

Variable Selection

©Copyright 2006
Eigenvector Research, Inc.
No part of this material may be
photocopied or reproduced in any form
without prior written consent from
Eigenvector Research, Inc.



Outline

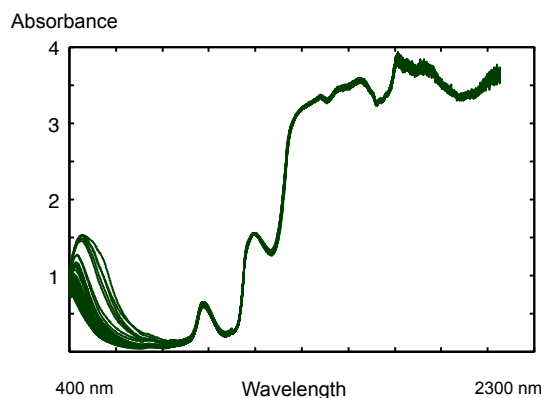
- **Why select variables?**
- **What methods are there?**
- **So which to choose?**
- **In practice**
- **Some examples**



Why select
What methods
Choose method
In practice

VIS/NIR spectra of 61 beers

Purpose: prediction of real extract



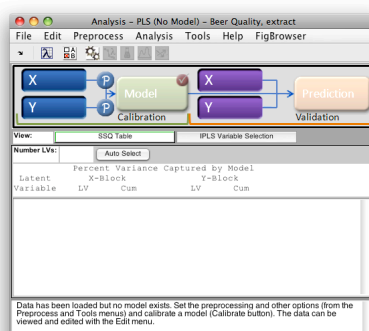
EIGENVECTOR
RESEARCH INCORPORATED

Why select
What methods
Choose method
In practice

VIS/NIR spectra of 61 beers

Try to make a PLS model for extract

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



EIGENVECTOR
RESEARCH INCORPORATED

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.

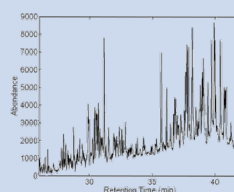
Data sorted by extract value, and a model independent test set was constructed by selecting every third sample of this full data set. There are thus two data sets: one for calibration (40 samples) and one for independent estimation of prediction error (20 samples).



5

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



6

Why select
What methods
 Choose method
 In practice

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

7

Why select
What methods
 Choose method
 In practice

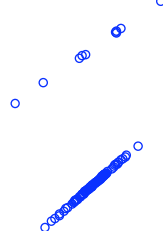
What methods?

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

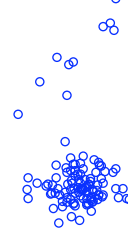
Best fitting variable
 of 1000 to ten
 calibration
 samples

New data (100
 samples)

Relevant data



Irrelevant data



8

Why select
What methods
Choose method
In practice

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

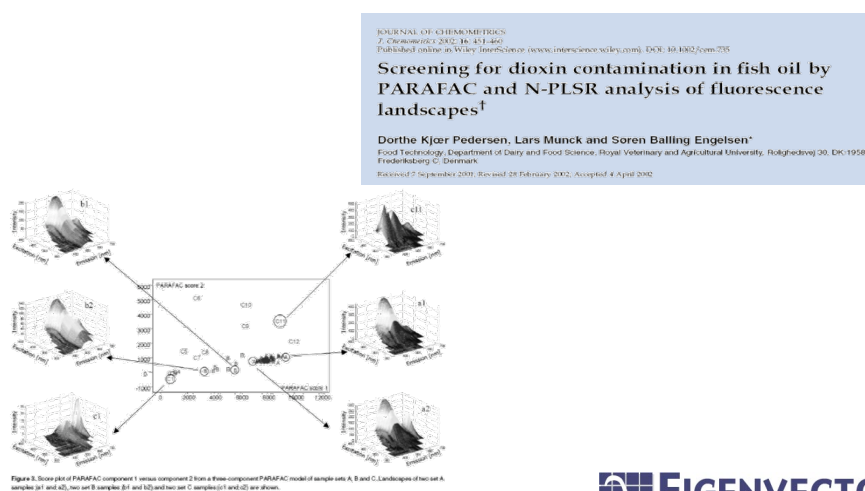
- There may be indirect correlations that you did not anticipate
- Don't choose too few variables a priori



