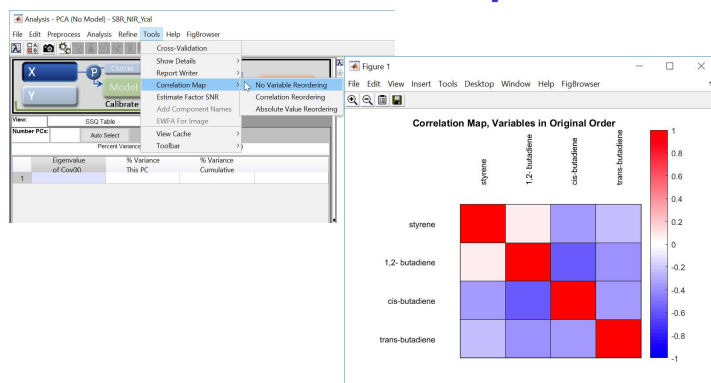
105

Try This:



106

Correlation Map



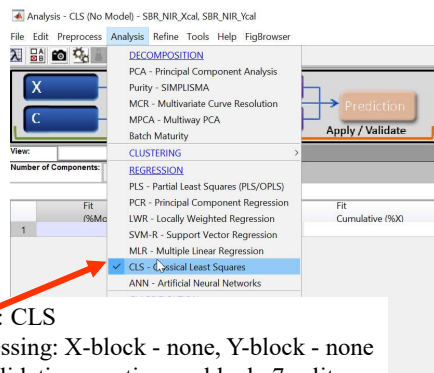
107

Loading Data

- Calibration data
 - Xcal as x-block
 - Ycal as y-block
- Validation data
 - Xtest as x-block
 - Ytest as y-block

108

CLS Regression

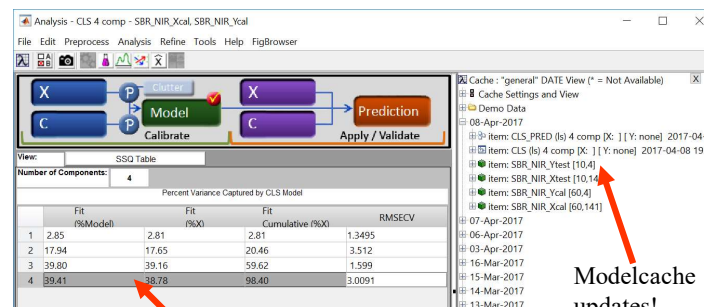


- 1 Analysis: CLS
- 2 Preprocessing: X-block - none, Y-block - none
- 3 Cross-validation: contiguous block, 7 splits
- 4 Calculate model



109

CLS results



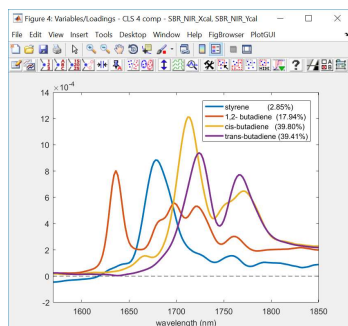
4 components determined "automatically"
(because of 4 Y variables)



110

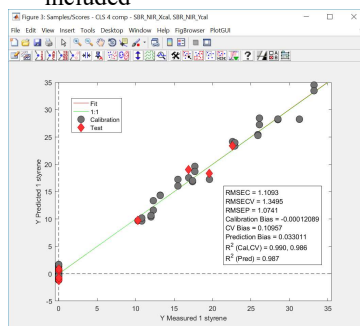
CLS Results

Loadings: all 4 components



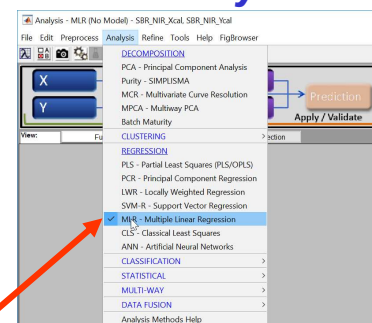
This is the S matrix!

