

Multivariate statistical analyses and machine learning for metabolomics

Roy Goodacre

School of Chemistry,
The University of Manchester,
Sackville Street, PO Box 88,
Manchester, M60 1QD



MANCHESTER
1824

Roy.Goodacre@manchester.ac.uk
www.biospec.net

Laboratory for Bioanalytical Spectroscopy - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://www.biospec.net/

Getting Started Latest Headlines

MANCHESTER
1824
The University of Manchester

Laboratory for Bioanalytical Spectroscopy
School of Chemistry, University of Manchester

[Home](#) | [Projects](#) | [News](#) | [Highlights](#) | [Members](#) | [Publications](#) | [Facilities](#) | [Research & Teaching Support](#) | [Links](#)

Welcome to our Home Page.

Research: The research theme in our group is predominantly directed towards developing metabolomic and proteomic technologies for the rapid accurate characterisation of biological systems. This is achieved via a tandem analysis using analytical instrumentation that are used to produce rapid physiological metabolome, proteome or 'holistic' whole-organism (phenotypic) fingerprints of bacteria, fungi, human and animal body fluids and plant materials. In order to analyse these high dimensional multivariate data we have been very active in the development of novel chemometric and machine learning techniques. In cognate projects, we are also developing molecular methods for the phylogenetic analysis of microorganisms.

In October 2005 we will move into the Manchester Interdisciplinary Biocentre (www.mib.ac.uk) a £35M building designed to bring together physical scientists, engineers, computer scientists and mathematicians to attack and solve complex, multidisciplinary biological problems.

Postdocs: Dr Catherine Winder, Dr Robert Cornell, Dr Roger Jarvis, Dr Mohammad Afzaal, A.N. Other x 5,

With Collabs: Dr Seetharaman Vaidyanathan, Dr John Fletcher

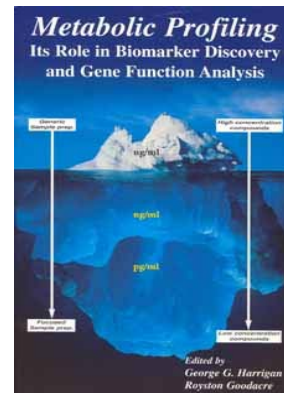
Research Technicians: Jo Ellis, Steffi Schuler + A.N. Other

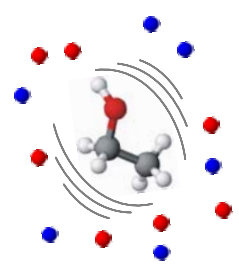
Research Students (PhD): Felicity Currie, Katherine Hollywood, Robben Jaber, Nicoletta Nicolaou, Soyab Patel, Ketan Patel, Emma Wharfe, Will Allwood, Will Cheung, Robert Coe.

From Oct 1: Nicola Wood, Dong Hyun Kim.

Research Students (MSc): Nicola Wood, Sadia Rabbini, Martin Coleston, Graham Mullard.

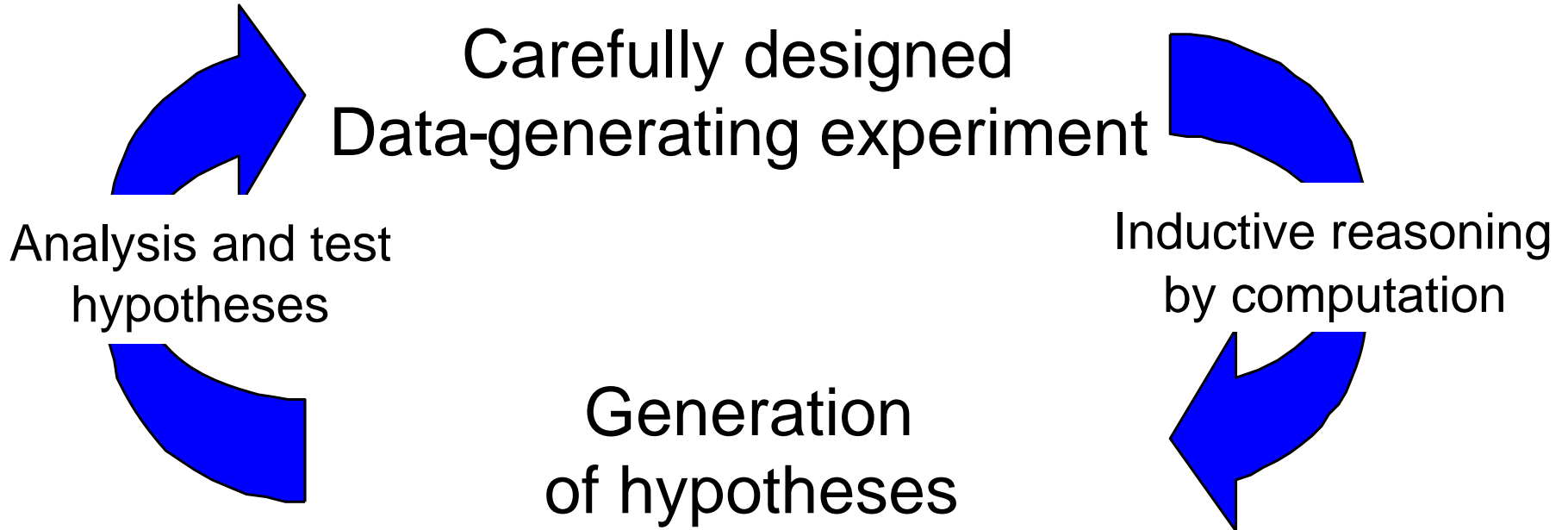
£s, €s, \$s: BBSRC, EPSRC, RSC, EU F6, DEFRA, NERC, ORS



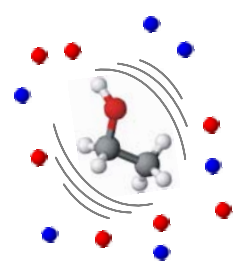


‘Top down’ metabolomics

Strategy: Inductive approach to knowledge discovery via holism



“Hiring a statistician after the data have been collected is like hiring a physician when the patient is in the morgue. He might be able to tell you what went wrong, but is unlikely to be able to fix it”