9

Why select
What methods
Choose method
In practice

Method: A Priori Example indirect relation



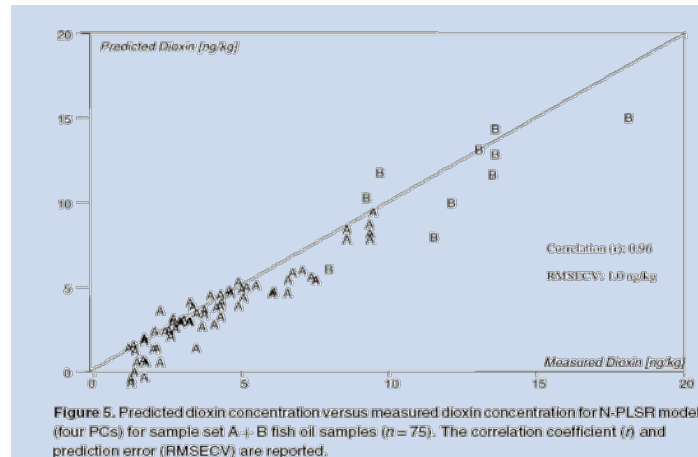
10

Method: A Priori

Example indirect relations

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

Why select
What methods
Choose method
In practice



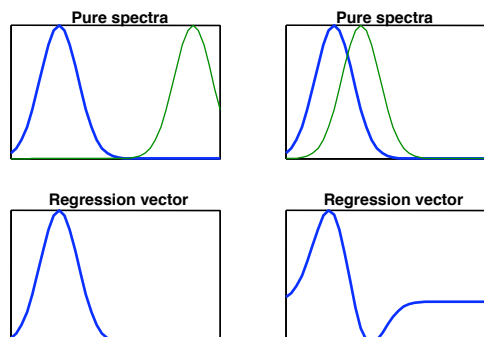
11

EIGENVECTOR
RESEARCH INCORPORATED

Method : A posteriori

Why select
What methods
Choose method
In practice

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



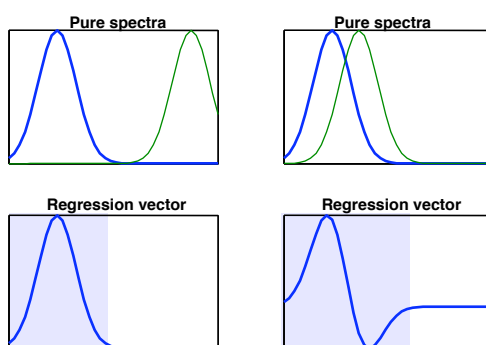
12

EIGENVECTOR
RESEARCH INCORPORATED

Why select
What methods
Choose method
In practice

Method : A posteriori

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



EIGENVECTOR
RESEARCH INCORPORATED

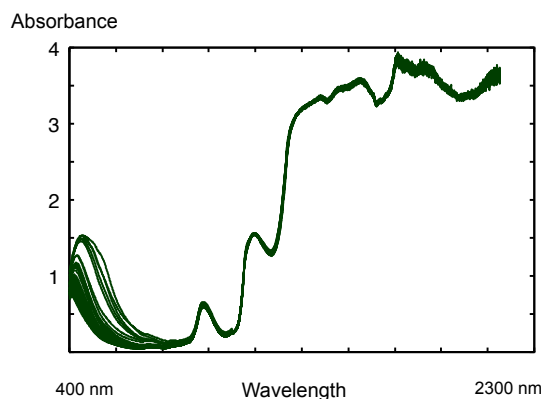
13

Why select
What methods
Choose method
In practice

Example: A posteriori

VIS/NIR spectra of 61 beers

Purpose: prediction of real extract



EIGENVECTOR
RESEARCH INCORPORATED

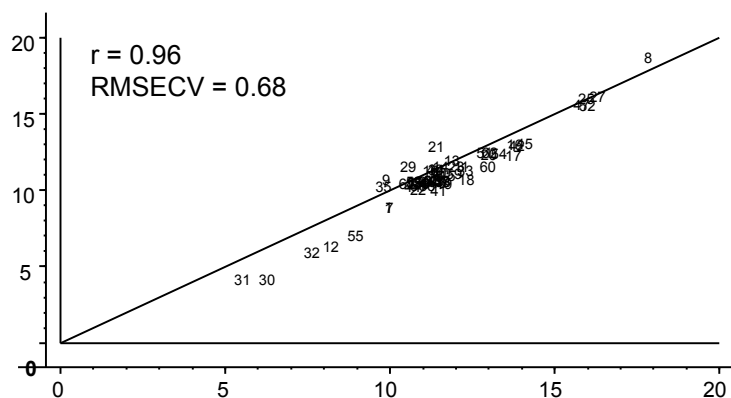


Example: A posteriori