Scores: Styrene predicted vs.
Styrene measured, Cal data
included



111

MLR- styrene



- 1 Analysis: MLR
- 2 Preprocessing: X-block - none, C-block - none
- 3 Cross-validation: contiguous block, 7 splits
- 4 Right-click calibration Y-block => Select Y Columns => Styrene
- 5 Stepwise Variable Selection button

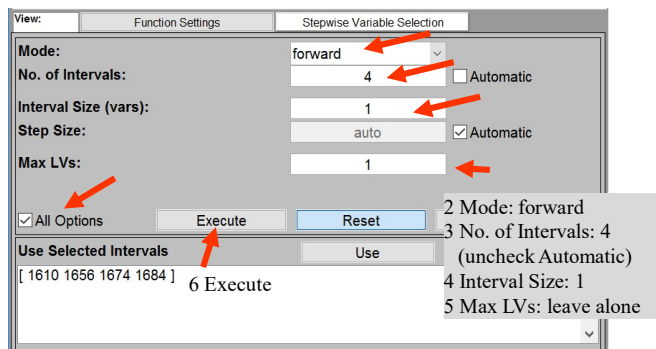


112

Stepwise Variable Selection

1 Check "All Options"

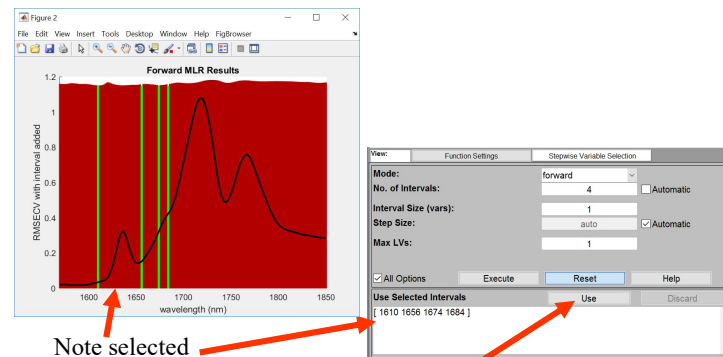
6 Execute



EIGENVECTOR
RESEARCH INCORPORATED

113

Forward MLR Results

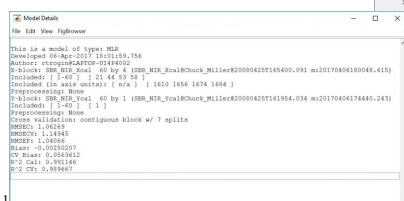
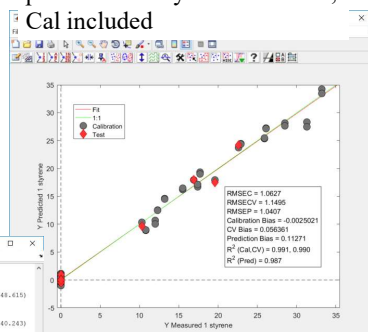
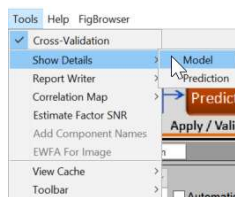


EIGENVECTOR
RESEARCH INCORPORATED

114

MLR results

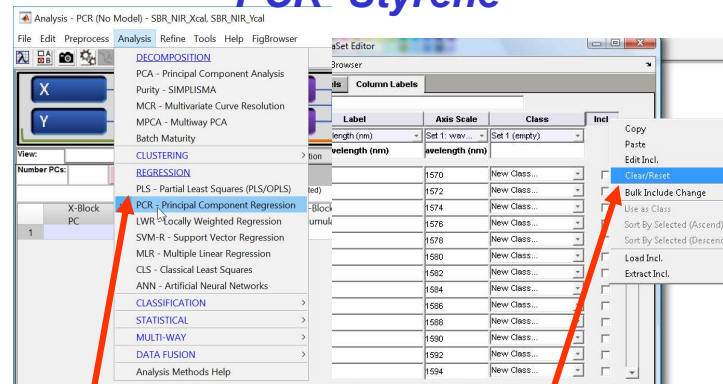
Scores button: Styrene predicted vs Styrene measured, Cal included



