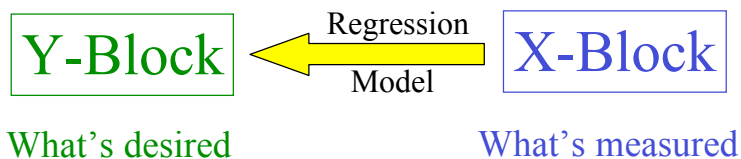


Building Predictive Models: Regression

©Copyright 1996-2008
Eigenvector Research, Inc.
No part of this material may be
photocopied or reproduced in any form
without prior written consent from
Eigenvector Research, Inc.



Regression



Regression analysis identifies the dependency between two blocks of data.

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



Outline

- Nomenclature and conventions
- Classical Least Squares (CLS)
- Inverse Least Squares (ILS) Models
- Multiple Linear Regression (MLR)
- Ridge Regression (RR)
- Principal Components Regression (PCR)
- Partial Least Squares Regression (PLS)
- Determining of the Number of Factors
- Outlier Detection and Model Diagnostics
- A Unifying Theme: Continuum Regression (CR)
- Summary
- Examples

3



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

4



Data Preprocessing

- Everything that was said about preprocessing for PCA goes double for regression
- Data should be linearized, if possible
- Data is almost always mean centered
- Variance scaling used when variables are in different units or greatly different magnitudes
- Outlier elimination is critical to regression models

5



Classical Least Squares

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

$$\mathbf{X} = \mathbf{CS}^T + \mathbf{E}$$

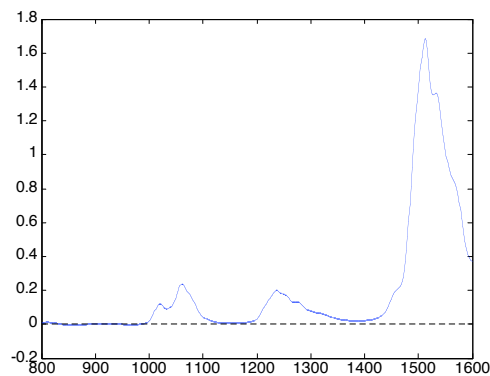
where \mathbf{X} ($m \times n$) is the measured response, \mathbf{S} ($n \times k$) is a matrix of pure component responses, \mathbf{C} ($m \times k$) is a matrix of weights (*i.e.* concentrations), and \mathbf{E} ($m \times n$) is noise or an error matrix

6



The CLS Model

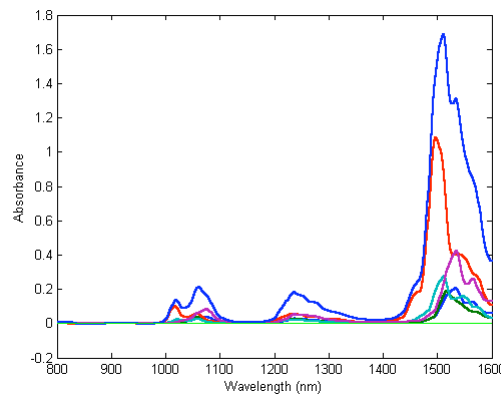
- Given known pure component spectra, how much of each does it take to make up the observed spectrum?



7

The CLS Model

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



8

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} ?

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 \mathbf{S}

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

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

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



11

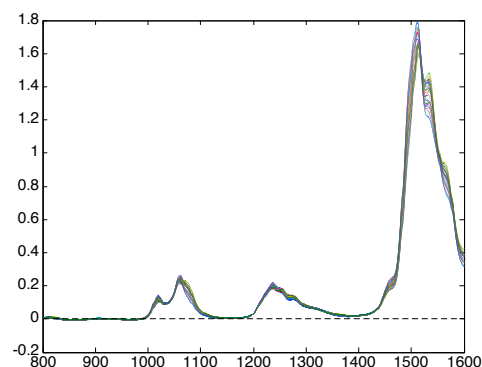
CLS Example

- NIR data of pseudo-gasoline samples
 - absorbance at 401 channels
 - 30 samples
 - 5 analytes

(in analysisGUI..)