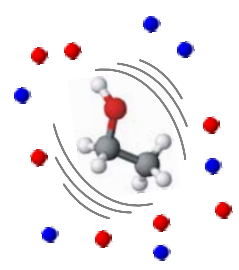
Data floods



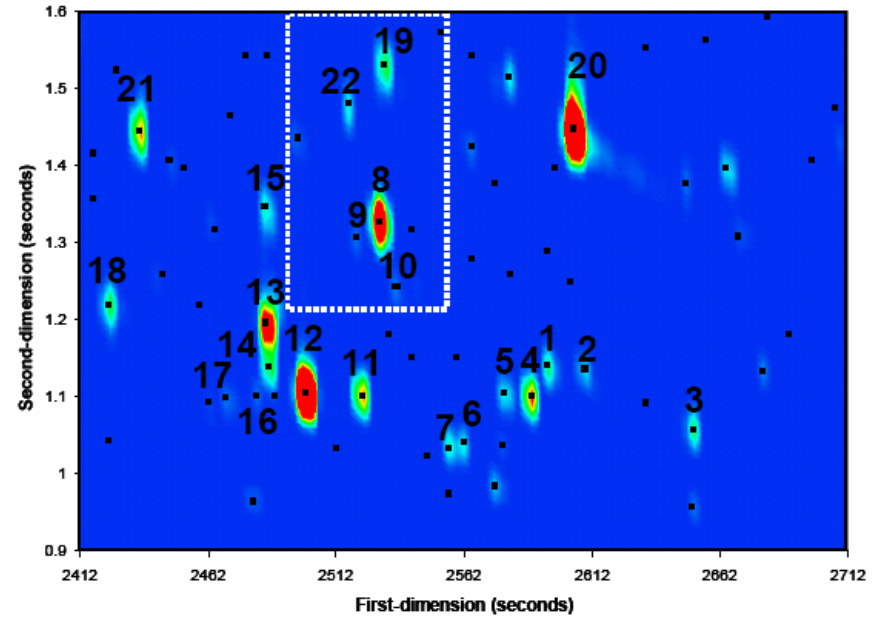
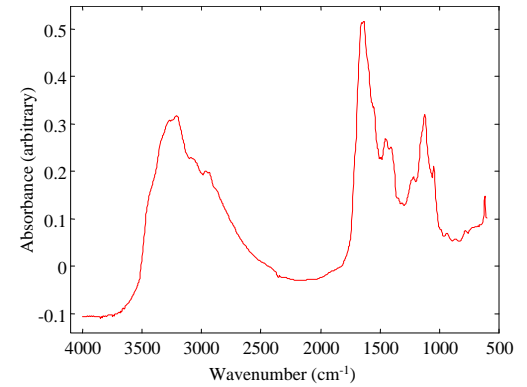
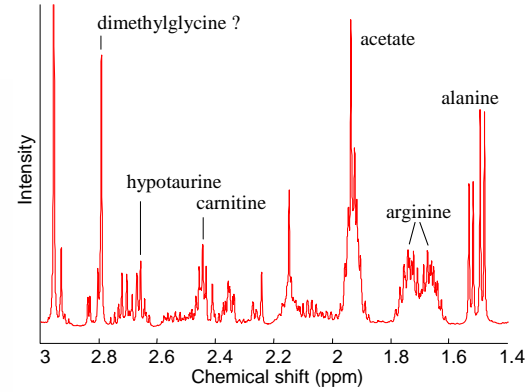
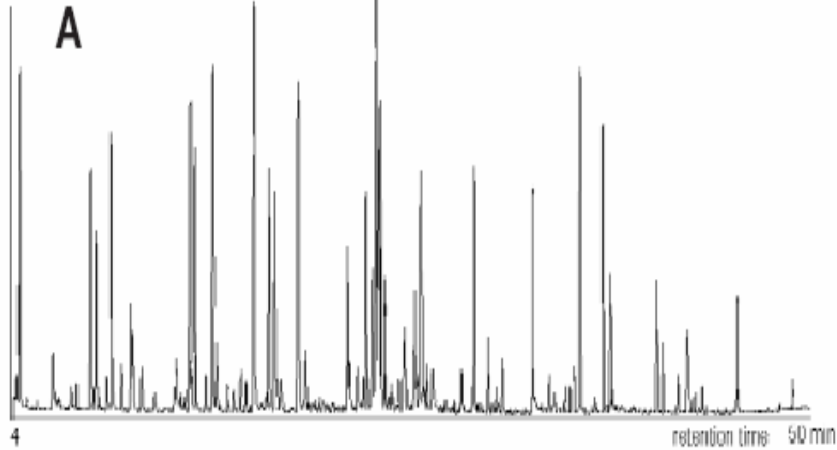
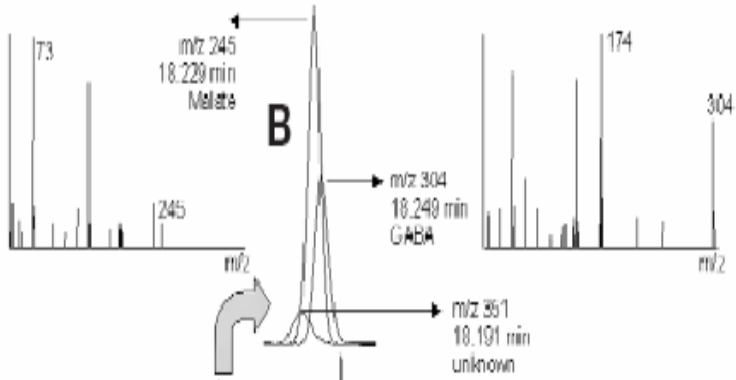
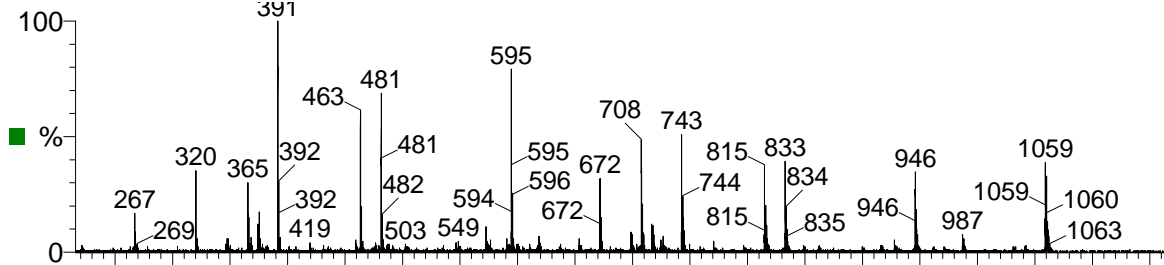
“Data does not equal information;
information does not equal knowledge;
and, most importantly of all, knowledge
does not equal wisdom.