EIGENVECTOR
RESEARCH INCORPORATED

1

PCR- Styrene

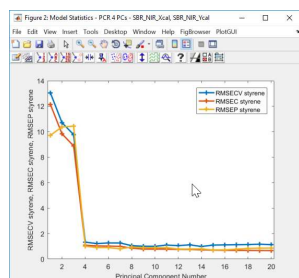


EIGENVECTOR
RESEARCH INCORPORATED

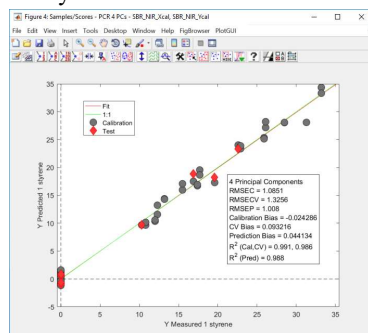
116

PCR Results

Eigenvalue button:
RMSE[C], [CV], and [P]



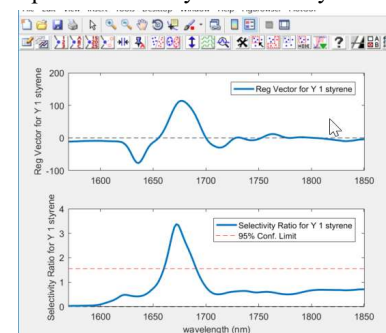
Scores: Styrene predicted vs.
Styrene measured



117

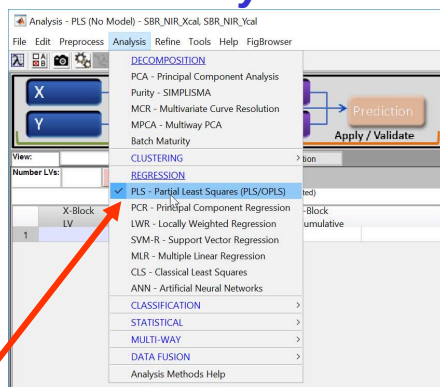
PCR Results

Loadings button: Regression Vector for Styrene
Type "2" (creates second plot)
Second plot: Selectivity Ratio for Styrene



118

PLS Styrene

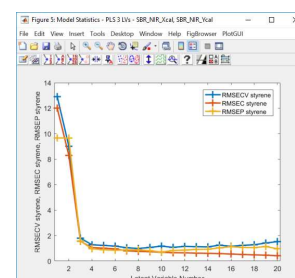


- 1 Analysis: PLS
- 2 Calculate

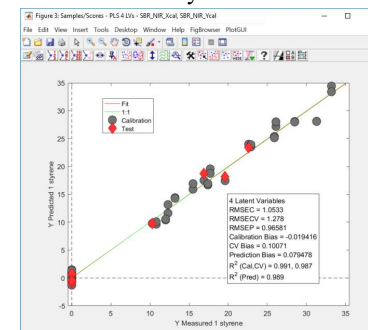
119

PLS Results

Eigenvalue button:
RMSE[C], [CV], and [P]