Full spectrum PLS-model

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

Why select
What methods
Choose method
In practice

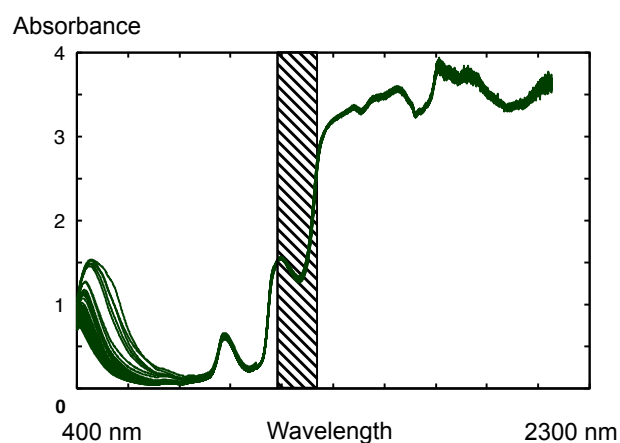


15

Example: A posteriori

Selected interval: 1218-1300 nm

Why select
What methods
Choose method
In practice

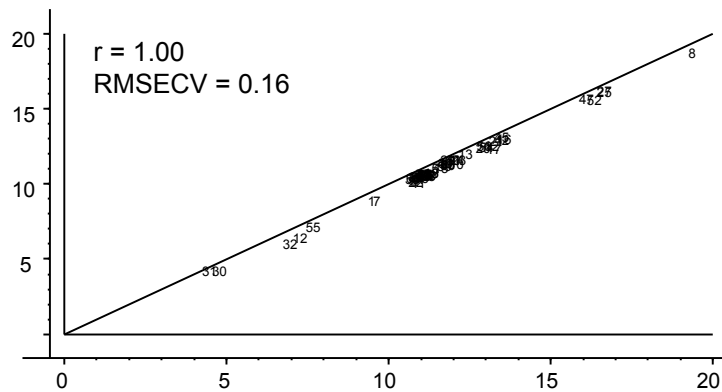


16

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

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

Why select
What methods
 Choose method
 In practice



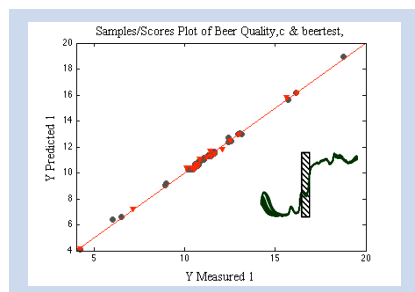
EIGENVECTOR
 RESEARCH INCORPORATED

17

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

Why select
What methods
 Choose method
 In practice

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



RMSEC = 0.13
 RMSECV = 0.15
 RMSEP = 0.17

EIGENVECTOR
 RESEARCH INCORPORATED

18

Model based selection

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

19



Model based selection *Important assumptions*

- **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

20

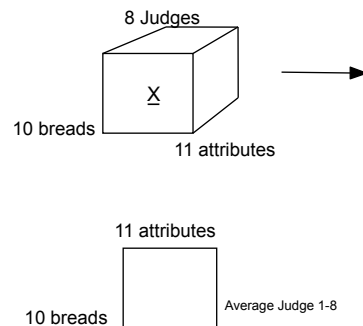