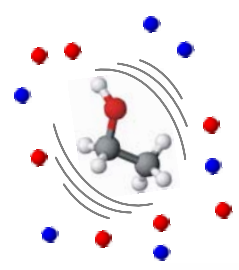
We have oceans of data, rivers of
information, small puddles of knowledge,
and the odd drop of wisdom”

Henry Nix, 1990

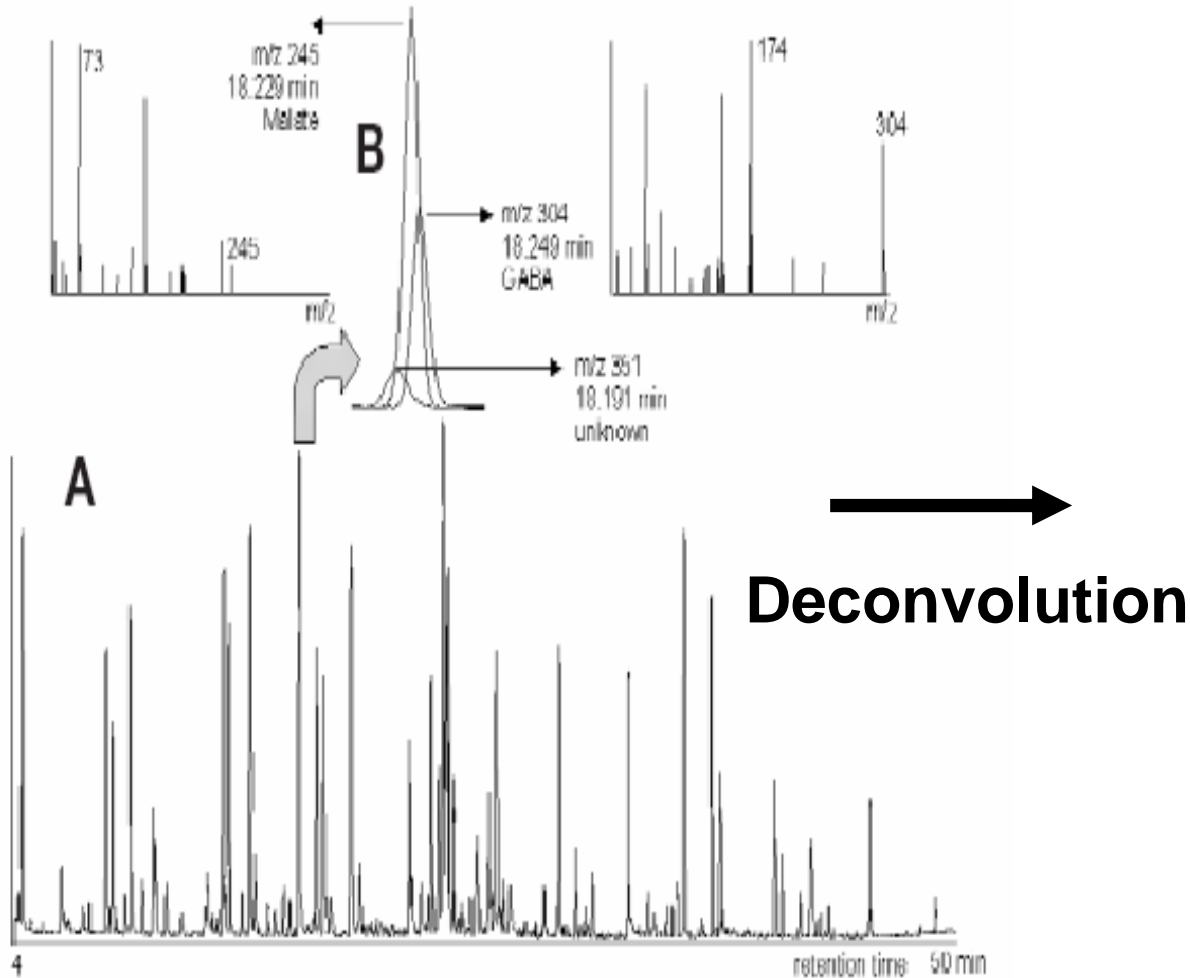


Data outputs

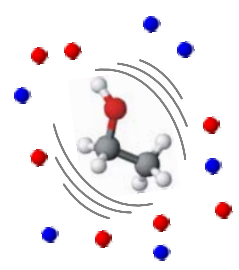




Typical metabolic profile



Metabolite	Conc
Glucose	0.1
Indole	0.001
Tryptophan	1.2
Ethanolamine	0.7
...	...
Metabolite #88	0.9
Metabolite #167	0.05



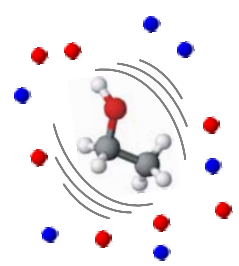
Data handling

Objects going down in different rows	X-var 1 Metabolite or peak 1	X-var 2 Metabolite or peak 2	X-var 3 Metabolite or peak 3
Sample 1			
Sample 2...			

Metabolite	Conc
Glucose	0.1
Indole	0.001
Tryptophan	1.2
Ethanolamine	0.7
...	...
Metabolite #88	0.9
Metabolite #167	0.05



Input data



Data handling

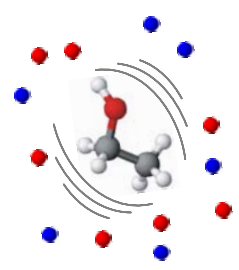
Objects going down in different rows		X-var 1 Metabolite or peak 1	X-var 2 Metabolite or peak 2	X-var 3 Metabolite or peak 3	Y-var 1 Lots of Metadata	Y-var 2 Diseased or Healthy (Levels)
Sample 1					Species Age M/F BMI	0 (control)
Sample 2...					sampling processing etc, etc...	1 (diseased)