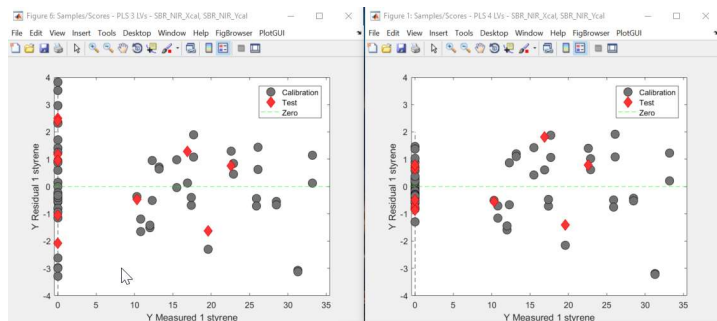


Scores button: Predicted vs.
Measured Styrene



120

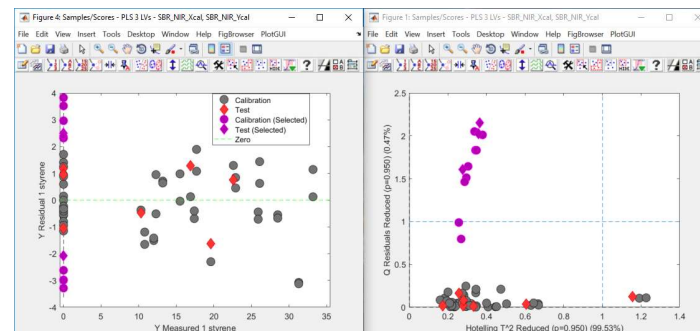
PLS – 3 vs 4 LVs



Percent Variance Captured by Model (T = suggested)					
	X-Block	X-Block	Y-Block	Y-Block	RMSEC
LV	Cumulative	IV	Cumulative	IV	
1	94.07	96.87	34.09	34.09	1.2100
2	2.31	99.12	33.02	68.91	0.0111
3	0.42	99.51	29.01	98.94	1.1917
4	0.19	99.93	0.02	99.51	1.1778
5	0.05	99.98	0.04	99.55	1.2302
6	0.01	99.99	0.05	99.60	1.1743

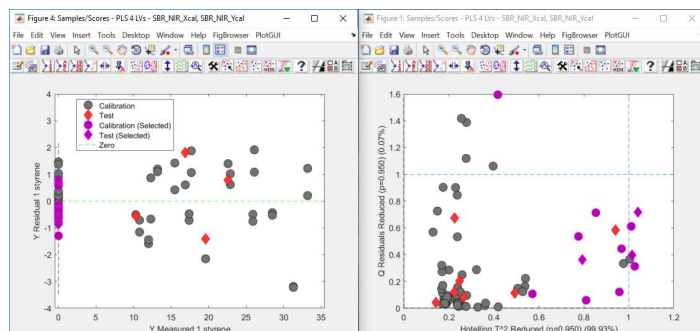
121

Back to 3 LVs



122

And at 4 LVs



123

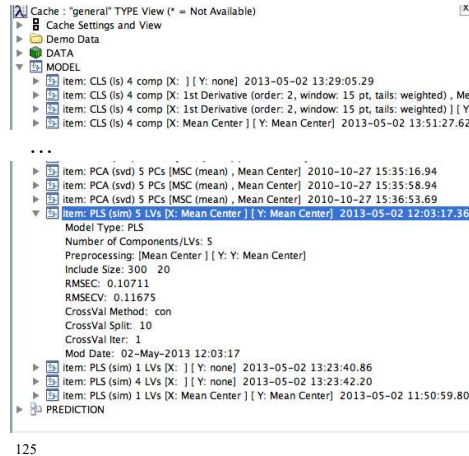
Regression Summary- Styrene

	RMSEC	CV	P	Comments
CLS	1.11	1.35	1.07	4 factors
MLR	1.06	1.15	1.04	4 stepwise-selected variables
PCR	1.08	1.33	1.01	4 factors
PLS	1.05	1.28	0.97	4 factors
CLS	1.39	1.61	1.33	4 factors
MLR	0.89	1.07	0.97	4 stepwise-selected variables
PCR	0.84	0.94	0.73	4 factors
PLS	0.84	0.95	0.73	4 factors