File/Load Data/Xblock:
nir_data, "spec1" array

Edit/Plot X-block (Data)



12

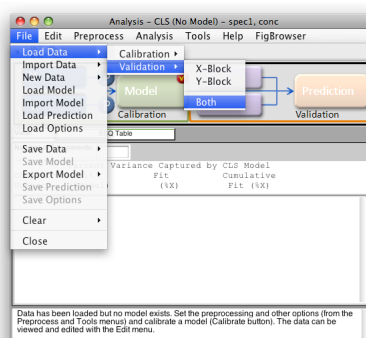
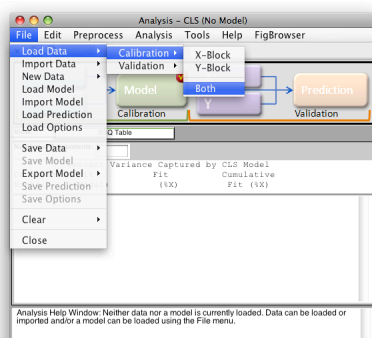
Set up Calibration and Validation Sets

```
>> load_nir_data
>> conc_cal = conc(1:24,:);
>> conc_val = conc(25:30,:);
>> spec_cal = spec1(1:24,:);
>> spec_val = spec1(25:30,:);
```

13



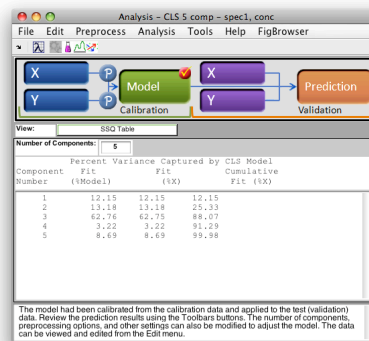
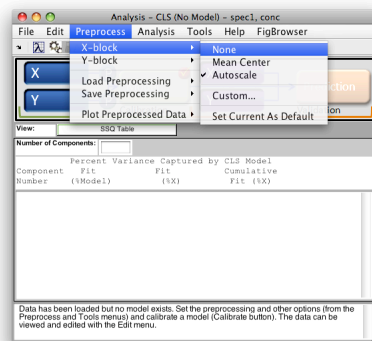
Load into Analysis Interface



14



Set Preprocessing, Calculate Model

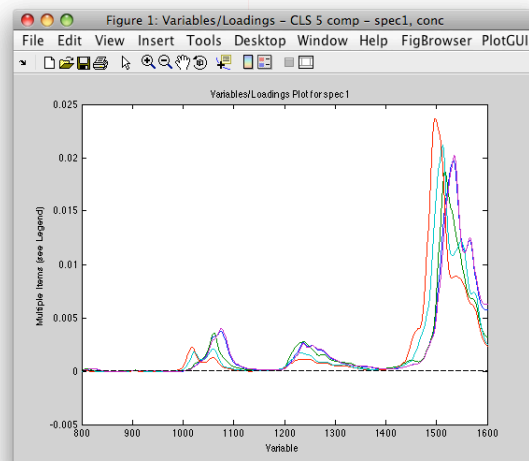
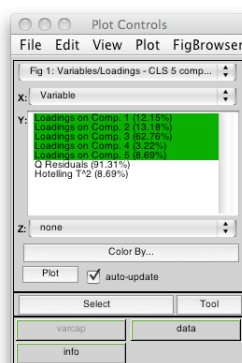