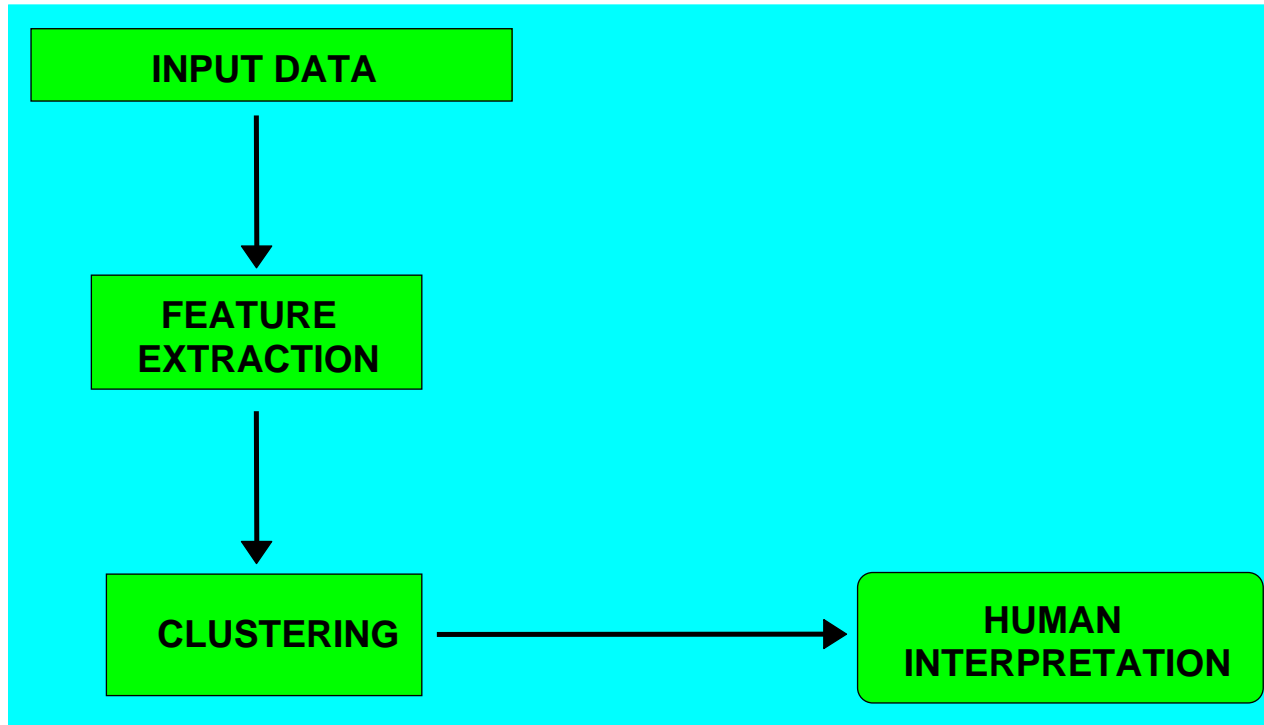
Input data



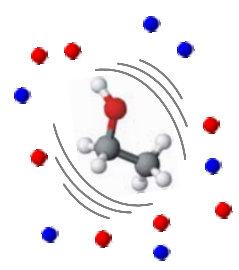
Output data



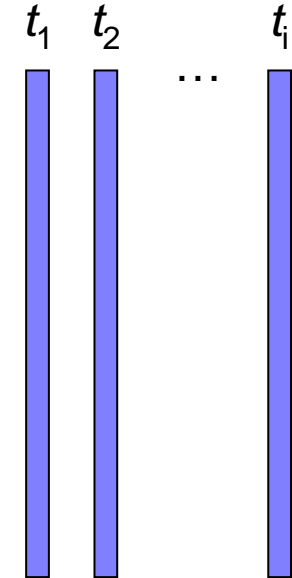
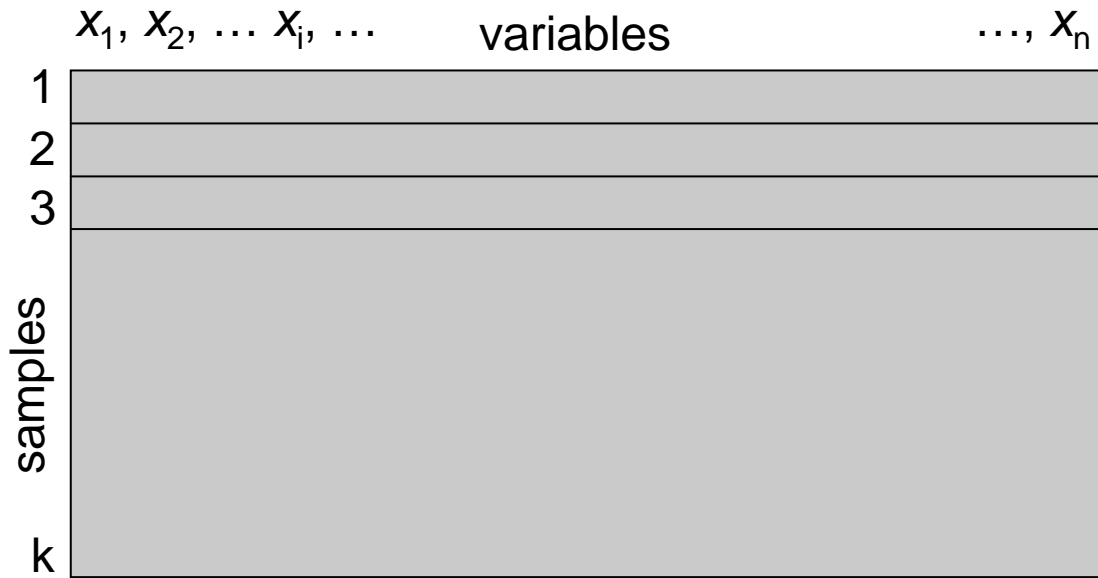
Unsupervised learning



- System is shown a set of inputs (spectra) and then left to cluster the spectra into groups
- This optimization procedure is usually simplification or dimensionality reduction.



