

Variable Selection

©Copyright 2007-2017
Eigenvector Research, Inc.
No part of this material may be
photocopied or reproduced in any form
without prior written consent from
Eigenvector Research, Inc.



2

Outline

- Why select variables?
- What methods are there?
- So which to choose?
- In practice
- Some examples



Course Materials

- These slides
- PLS_Toolbox or Solo 6.7 or later
- Data sets
 - From DEMS folder (distributed with software)
 - nir_data.mat (optional)
 - From EVRIHW folder (additional data sets)
 - beer.mat, nir_shootout_2002.mat (optional)

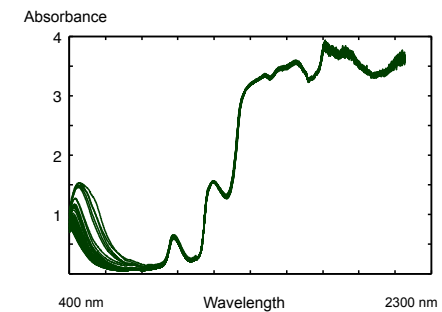
3



VIS/NIR spectra of 61 beers

Purpose: prediction of real extract

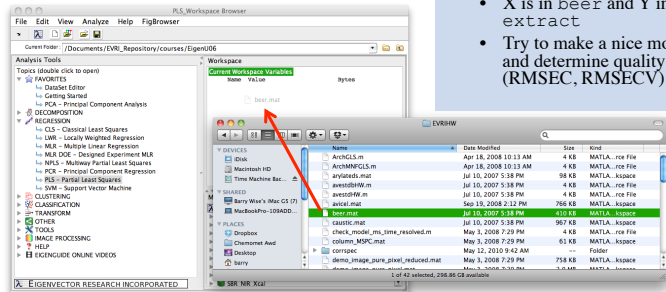
Why select
What methods
Choose method
In practice
Some examples



Why select
What methods
Choose method
In practice
Some examples

VIS/NIR spectra of 61 beers

Try to make a PLS model for extract

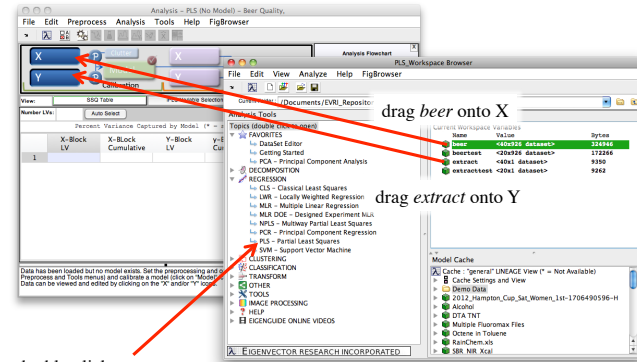


- Load beer.mat
- X is in beer and Y in extract
- Try to make a nice model and determine quality (RMSEC, RMSECV)

5



Start PLS, Load Beer



double-click to
start PLS

6



Why select
What methods
Choose method
In practice
Some examples

Exercise Data

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

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

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

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

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

7



Why select
What methods
Choose method
In practice
Some examples

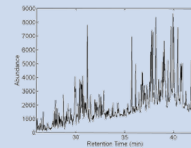
Why variable selection?

- Improvement of the model
 - Remove irrelevant, unreliable or noisy variables
 - Improve predictions
 - Improve statistical properties

- Interpretation
 - Obtain a model that is easier to understand

- Costs
 - Use less measurements to replace expensive or time-consuming one

- Development of fast instruments/routines for on-line control
 - Find wavelength ranges for a filter-based instrument



8



What methods available?

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



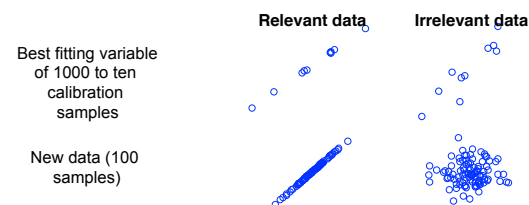
EIGENVECTOR
RESEARCH INCORPORATED

9

Why select
What methods
Choose method
In practice
Some examples

What methods?