Set preprocessing to “none” in both X and Y
Click “Model” to calculate



15

Pure Component Spectra

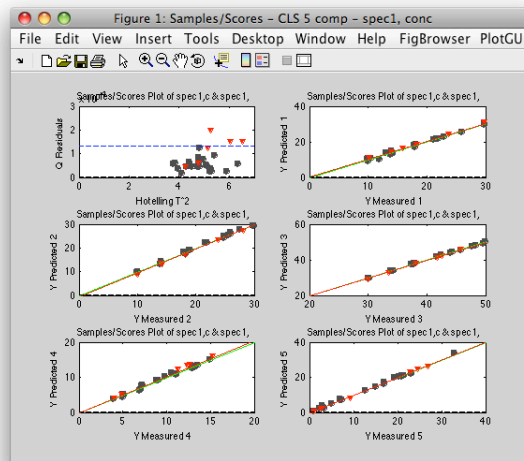


S, estimated from mixtures,
using known concentrations
of all 5 analytes



16

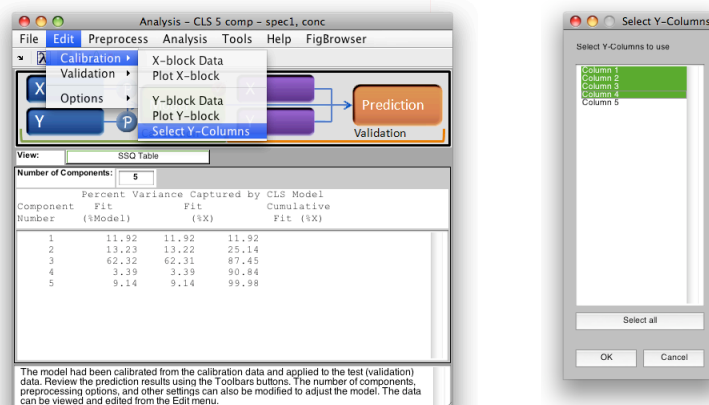
Estimate for Unused (Test) Samples



CLS Problem

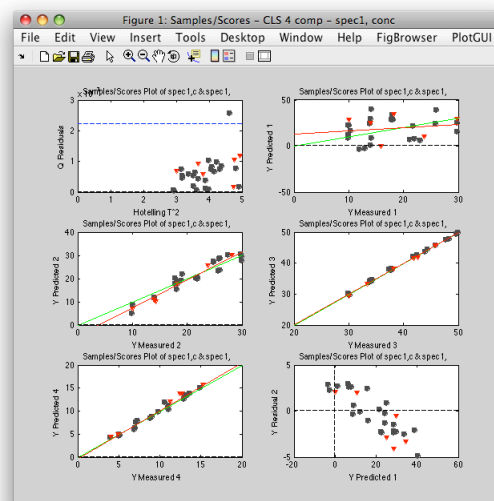
- What if the concentration of 1 analyte was unknown?
- Repeat the CLS procedure using only the first 4 (of 5) analytes
- Attempt to predict concentrations of unused (test) samples

Remove 5th Analyte



19

CLS Solution-Missing Analyte



20

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} ($m \times n$) is the measured response, \mathbf{e} ($m \times 1$) is an error vector, and \mathbf{b} ($n \times 1$) is a vector of coefficients

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

21



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 most obvious is multiple linear regression (MLR), a.k.a. Ordinary Least Squares (OLS)
- In this case \mathbf{X}^+ is estimated from

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

22



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 $n \times n$ and square

23



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

24



Problem with MLR

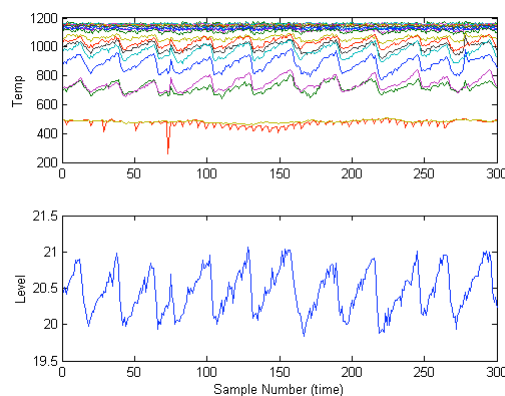
- Inverse of $\mathbf{X}^T \mathbf{X}$ only exists if
 - $\text{Rank}(\mathbf{X}) = nx$, but $\text{rank}(\mathbf{X}) \leq \min(m, nx)$
 - \mathbf{X} has more samples than variables *i.e.* if $m \geq nx$
 - problem with spectra
 - Columns of \mathbf{X} are not co-linear
- Inverse may exist but be highly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

25



MLR Example

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



26



Load and Edit Data

```
>> clear
>> load plsdata
>> whos
```

Name	Size	Bytes	Class
xblock1	300x20	55240	dataset object
xblock2	200x20	38440	dataset object
yblock1	300x1	9008	dataset object
yblock2	200x1	7408	dataset object

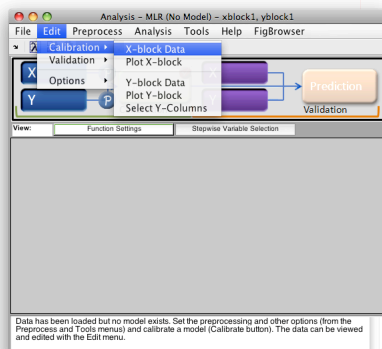
```
>> x = xblock1.data(de\samps([1:300]',[73 167 278 279]),:);
>> y = yblock1.data(de\samps([1:300]',[73 167 278 279]),:);
>> [mx,mnx] = mncn(x); %center the data
>> [my,mny] = mncn(y);
```

27

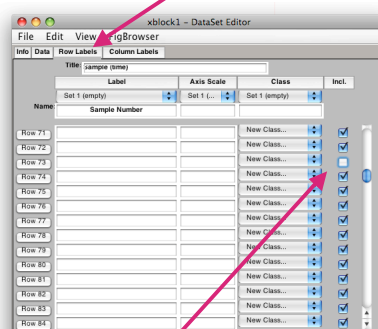


...or via the Analysis GUI

- 1-Load x and y for both calibration and validation
- 2-Select Edit/Calibration/X-block Data



- 3-Select Row Labels tab



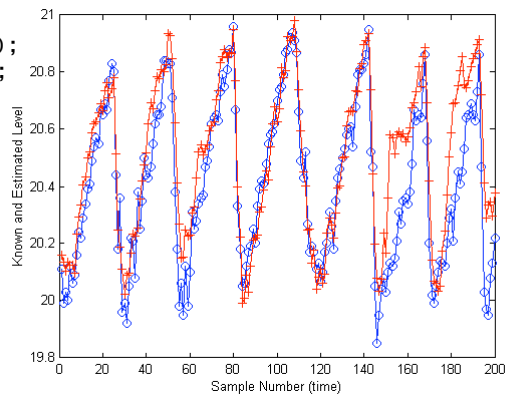
- 4-Deselect samples 73, 167, 278 & 279, then close