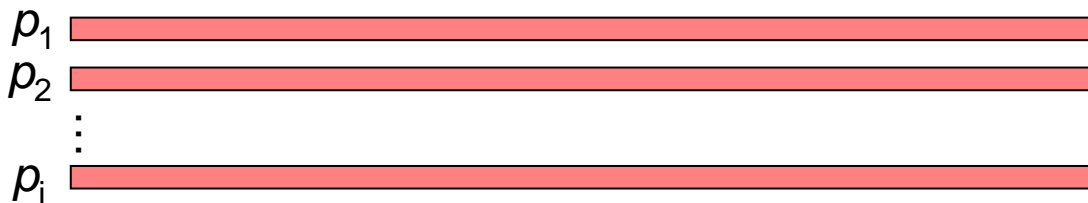
Projection of the data



scores (t)
summarise
variation in
samples

uncorrelated
orthogonal

variance: $t_1 \geq t_2 \geq \dots \geq t_i$

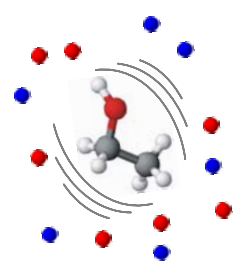


loadings (p)

summarise variation in variables

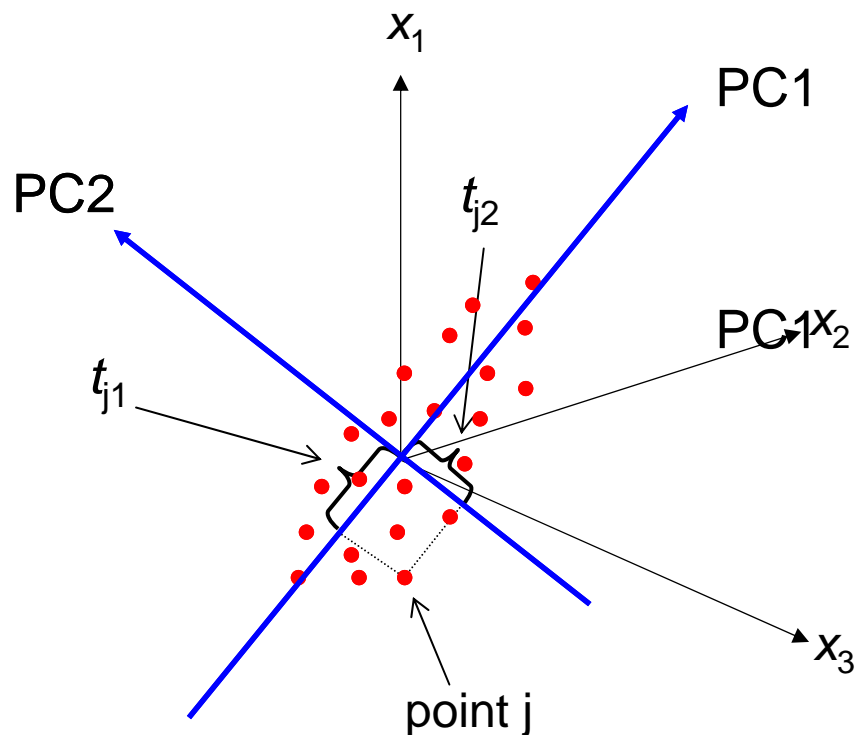
scores = loadings × data

$$t_1 = p_1x_1 + p_2x_2 + \dots + p_ix_i + \dots + p_nx_n$$



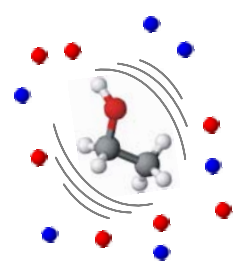
PCA

- Finds structure in data
- Rotate to uncover *maximum* correlations with respect to *natural* variation

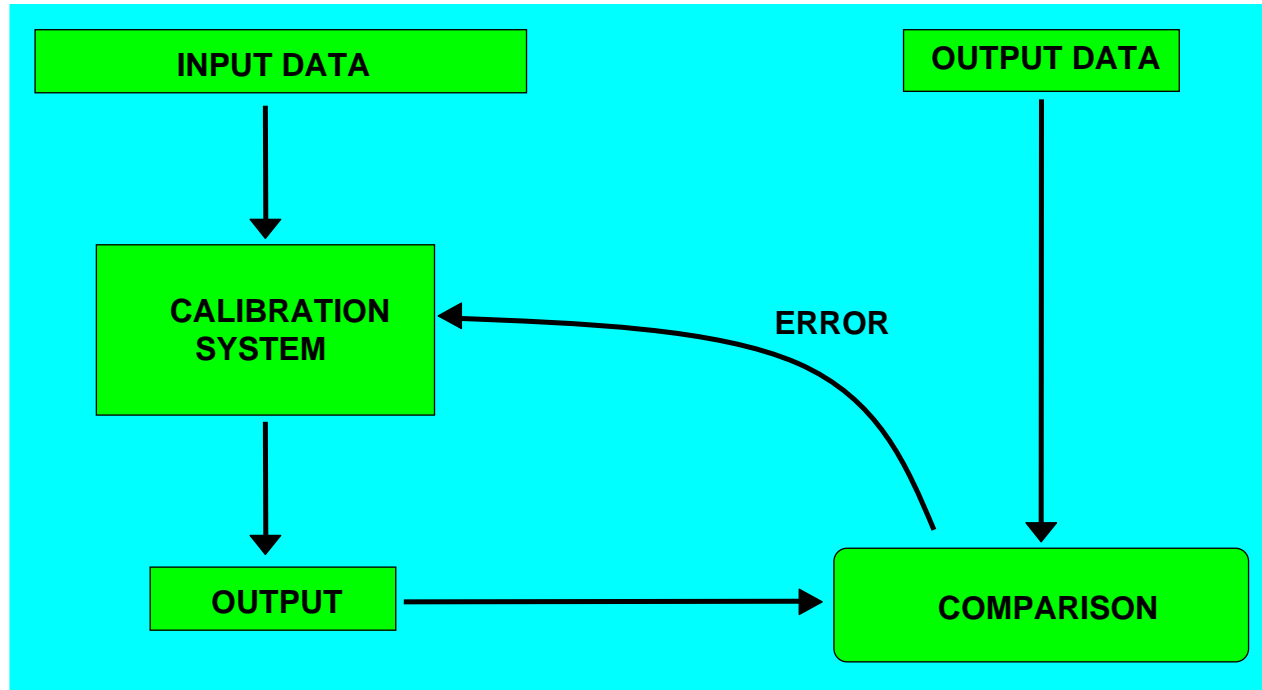


PC1 = 1st principal component
describes largest variance
goes through variable origin space
 t_{j1} = score for point j = distance
from projection of that point
onto PC1 from origin