- Why not just choose the best variables?
 - Highly nonlinear problem
 - Exhaustive search not possible
 - How to validate what's good?



10

Why select
What methods
Choose method
In practice
Some examples

Method : A priori Choose the right measurements

- The most important of all
- Beyond the scope of this course as we assume the data are already available/fixed

Important assumptions

- There may be indirect correlations that you did not anticipate
- Don't choose too few variables a priori

EIGENVECTOR
RESEARCH INCORPORATED

11

Why select
What methods
Choose method
In practice
Some examples

Method: A Priori Example indirect relation

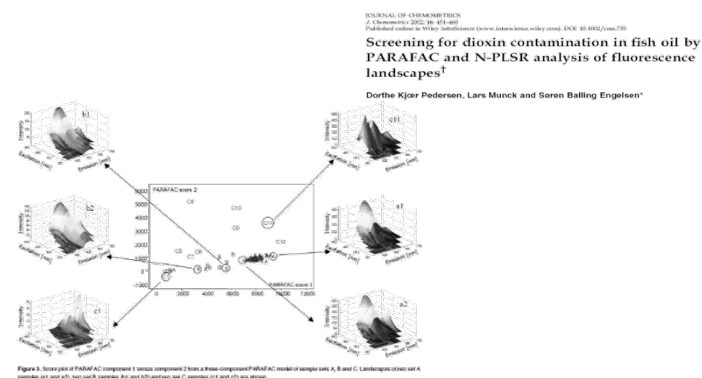


Figure 3. Score plot of PARAFAC component 1 versus component 2 from a three component PARAFAC model of sample sets A, B and C. Landscapes of these set A samples (01 and 02), set B samples (01 and 02) and set C samples (01 and 02) are shown.

EIGENVECTOR
RESEARCH INCORPORATED

12

Method: A Priori Example indirect relations

Even though not direct link – fluorescence works well!!!

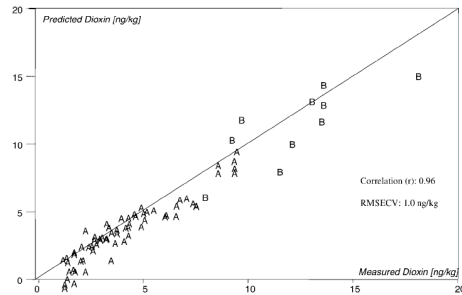


Figure 5. Predicted dioxin concentration versus measured dioxin concentration for N-PLSR model (four PCs) for sample set A + B fish oil samples ($n = 75$). The correlation coefficient (r) and prediction error (RMSECV) are reported.



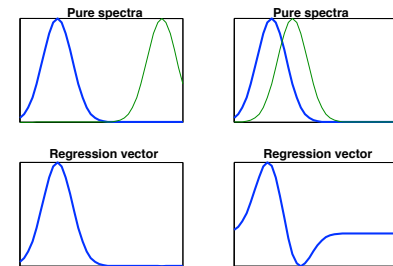
13

Why select
What methods
Choose method
In practice
Some examples

Method : A posteriori

Why select
What methods
Choose method
In practice
Some examples

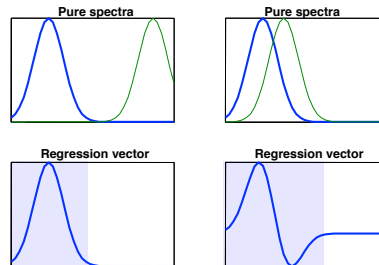
- Remember that a regression model is not only depending on the analyte directly
- Also has to adjust for overlapping signals



14

Method : A posteriori

- Remember that a regression model is not only depending on the analyte directly
- Also has to adjust for overlapping signals



15

Why select
What methods
Choose method
In practice
Some examples

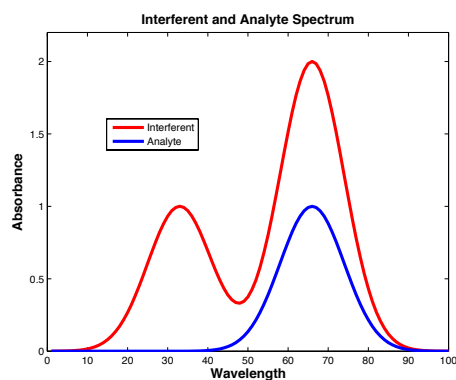
Why is Variable Selection Hard?