Why select
What methods
Choose method
In practice

Model based selection Example – Sensory analysis

■ Sensory profiling of bread

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

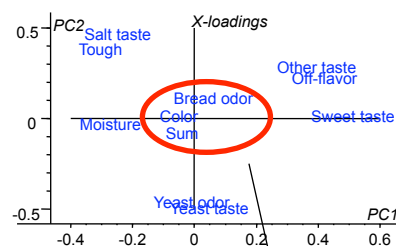


21



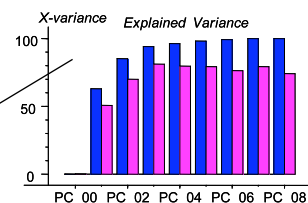
Why select
What methods
Choose method
In practice

Model based selection Example



Loadings indicate some variables not important

Seems trustworthy as the model is otherwise well-behaved



22



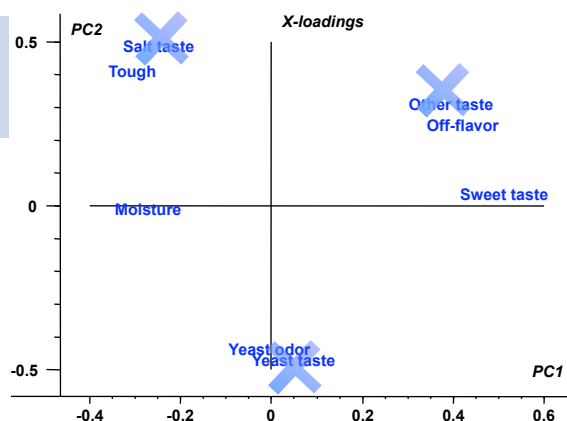
Model based selection

Example - Many ways to skin a cat

Why select
What methods
Choose method
In practice

Want to get rid of more?

Remove redundant variables



EIGENVECTOR
RESEARCH INCORPORATED

23

Model based selection

Automating it a bit

Why select
What methods
Choose method
In practice

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

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

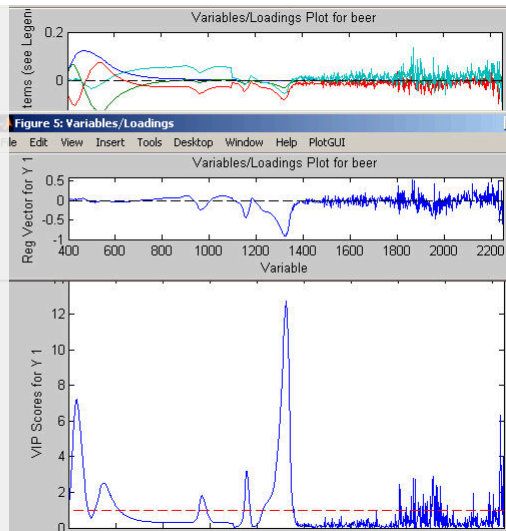
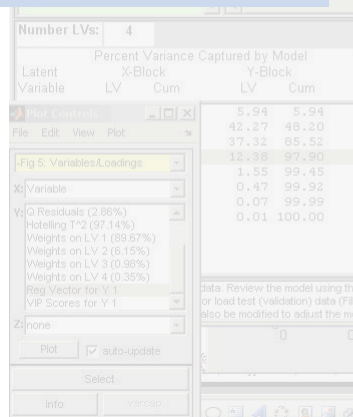
EIGENVECTOR
RESEARCH INCORPORATED