Mean centered Not centered

124

Modelcache

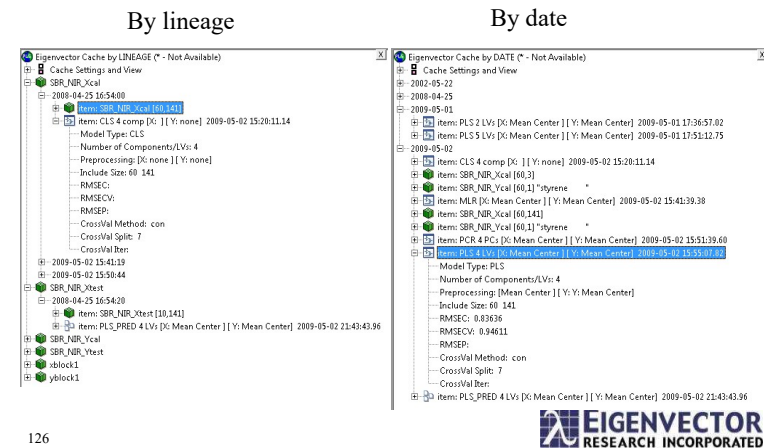


- Can list recent work in 3 ways:
 - By type (shown here)
 - By lineage
 - By date
- Also, can load/show/save any item in the cache

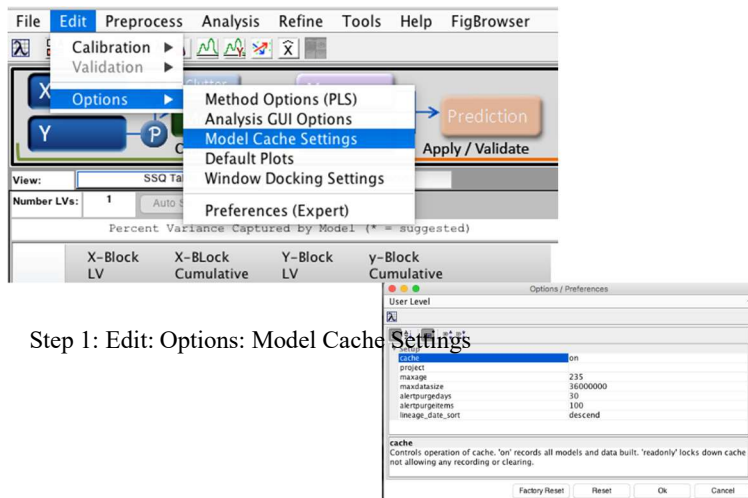


125

Modelcache by lineage and date



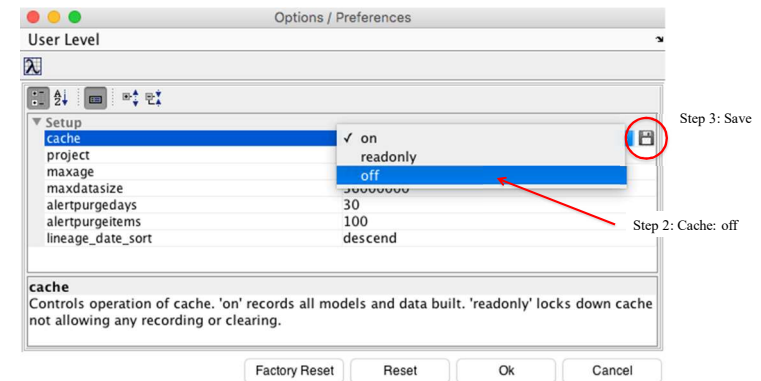
126



Step 1: Edit: Options: Model Cache Settings



127



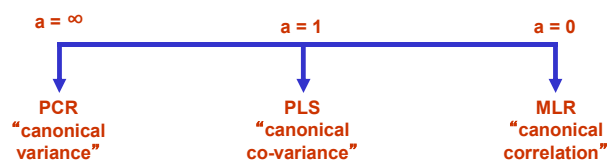
Step 2: Cache: off
Step 3: Save



128

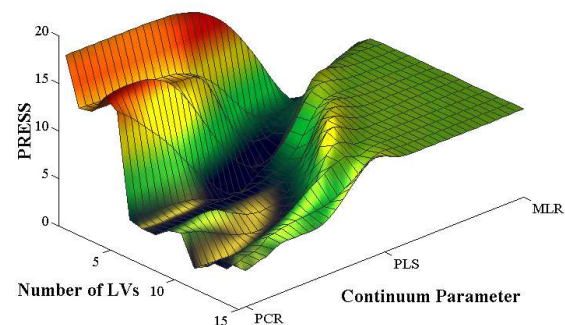
Continuum Regression

- PCR, PLS and MLR can be unified under the single technique Continuum Regression (CR)
- CR is continuously adjustable and encompasses PLS and includes PCR and MLR at the extremes