- Is much easier when analyte of interest has “pure variables,” *i.e.* variables where the analyte has a signal but no interferences do
- When interferences overlap with analyte, important variables can be positively correlated with concentration, negatively correlated, or uncorrelated!



16

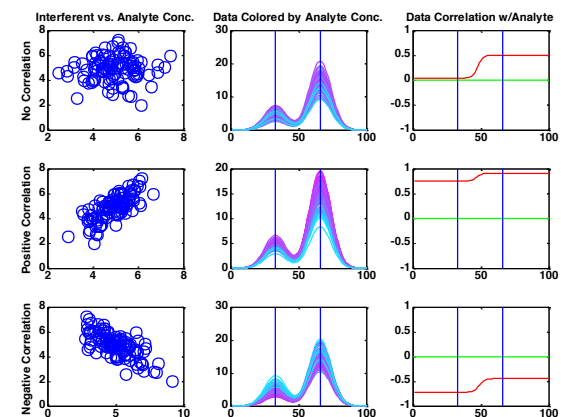
Simple System Example



EIGENVECTOR
RESEARCH INCORPORATED

17

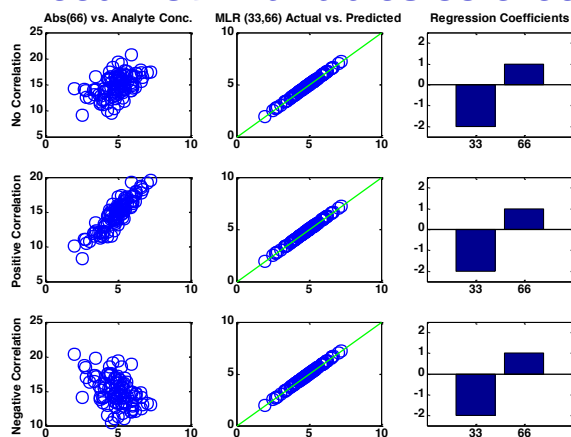
Consider 3 Cases



EIGENVECTOR
RESEARCH INCORPORATED

18

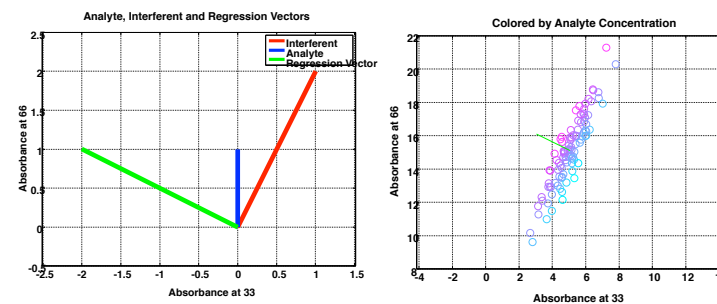
Need Both Variables 33 & 66!



EIGENVECTOR
RESEARCH INCORPORATED

19

Regression Coeffs and Data



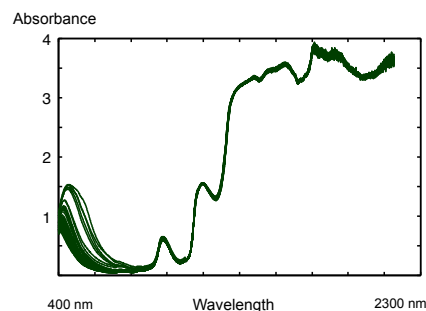
Model needs signal at 33 to be orthogonal to interferent!