24

Model based selection Automating it a bit

Why select
What methods
Choose method
In practice

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



25

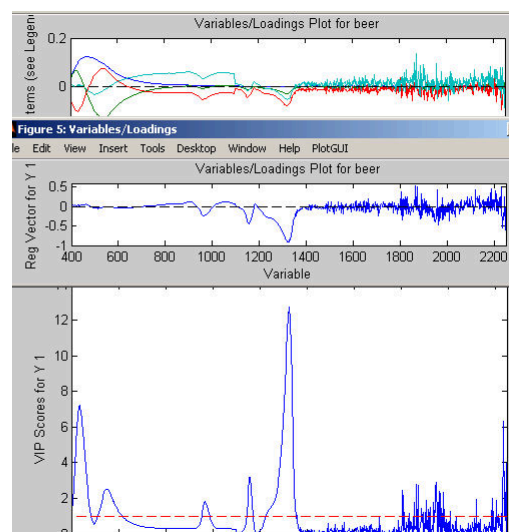
RED

Model based selection Automating it a bit

Why select
What methods
Choose method
In practice

Try yourself on the beer data.

Is a one-shot selection optimal?



26

RED

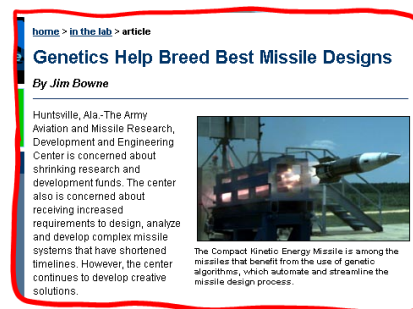
Why select
What methods
Choose method
In practice