129

CR Press Surface



130

Missing Data

- MDCHECK
 - Checks data sets for 'NaN' and 'inf' and replaces with values consistent with a PCA model (if desired)
 - e.g., see the ISFINITE function
 - This is an iterative method
 - Example, use some data from SFCM

```
>> x = xblock1.data(1:50,[6:9 16:19]);
>> x2 = x(:,2);
>> x(2:4:50,2) = NaN;
```

← every 4th sample of column 2 removed

131

Call MDCHECK

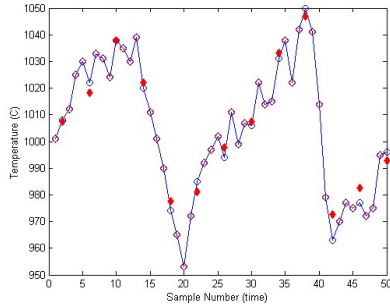
- Change the options to reduce the number of PCs

```
>> options = mdcheck('options')
options =
    max_pcs: 5
   frac_ssq: 0.9500
  meancenter: 'yes'
      output: 'no'
   tolerance: [1.0000e-006 100]
   max_missing: 0.4000
>> options.max_pcs = 3;
>> [flag,mismatch,xfill] = mdcheck(x,options);
```

132

MDCHECK Results

```
>> plot(1:50,x2,'ob-',1:50,xfill(:,2),'rd'), hold on
>> plot(2:4:50,xfill(2:4:50,2),'rd','markerfacecolor',[1 0 0])
```



Model Development

- Developing PCR or PLS models
 - center and scale the data (as appropriate)
 - cross-validate to determine number of factors
 - check **X**-block Q , T^2 , leverage, and **Y**-block residuals for outliers
 - remove / explain outliers
 - check RMSEC and RMSECV values for overfit
 - repeat as necessary
- PCR or PLS models consist of
 - mean and scaling vectors
 - **X**-block loadings **P**, scores **T**, and weights **W** (if PLS)
 - **Y**-block loadings **Q**, and scores **U**
 - inner coefficients **b**
 - all of this can be reduced to $y = \mathbf{x}\mathbf{b} + a$ form for prediction with new data

Regression Summary

- Regression models can be divided into CLS (used when pure analyte spectra are available) and ILS models (MLR, PCR, PLS, RR, CR, ...)
- PCR and PLS work with ill-conditioned data by reducing to a smaller number of factors
 - has advantage of signal averaging
- Cross-validation is used to determine number of factors
- Fit and Prediction are two different things

Model Application

- A PCR or PLS model is applied by
 - centering and scaling to the model mean and variance
 - multiply measurements by regression vector to get scaled predictions
 - rescale the predictions back to original units using model mean and variance
- Prediction outliers can be found by
 - calculating Q and T^2 values for new samples
- All the modeling and application is packaged:
 - the model is an object that contains all the parameters
 - validation e.g.:
`valid = pcr(x,y,model,options); %pred's with new X- & Y-block`
 - prediction e.g.:
`pred = pcr(x,model,options); %predictions with a new X-block`

Model Application – Object Form

- All the modeling and application is packaged:
 - the model is an object that contains all the parameters
 - validation *e.g.*:

```
valid = model.apply(x,y); % this creates a prediction object
pred = model.apply(x); % this does too

valid_scores = valid.plotscores(options);
pred_scores = pred.plotscores(options);
```

This will create dataset objects containing all of the information that you'll see in a scores plot when using the Analysis interface