28



MLR Regression and Results

```
>> bmlr = mx\my;
>> sx = scale(xblock2.data,mnx);
>> ymlr = rescale(sx*bmlr,mny);
```



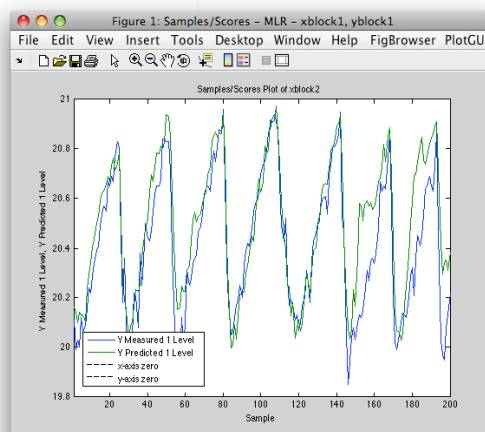
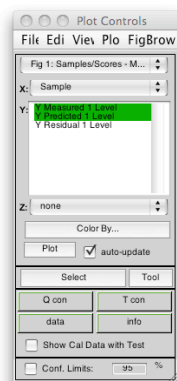
```
>> plot(1:200,yblock2.data,'o-b',1:200,ymlr,'-r+')
>> xlabel('Sample Number (time)')
>> ylabel('Known and Estimated Level')
```

29



MLR Results

- 1-Set X-block preprocessing to Mean Center
- 2-Click Model



- 3-Select Scores

30



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:

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

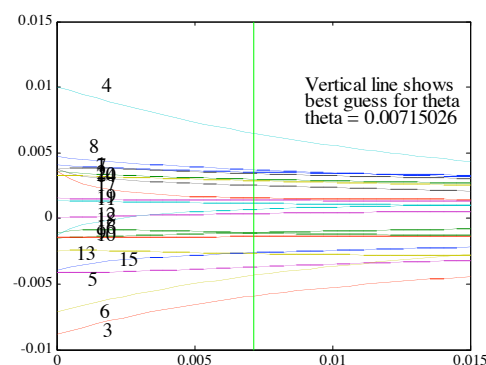
- The addition of Θ stabilizes the inverse and shrinks the values of the coefficients

31



RR Shrinkage

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

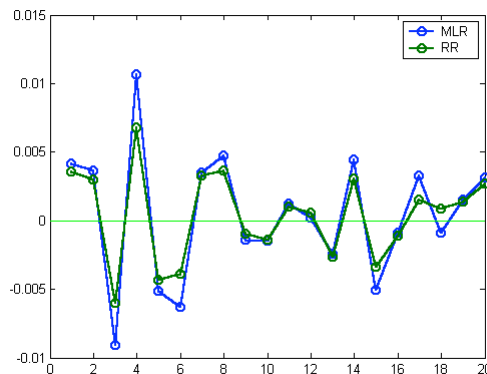


32



RR and MLR Regression Vectors

```
>> plot([bm1r,brr],'o-','linewidth',2)
>> legend('MLR','RR')
>> hline
```



33



Problem with MLR and RR

- RR helps stabilize the inverse but still has problems with strong co-linearity
- Neither MLR nor RR work when $m < nx$
- Possible solution: eliminate variables, *e.g.* stepwise regression or other variable selection
 - how to choose which variables to keep?
 - 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

34



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

35



Principal Components Regression

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

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

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

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

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

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

- Note that $\mathbf{T}_k^T \mathbf{T}_k$ is $k \times k$ and square

36



Cross-Validation

- Divide data set into j sample subsets
- For *each subset* (j times):
 - Build PCA model using all samples in the *remaining* subsets
 - Apply the model to the subset samples
 - Calculate PRESS (Predictive Residual Sum of Squares) for the subset samples:

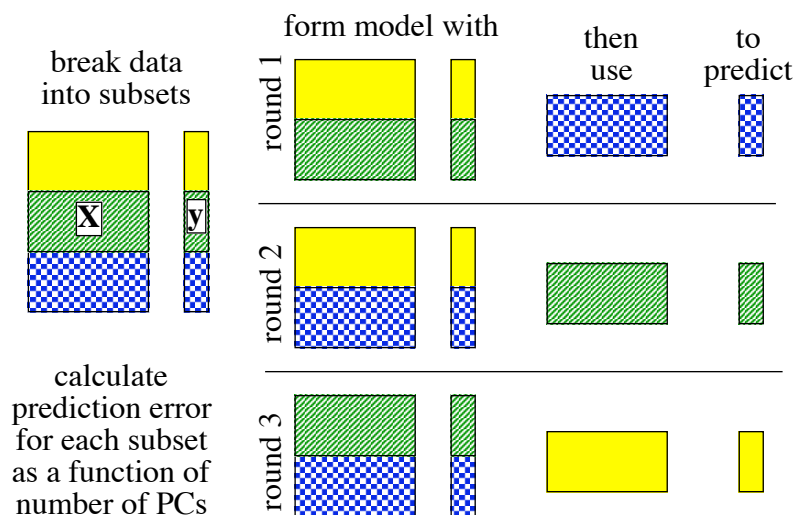
$$e^2 = (y - \mathbf{Xb})^2$$

- Look for minimum or “knee” in PRESS curve

37



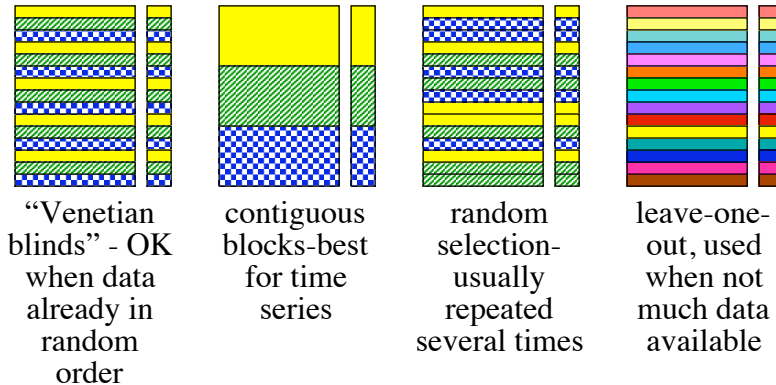
Cross-validation Graphically



38



Formation of Test Sets



What else?

Custom selection, based on prior knowledge!



39

Cross-validation Considerations

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



40

Cross-validation Usage Matrix-I

	Venetian Blinds	Contiguous Blocks	Random Subsets	Leave-One Out	Custom
General 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 determine/construct cross validation array
Small data sets (<~20 objects)	•	•	<ul style="list-style-type: none"> • OK, if many iterations done • Good choice 	<ul style="list-style-type: none"> • Good choice.... •unless designed/DOE data 	<ul style="list-style-type: none"> • often needed to avoid the <i>external subset selection trap</i>
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"> • Can take a while with large m, many iterations 	<ul style="list-style-type: none"> • OK, but.... • Can take a while with large m 	•

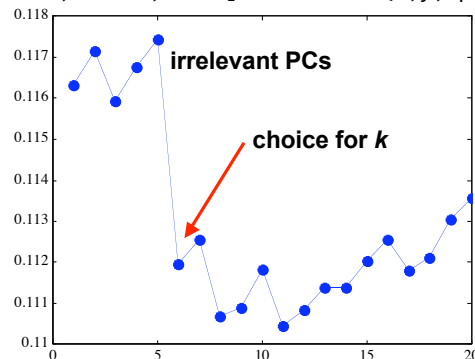
Cross-validation Usage Matrix-II

	Venetian Blinds	Contiguous Blocks	Random Subsets	Leave-One Out	Custom
time-series data	<ul style="list-style-type: none"> • Useful for assessing <i>NON-temporal</i> model errors • Can be optimistic with low number of data splits 	<ul style="list-style-type: none"> • Useful for assessing <i>temporal</i> stability of model 	•	•	•
Batch data	<ul style="list-style-type: none"> • Useful for assessing predictability <i>within</i> batches/parts of batches 	<ul style="list-style-type: none"> • Useful for assessing predictability <i>between</i> 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 <i>replicate sample trap</i> (optimistic results)! 	<ul style="list-style-type: none"> • Good way to avoid <i>replicate sample trap</i> • Beware the <i>external subset selection trap</i>! 	<ul style="list-style-type: none"> • Can use to avoid <i>replicate sample trap</i> (high number of splits) 	<ul style="list-style-type: none"> • overly optimistic results, due to <i>replicate sample trap</i> 	•
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 (<i>external subset selection trap</i>) 	<ul style="list-style-type: none"> • often needed to avoid the <i>external subset selection trap</i>

PCR Cross-Validation Example

- block CV since data is time series
- mean centered, 20 PCs, split 10 times

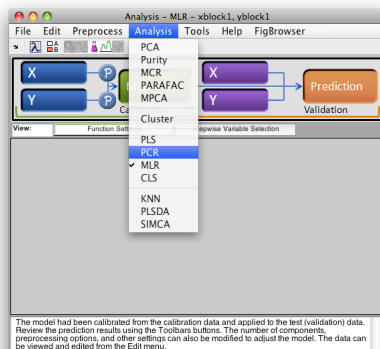
```
>> s = preprocess('meancenter');
>> opts = crossval('options');
>> opts.preprocessing = {s s};
>> [press,cumpress,rmsecv,rmsec] = crossval(x,y,'pcr',{'con',10},20,opts);
```