“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



27

Why select
What methods
Choose method
In practice

Genetic algorithm

- **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.*



28

Genetic algorithm

Why select
What methods
Choose method
In practice

- **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

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

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



29

Random methods Assumptions

Why select
What methods
Choose method
In practice

- "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



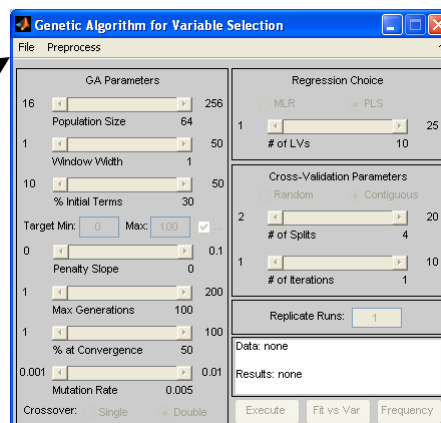
30

Genetic algorithms Exercise

Why select
What methods
Choose method
In practice

In MATLAB
>> load beer
>> genalg

- Load calibration data (beer and extract)
- Execute

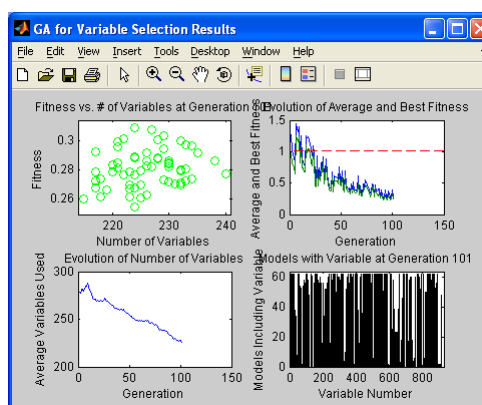


31

Why select
What methods
Choose method
In practice

Genetic algorithms Exercise

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

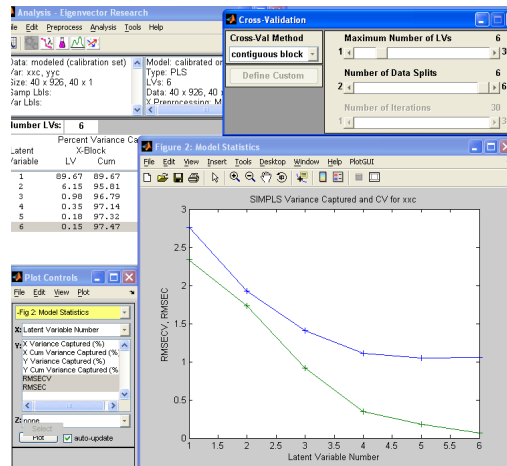


32

Genetic algorithms Exercise

Why select
What methods
Choose method
In practice

- **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



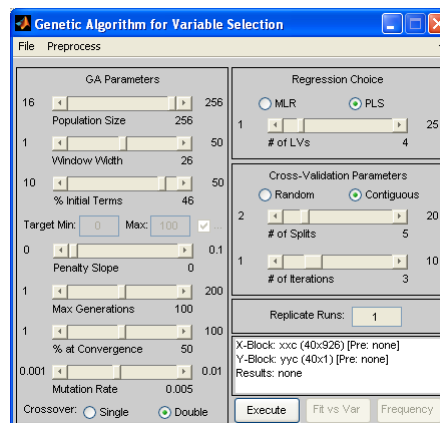
EIGENVECTOR
RESEARCH INCORPORATED

33

Genetic algorithms Exercise

Why select
What methods
Choose method
In practice

- **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

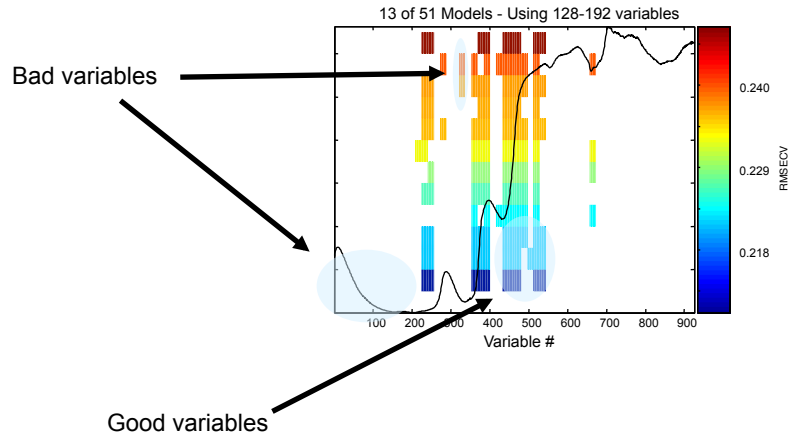


EIGENVECTOR
RESEARCH INCORPORATED

34

Genetic algorithms Exercise

Why select
What methods
Choose method
In practice

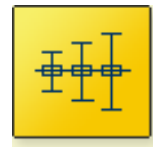


35

Classical methods

Why select
What methods
Choose method
In practice

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



36

Classical methods

Statistical significance tests

Why select
What methods
Choose method
In practice

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

37

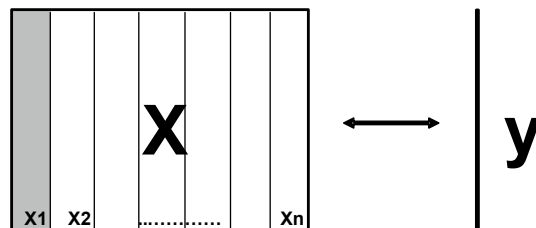


Intermezzo

Why select
What methods
Choose method
In practice

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.