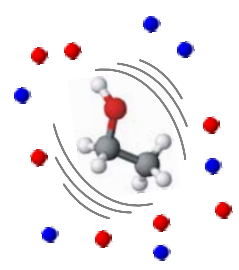
PC2 = 2nd principal component
describes 2nd largest variance
goes through variable origin space
 t_{j2} = score for point j



Supervised learning

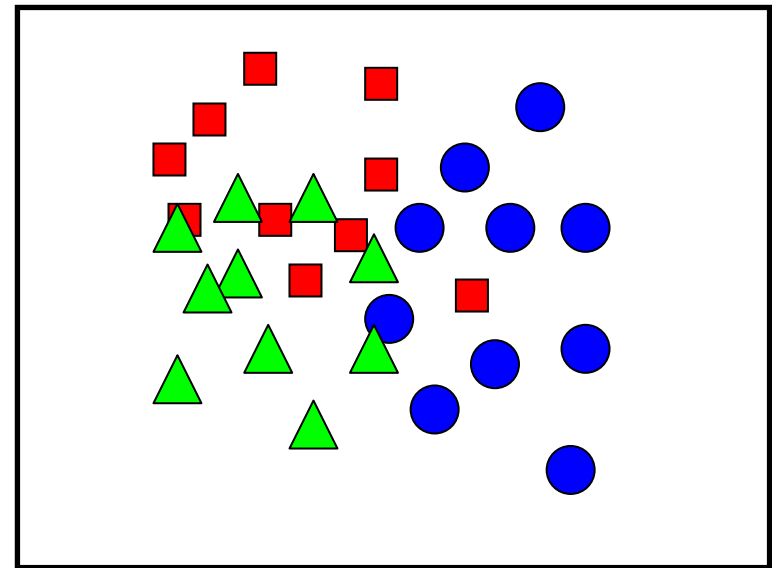
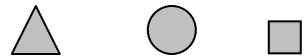


- The goal is to find a mathematical model that will correctly associate the inputs with the targets
- Usually achieved by minimising the error between the target and the model's response (output).

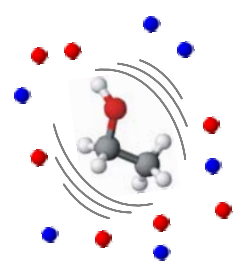


Discriminant function analysis (aka, canonical variates analysis)

- Uses uncorrelated inputs
a priori information
- Projection based on:
 - Minimises within group variance
 - Maximises between group variance
- Test by projection of 'unknown' samples

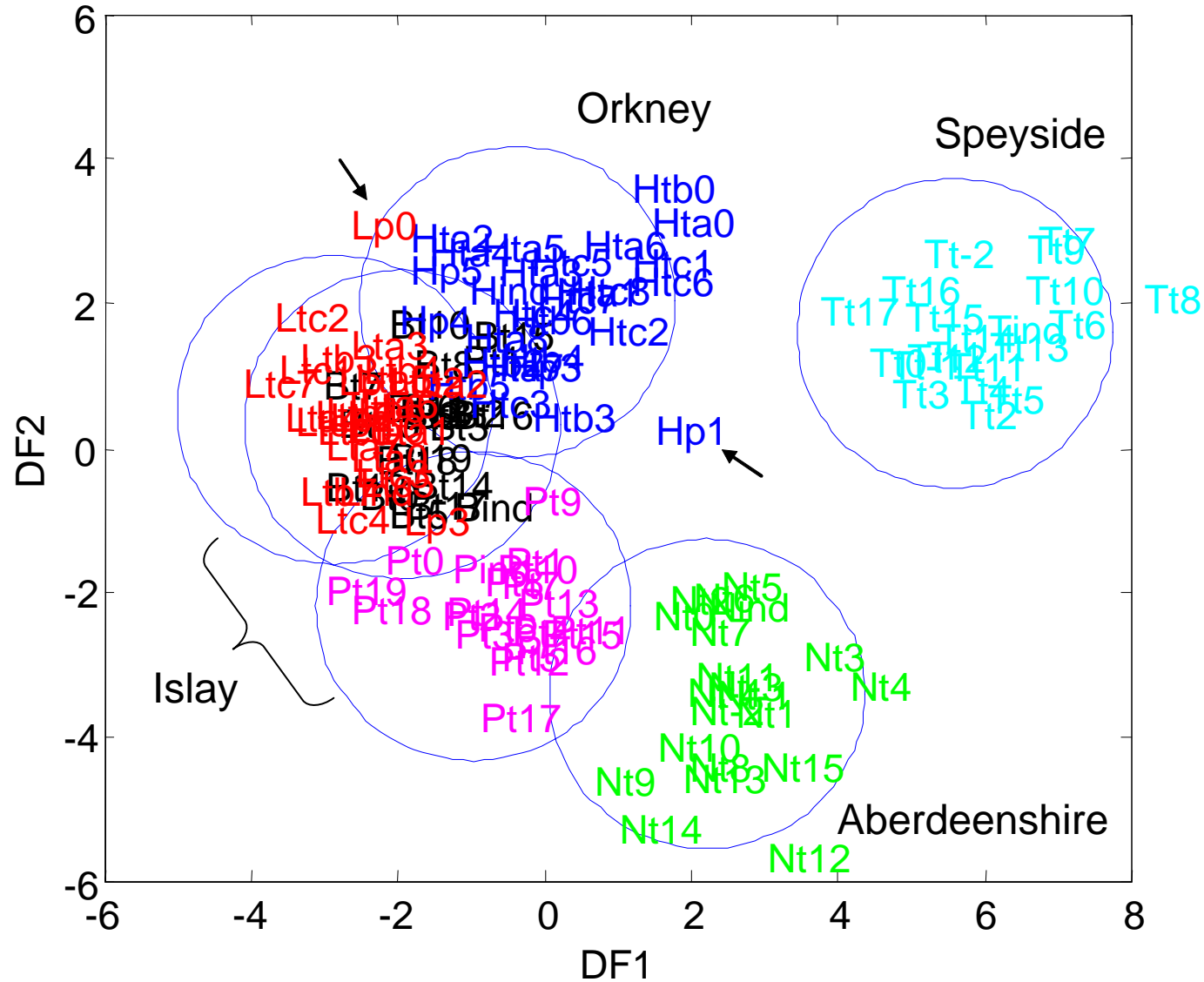


- Statistical significance:
 χ^2 confidence limits



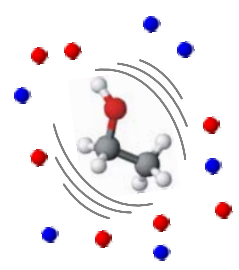
Peat

- 6 groups
- Circles = 95% χ^2 confidence limits
- Arrows represent outlier samples that were from upper horizon of the peat depth profile



Harrison, B., Ellis, J., Broadhurst, D., Reid, K., Goodacre, R. & Priest, F.G. (2006)

Differentiation of peats used in the preparation of malt for Scotch whisky production using Fourier transform infrared spectroscopy, *Journal of the Institute of Brewing*, submitted.