EIGENVECTOR
RESEARCH INCORPORATED

20

Example: A posteriori
VIS/NIR spectra of 61 beers
 Purpose: prediction of real extract

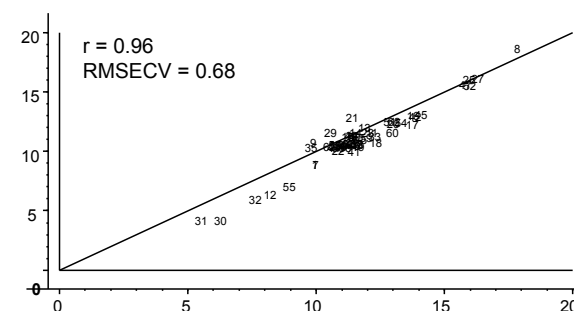
Why select
 What methods
 Choose method
 In practice
 Some examples



EIGENVECTOR
 RESEARCH INCORPORATED

Example: A posteriori
Full spectrum PLS-model
 Cross validated prediction error (61 samples, 6 segments)
 Five PLS factors

Why select
 What methods
 Choose method
 In practice
 Some examples

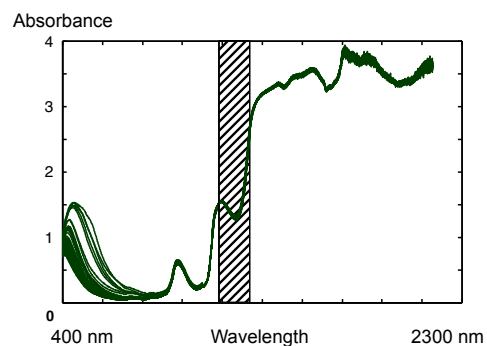


EIGENVECTOR
 RESEARCH INCORPORATED

22

Example: A posteriori
Selected interval: 1218-1300 nm

Why select
 What methods
 Choose method
 In practice
 Some examples



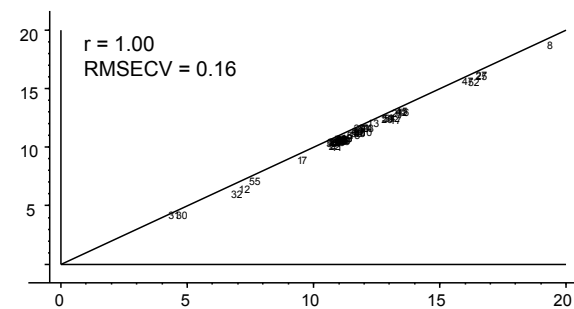
EIGENVECTOR
 RESEARCH INCORPORATED

23

Example: A posteriori
PLS-model based on 1218-1300 nm

Why select
 What methods
 Choose method
 In practice
 Some examples

Cross validated prediction error (61 samples, 6 segments)
 Three PLS factors



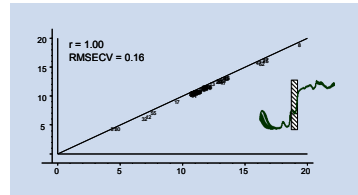
EIGENVECTOR
 RESEARCH INCORPORATED

24

Example: A posteriori PLS-model based on 1218-1300 nm

Why select
What methods
Choose method
In practice
Some examples

This course is about how to achieve results as these even in situations where such detailed background knowledge is not accessible.



25

Model based selection Important assumptions

Why select
What methods
Choose method
In practice
Some examples

- **Model is reasonable**
 - if 900 out of 1000 variables irrelevant, the model may be reflecting those and hence the relevant ones look insignificant
- **Model is certain**
 - Few samples or noisy measurements =>
 - Statistical uncertainty high =>
 - Do not trust the model parameters too much



27

Model based selection

Why select
What methods
Choose method
In practice
Some examples

- Simply use the visual appearance of the model
- E.g. small loadings, low regression coefficients etc.