38



Classical methods

Forward selection

Why select
What methods
Choose method
In practice

- **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
 - Works well sometimes!
- **Bad**
 - Disregards interactions to some extent

39



Forward selection

Exercise

Why select
What methods
Choose method
In practice

- **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

Use the algorithm **IPLS** to do variable selection

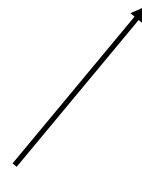
40



Forward selection Exercise

Why select
What methods
Choose method
In practice

```
In MATLAB
>> load beer
>> help ipls
>> width = 20;
>> lv = 4;
>> [use,fit,lvs,int] = ipls(beer,extract,width,lv);
```



Divide into intervals of 20 variables width

41



Forward selection Exercise

Why select
What methods
Choose method
In practice

```
In MATLAB
>> load beer
>> help ipls
>> width =
>> lv = 4;
>> [use,fit,lvs,int] = ipls(beer,extract,width,lv);
```

Maybe more than one interval could be useful.

Use the options in ipls to select more than one interval.



Divide into intervals of 10 variables width

42

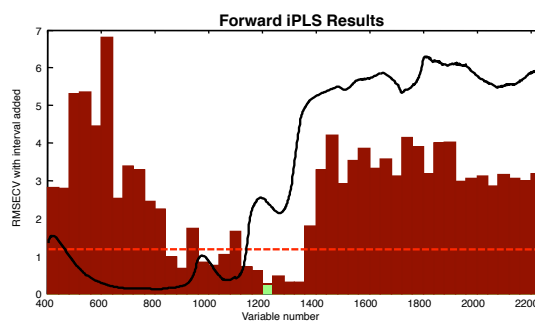


Why select
What methods
Choose method
In practice

Forward selection Exercise

In MATLAB

```
>> myopt = ipls('options');  
>> myopt.numintervals = 3;  
>> myopt.preprocessing = 'meancenter';  
>> [use,fit,lvs,intervals] = ...  
    ipls(beer,extract,width,lv,myopt);
```



43

Why select
What methods
Choose method
In practice

Forward selection Exercise

Number of intervals

Should it perhaps be 2 or 8?

Let ipls check it

In MATLAB

```
>> myopt.numintervals = Inf;  
>> [use,fit,lvs,intervals] = ...  
    ipls(beer,extract,width,lv,myopt);
```

44

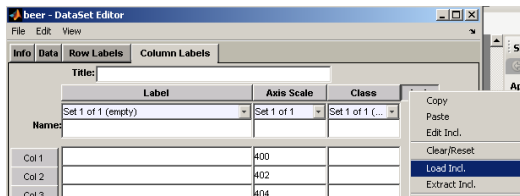
Forward selection Exercise

Why select
What methods
Choose method
In practice

Use the solution and check in PLS

Open beer in data set editor and load use as new include field for variables and save e.g. as beer2. Do a PLS model and see if predictions are fair

```
>> [use, fit, lvs, interval  
    ipls (beer, extract,
```



Classical methods Backward selection

Why select
What methods
Choose method
In practice

- **Principle**
 - Make full model and select the variable contributing the least to the fit
 - Repeat that
- **Good**
 - Takes interactions into account
 - Reasonably fast
- **Bad**
 - "Random" removal for large data sets
 - Often works bad for many irrelevant variables

Forward selection
is usually better



Choosing method

If there are many irrelevant variables

Why select
What methods
Choose method
In practice

- Do not use tests based on statistics of overall model
 - Jack-knife, significance test, etc.
- 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 or forward selection are good choices
- Refine solution using background knowledge and other variable selection tools

47



Choosing method

If refining an already ok model

Why select
What methods
Choose method
In practice

- 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 jack-knife, I-PLS, etc.

48



Choosing method

If you have continuous or smooth data

Why select
What methods
Choose method
In practice

- Do not select individual variables but rather windows of variables
- Use I-PLS, genetic algorithms with window selection etc.

49



Why select
What methods
Choose method
In practice

How to do it in practice

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

50

Alternative – robust methods



How to do it in practice

- **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