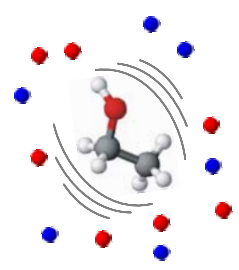
Target encoding for PLS, ANNs, etc...

- Usually binary encoded:

Known bacteria

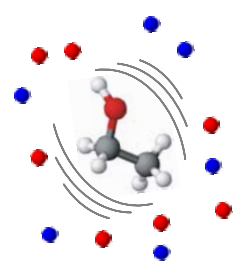
	A	B	C	identity
New isolates				

Easy look up table



Supervised methods are powerful...

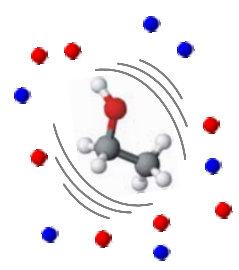
- Learn from experience
- Generalise from previous examples to new ones
- Perform pattern recognition on complex multivariate data.
- Make errors
 - usually because of badly chosen data
 - tanks from trees...



What have I measured that is important

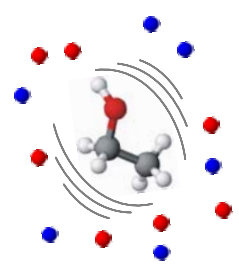
Consider healthy vs. diseased

- 2 class problem
- Collect metabolite data from both classes:
 - Encompass biological variation
 - ❖ Even distribution
 - Try to keep number of samples the same
 - ❖ Statistics are easier / more rigorous



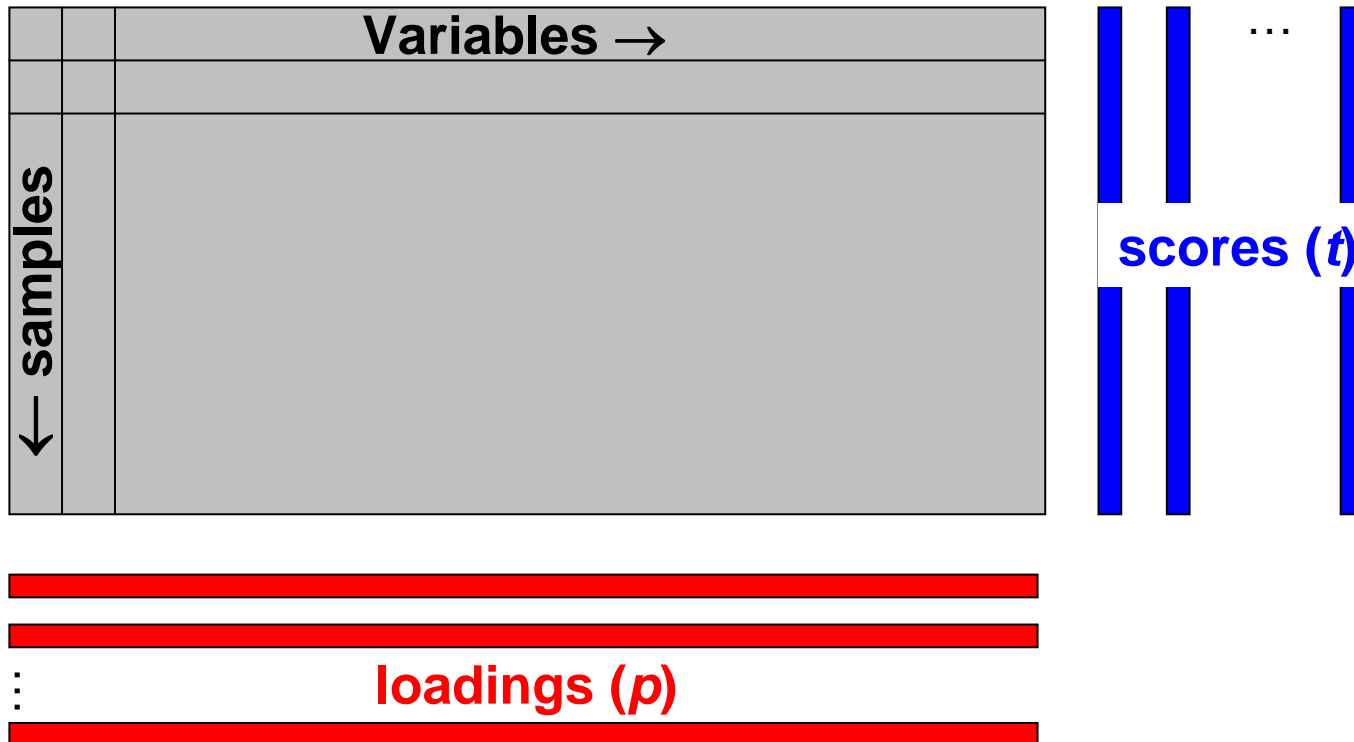
Easy strategies first

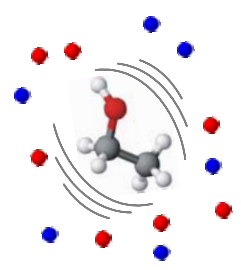
- Stare and compare!
 - Always wise to actually look at the data...
- Difference spectra
 - $\text{Avg}(\text{diseased}) - \text{Avg}(\text{healthy})$
- Univariate analyses
 - Analysis of variance (ANOVA)
 - T-test, etc