43

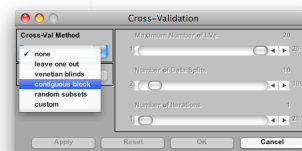
...or via the GUI

1-Reset Analysis to “PCR”



2-Select Tools/Crossvalidation

3-Set to “Contiguous Block”



44

PCR Variance Captured

```
>> options = pcr('options');
>> options.preprocessing = {s s}; % s is mean centering
>> modelpcr = pcr(x,y,6,options);
```

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

45



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**

46



Partial Least Squares

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

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

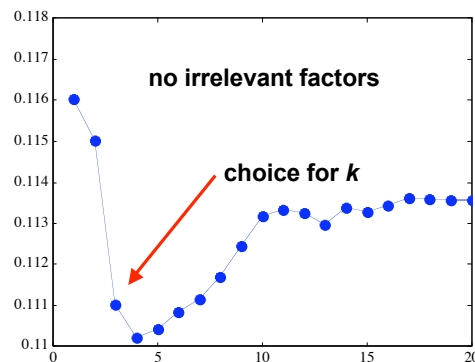
47



PLS Cross-Validation Example

- use block CV since data is time series
- mean centered, 20 PCs, split 10 times

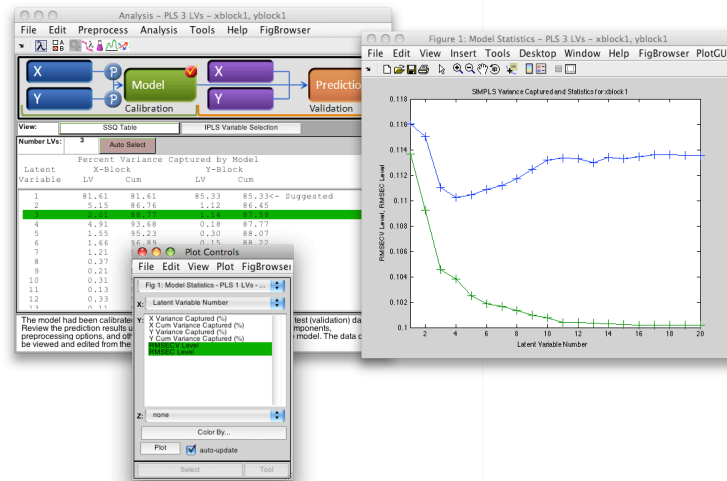
```
>>[press,cumpress,rmsecv,rmsec] = crossval(x,y,'sim',{'con',10},20,opts);
```



48



...or the easy way



49

PLS Variance Captured

```
>> options = pls('options');
>> options.preprocessing = {s s}; % s is mean centering
>> modelpls = pls(x,y,3,options);
Percent Variance Captured by PLS Model
```

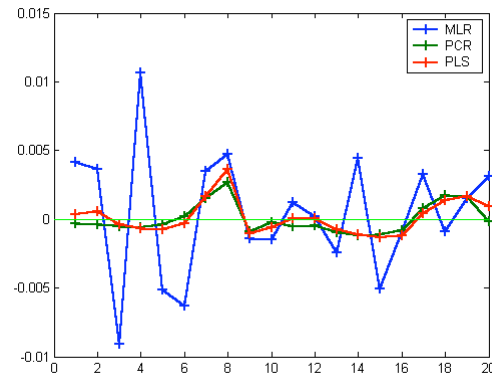
LV #	-----X-Block-----		-----Y-Block-----	
	This LV	Total	This LV	Total
1	81.61	81.61	85.33	85.33
2	5.15	86.76	1.12	86.45
3	2.01	88.77	1.14	87.59
4	4.91	93.68	0.18	87.77
5	1.55	95.23	0.30	88.07
6	1.66	96.89	0.15	88.22
7	1.21	98.10	0.05	88.27
8	0.37	98.47	0.07	88.34
9	0.21	98.68	0.08	88.42
10	0.31	99.00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



50

Regression Vectors

```
>> plot([bmlr modelpcr.reg modelpls.reg],'+-','linewidth',2)
>> legend('MLR','PCR','PLS'), hline
```



51



PLS NIPALS 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

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

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

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

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

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

52



Other PLS Algorithms

- It can be shown that \mathbf{w}_1 is given by

$$\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}_1 = \lambda \mathbf{w}_1$$

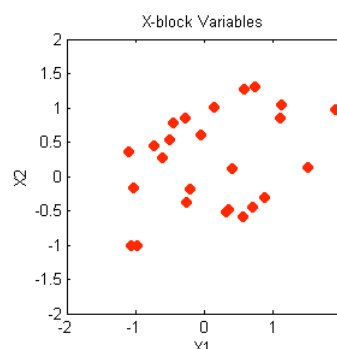
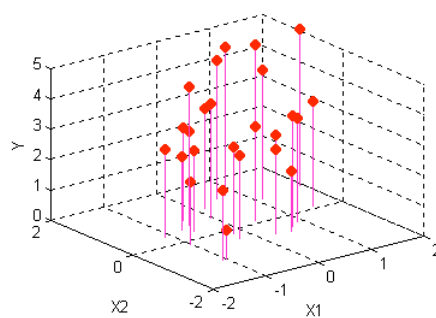
- The SIMPLS algorithm uses an orthogonalization of a Krylov sequence (faster than NIPLS algorithm)
- The important thing to remember is:

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

53



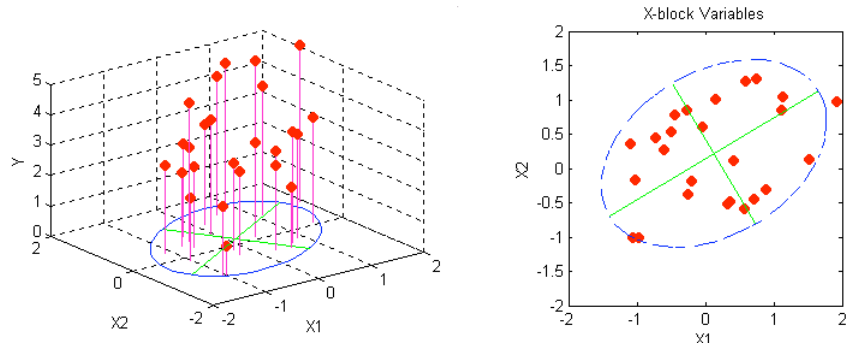
Y Projected onto X Plane



54



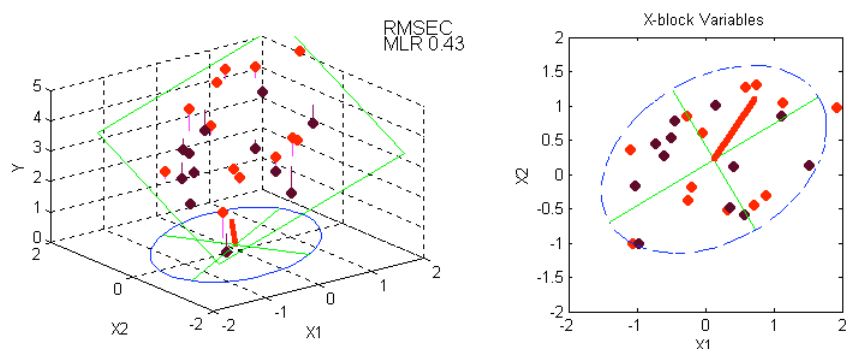
PCA of X-Block



55



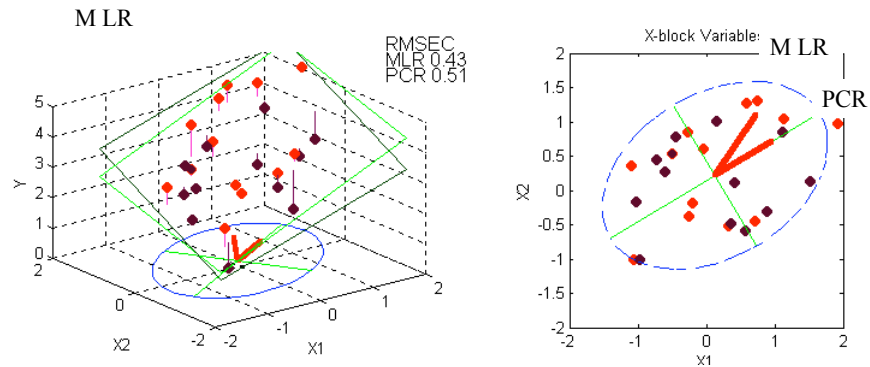
MLR Regression Vector and Surface



56



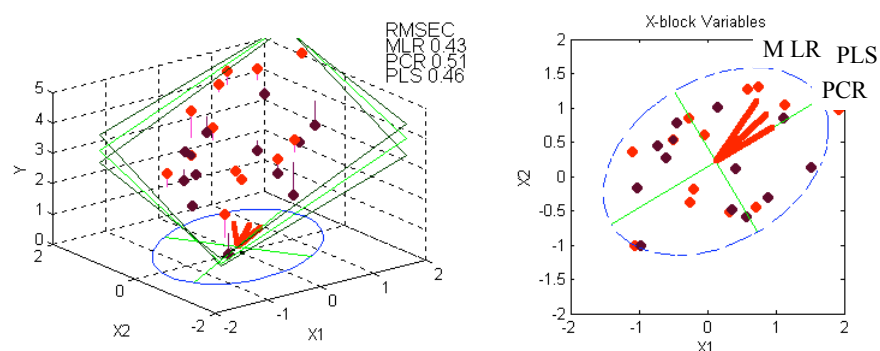
PCR Regression Vector and Surface



57



PLS Regression Vector and Surface

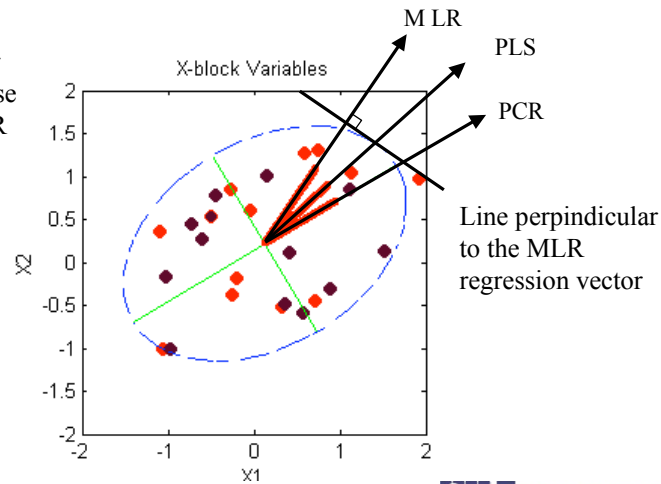


58



Geometric Relationship of MLR, PCR, and PLS

PLS is the vector on the PCR ellipse upon which MLR has the longest projection



PLS for Multivariate \mathbf{Y}

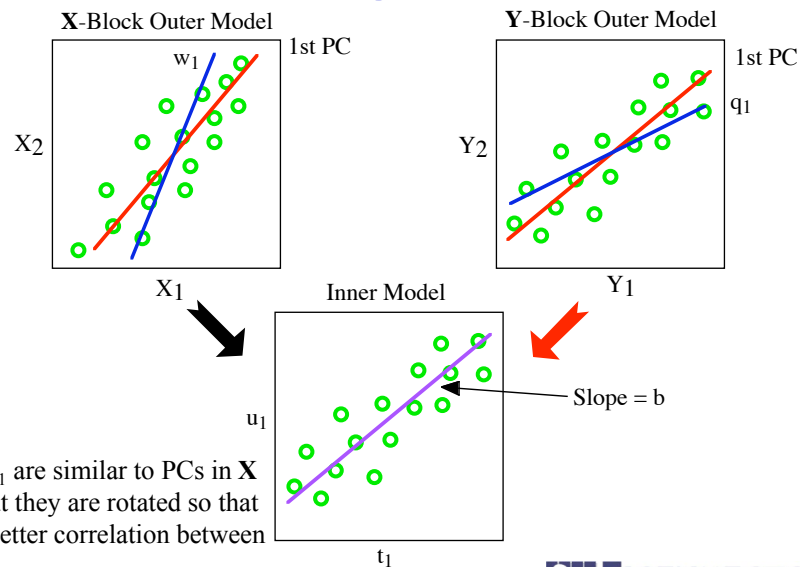
- PLS can be used to relate multivariate \mathbf{X} to multivariate \mathbf{Y} (*a.k.a.* PLS2)
 - outer relationships

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

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

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

PLS2



Model Quality Measures

- Root Mean Square Error (RMSE) Metrics
 - RMSE C
 - RMSE CV
 - RMSE P
 - In units of the Y variable!
- Correlation Coefficient (r)
 - Unit-less
 - Considers the range of Y