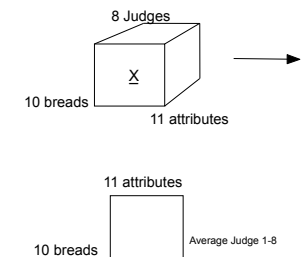


26

Model based selection Example – Sensory analysis

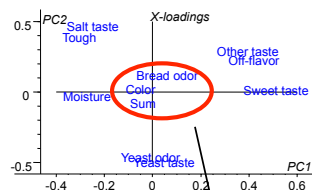
Why select
What methods
Choose method
In practice
Some examples

- Sensory profiling of bread
 - 10 breads (replicates) \times 11 attributes \times 8 judges
 - Average over judges: 10 \times 11 attributes
 - Data from Magni Martens



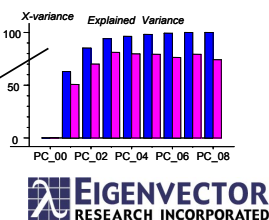
28

Model based selection Example



Loadings indicate some variables not important

Seems trustworthy as the model is otherwise well-behaved

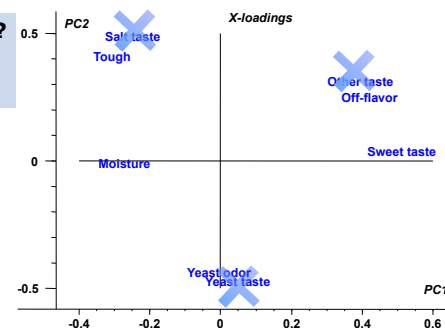


Why select
What methods
Choose method
In practice
Some examples

29

Model based selection Example - Many ways to skin a cat

Want to get rid of more?
Remove redundant variables



EIGENVECTOR
RESEARCH INCORPORATED

30

Model based selection Automating it a bit

Variable importance for projection (VIP)

Relative weighted sum of squares of PLS-weights, weighted by components importance for predicting

- Assumes valid model
- Hence only remove few variables at a time
- And check for outliers etc. along the way
- VIP smaller than one indicates low importance

Why select
What methods
Choose method
In practice
Some examples

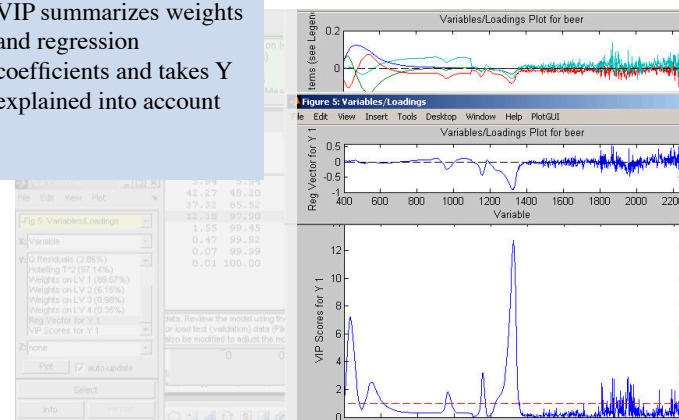
31

S. Wold, E. Johansson, M. Cocchi, 3D QSAR in Drug Design: Theory, Methods, and Applications, ESCOM, Leiden, Holland, 1993, pp. 523-550.

EIGENVECTOR
RESEARCH INCORPORATED

Model based selection Automating it a bit

VIP summarizes weights and regression coefficients and takes Y explained into account



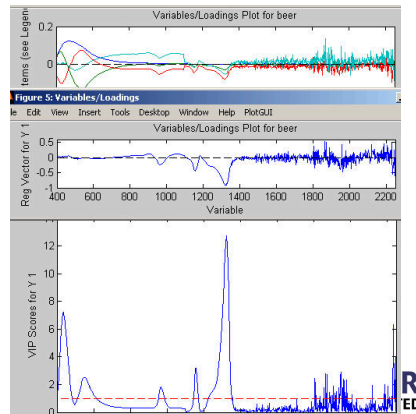
32

Model based selection Automating it a bit

Why select
What methods
Choose method
In practice
Some examples

Try yourself on the
beer data.

Is a one-shot
selection optimal?



33

Genetic algorithm

Why select
What methods
Choose method
In practice
Some examples

- **Method**
 - Survival of the fittest (best fit)
- **Principle**
 - Every combination of variables (a model) is defined by an index of which variables to use *and* the goodness of this model
 - Example: One species is given by the calibration model with variable 1, 3, 14 & 27 yields and RMSEP of 1.23.
 - Find better models by “mating” species so that the good ones mate more than bad ones (survival of the fittest).
 - Example: Calibration models with variable 8 are generally better and therefore increasingly part of new calibration models.



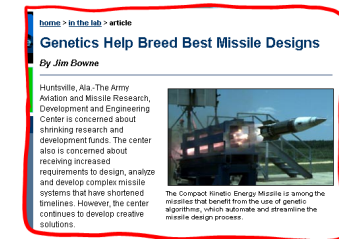
35

“Random” variable selection

Random selection means choosing without considering improvement in fit directly

Increases chance of obtaining unforeseen and complex interactions

Genetic algorithms is a good example



www.rdecom.army.mil



34

Genetic algorithm

Why select
What methods
Choose method
In practice
Some examples

- **Terminology**
 - Population = Set of individuals
 - One individual = Model with a given set of variables
 - One gene = Codes for one variable (in/out)
- **Algorithm**
 - Make start population (e.g. 50 different models with different variables included)
 - Evaluate each model (RMSEP or similar)
 - Have a party and let the best ones have most fun
 - Fun: two models mate and make a child which has similar variables
 - Arrange a new party for the 50 children and continue

$$\begin{bmatrix} y \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} X \end{bmatrix}$$

Gene for one model defines which variables are in. RMSEP defines the quality of the individual

Riccardo Leardi has written many papers on how to make genetic algorithms work



36

Random methods Assumptions

Why select
What methods
Choose method
In practice
Some examples

- "Random" methods are based on combining a random search (individuals) with a guided search (mating)
 - Mostly used because exhaustive search is too expensive
 - Good and sound principle
 - Excellent for getting ideas
 - Not good, though, for refining

37

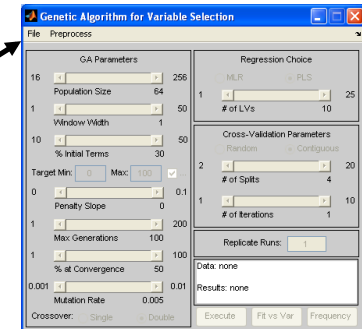


Genetic algorithms Exercise

Why select
What methods
Choose method
In practice
Some examples

In MATLAB
>> load beer
>> genalg

- Load calibration data (beer and extract)
- Execute



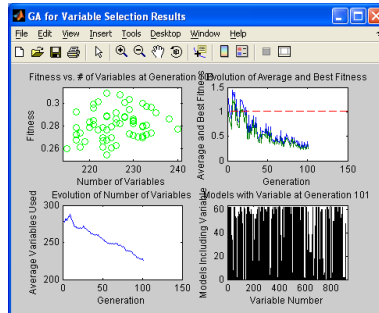
38



Genetic algorithms Exercise

Why select
What methods
Choose method
In practice
Some examples

- Does the result look nice?
- Maybe try increasing iterations?



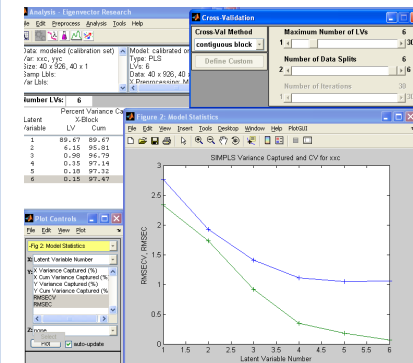
39



Genetic algorithms Exercise

Why select
What methods
Choose method
In practice
Some examples

- Generally
- Lower the number of components based on initial PLS analyses but use a little more
- Keep ending criteria sensitive. The more iterations which occur, the more feedback from the cross-validation information, thus more likely over-fitting,
- Use random cross-validation and multiple iterations if practical,
- Repeat the GA run multiple times and watch for general trends
- For data with many variables and fewer samples, increase the window width