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

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

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)

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

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

63



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

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

j's refer to different CV subsets

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

64



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)

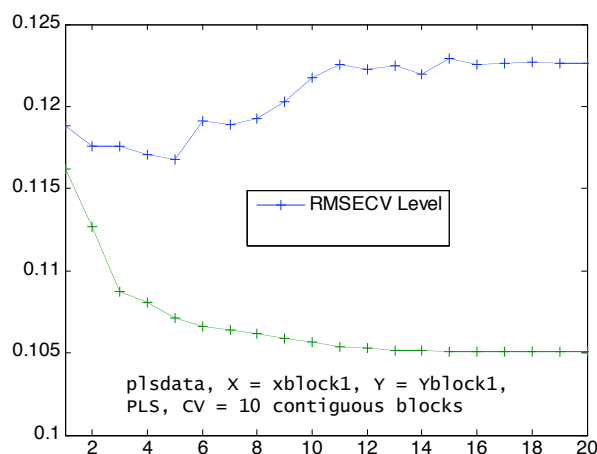
$$\text{RMSEP} = \sqrt{\frac{\sum_{i=1}^{m_p} (y_i - \hat{y}_i)^2}{m_p}}$$

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

65



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

66



Comparison of Models

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

	<u>MLR</u>	<u>PCR</u>	<u>PLS</u>
RMSEC	0.0991	0.1059	0.1034
RMSECV	0.1122	0.1108	0.1098
RMSEP	0.1496	0.1366	0.1396

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

67



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 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?

68