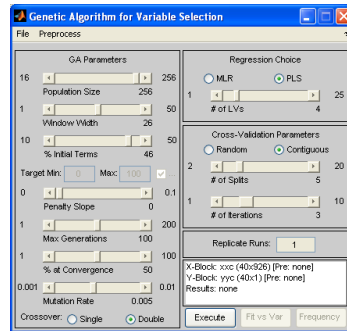
40



Genetic algorithms Exercise

Why select
What methods
Choose method
In practice
Some examples

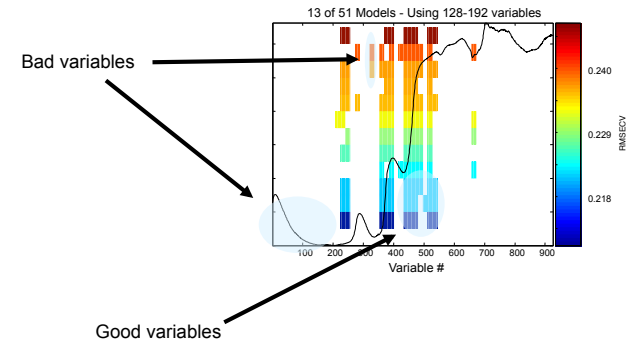
- **In PLSToolbox**
- **Size of Population** Larger populations provide a better representation of different variable combinations.
- **Window Width** When adjacent variables contain correlated information the original variables can be included or excluded in "blocks".
- **% Initial Terms** Appr. # variables included in the initial subsets. Few initial terms will make identification of useful variables more difficult, but will bias the end solution towards models with fewer variables.
- **Target Min/Max & Penalty slope** For guiding towards a specific number of variables



41

Genetic algorithms Exercise

Why select
What methods
Choose method
In practice
Some examples

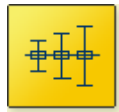


42

Classical methods

Why select
What methods
Choose method
In practice
Some examples

- Subset selection
 - Statistical significance tests
 - Forward selection
 - Backward selection
 - Best subset selection



43

Classical methods Statistical significance tests

Why select
What methods
Choose method
In practice
Some examples

Principle

- Do regression
- Eliminate variables with non-significant regression coefficients

Properties

- Very fast
- Assumes a statistically valid model
- Assumes statistically valid significance tests
- Hence only works for models with insignificant amount of irrelevant variables

The function Calibsel selects variables based on significance



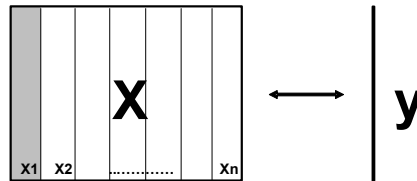
44

Intermezzo

Why select
What methods
Choose method
In practice
Some examples

iPLS: Interval PLS

Local models in n intervals. Very intuitive and useful approach that can be easily combined with variable selection



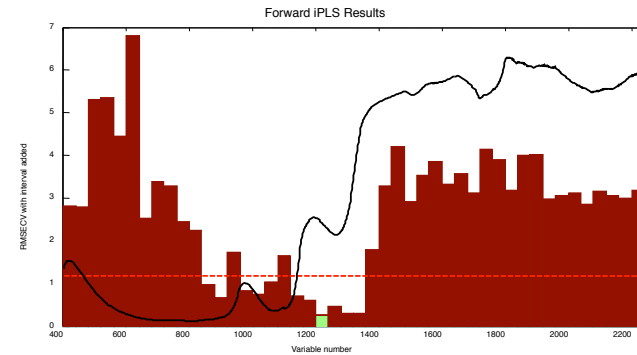
L. Nørgaard, A. Saudland, J. Wagner, J. P. Nielsen, L. Munck, S. B. Engelsen. Interval partial least-squares regression (iPLS). *Appl. Spectrosc.* 54 (3):413-419, 2000.



45

iPLS result on beer data

iPod introduced 2001
iPLS introduced 2000
Lawsuit?!



Classical methods Forward selection

Principle

- Select best fitting variable (or better cross-validated)
- Regress y on this and select variable fitting best on residual
- Regress y on these two and select best fitting on the residual
- Etc.

Good

- Fast
- Handles many irrelevant variables
 - If many samples or test set evaluation
- Often works well

Bad

- Disregards interactions to some extent



Forward selection Exercise

Try forward selection on NIR data

Many variables – computationally expensive

- If correlation between neighbors use windows instead of individual variables.
- E.g. use every 10 neighbors as one set and ex/ include them all together

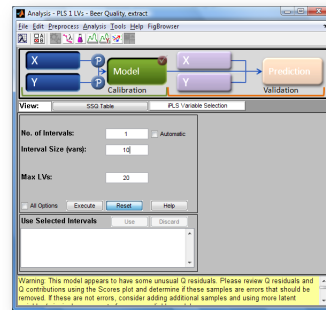


Forward selection Exercise

Maybe more than one interval could be useful?

Maybe bigger or smaller intervals are better?

Investigate ..

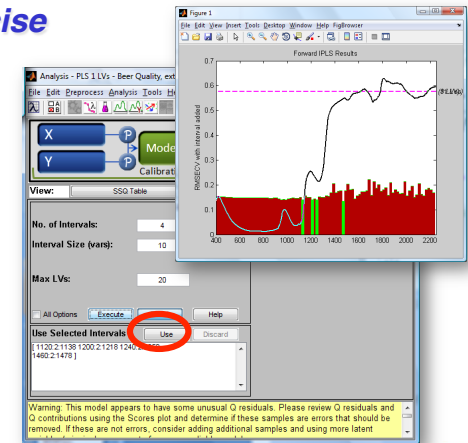


EIGENVECTOR
RESEARCH INCORPORATED

Forward selection Exercise

Once you have a good model, either

- 1) Quick and dirty: use the good intervals
- 2) Better: Remove the clearly bad intervals



EIGENVECTOR
RESEARCH INCORPORATED

Backward selection Exercise

Rather than forward selection, you can choose advanced options and rather than *selecting* variables, you can *de-select* variables by choosing the reverse mode

Principle

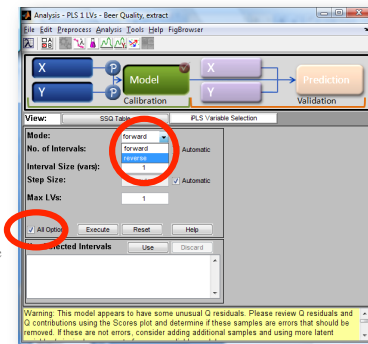
- Make full model and remove the variable contributing the least to the fit
- Repeat that

Good

- Takes interactions into account

Bad

- Often less efficient than forward selection



EIGENVECTOR
RESEARCH INCORPORATED

Choosing method If there are many irrelevant variables

Why select
What methods
Choose method
In practice
Some examples

- Do not use statistically based tests
- The fewer samples, the more important the independence of the optimization criterion is
 - Do not use fit values but rather cross-validation or test set
- Genetic algorithms and forward selection are good choices
- Refine solution using background knowledge and other variable selection tools

EIGENVECTOR
RESEARCH INCORPORATED

Choosing method If refining an already ok model

Why select
What methods
Choose method
In practice
Some examples

- The more samples (to variables), the less important optimization criterion is
- Statistics makes sense
 - Significance is very helpful, model parameters can be interpreted meaningfully
- Use *i*PLS, backward selection, model statistics such as VIP etc.

53



How to do it in practice

Why select
What methods
Choose method
In practice
Some examples

By now it is clear that most methods do not like outliers:

Remove outliers!!

Even if they are not finally removed, meaningful variable selection is not possible with even slightly significant outliers (unless very many samples).

So get rid of them and add them again afterwards

54

Alternative – robust methods



How to do it in practice

Why select
What methods
Choose method
In practice
Some examples

- **Never trust results**
 - Use several methods as inspiration
 - Evaluate selected variables
 - Add some, remove some, check with a priori knowledge
- **The better the initial model, the better results from variable selection**
 - Remove obviously irrelevant variables
 - Remove outliers
 - Bin similar variables

55