Loadings plots

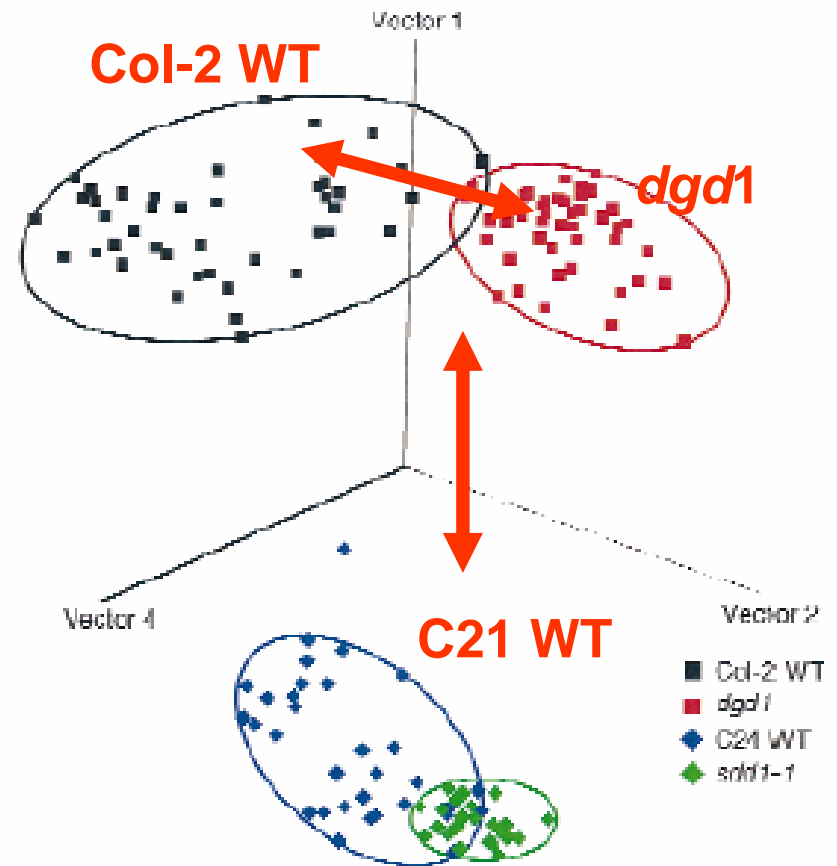
- Use multivariate analyses and inspect loadings



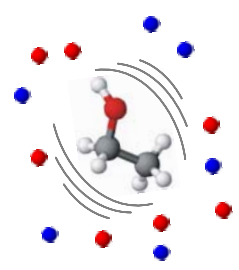


Types of loadings plots

- Unsupervised
 - PCA
 - SOFM
- Supervised: discriminatory
 - DFA (CVA)
 - PLS-DA
- Supervised: regression
 - PLS
- Supervised: tree-based
 - CART



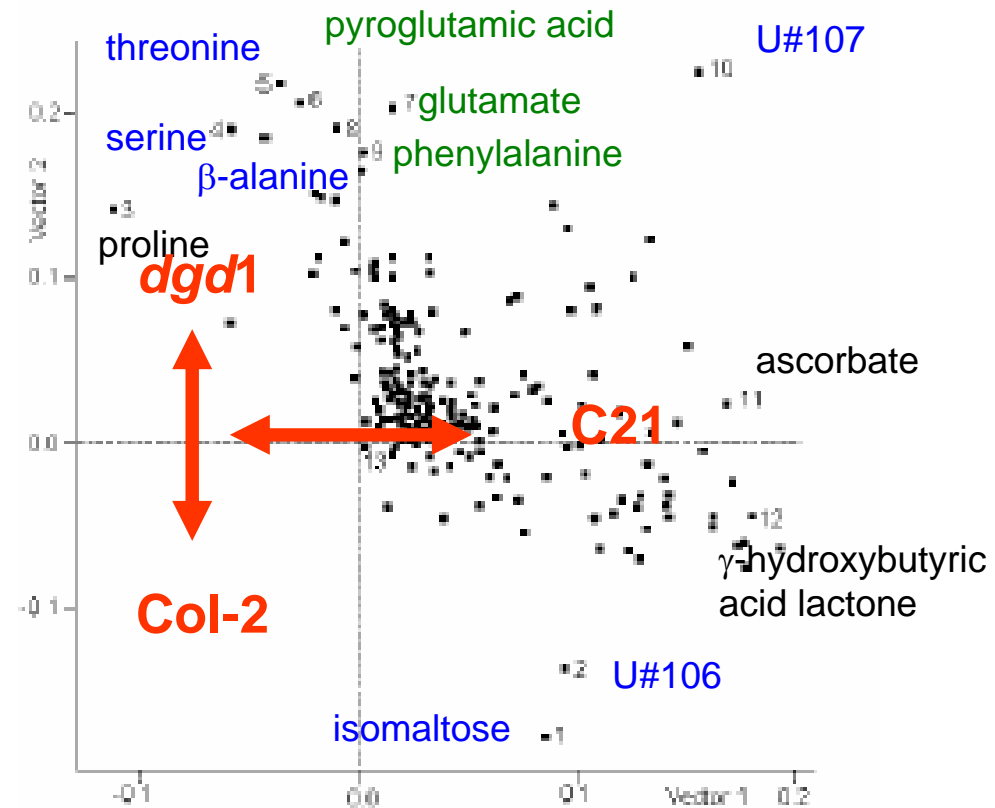
Fiehn *et al.* (2000) Metabolite profiling for plant functional genomics.
Nature Biotechnology **18**, 1157-1161.



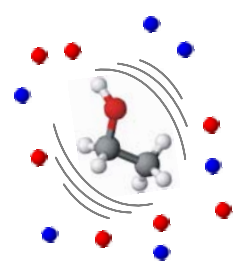
Types of loadings plots

Univariate
 $p < 0.01$ *t*-test
 p not signif.

- Unsupervised
 - PCA
 - SOFM
- Supervised: discriminatory
 - DFA (CVA)
 - PLS-DA
- Supervised: regression
 - PLS
- Supervised: tree-based
 - CART

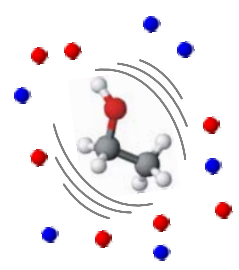


But not always that easy...



Modest 250 metabolites

- 2 class problem
 - Healthy vs. Diseased
- To use or not to use?
- Yes / No, 250 times = 2^{250} or 1.8×10^{75}
- PC does say 10×10^6 orderings every second it would still take $> 3 \times 10^{62}$ years

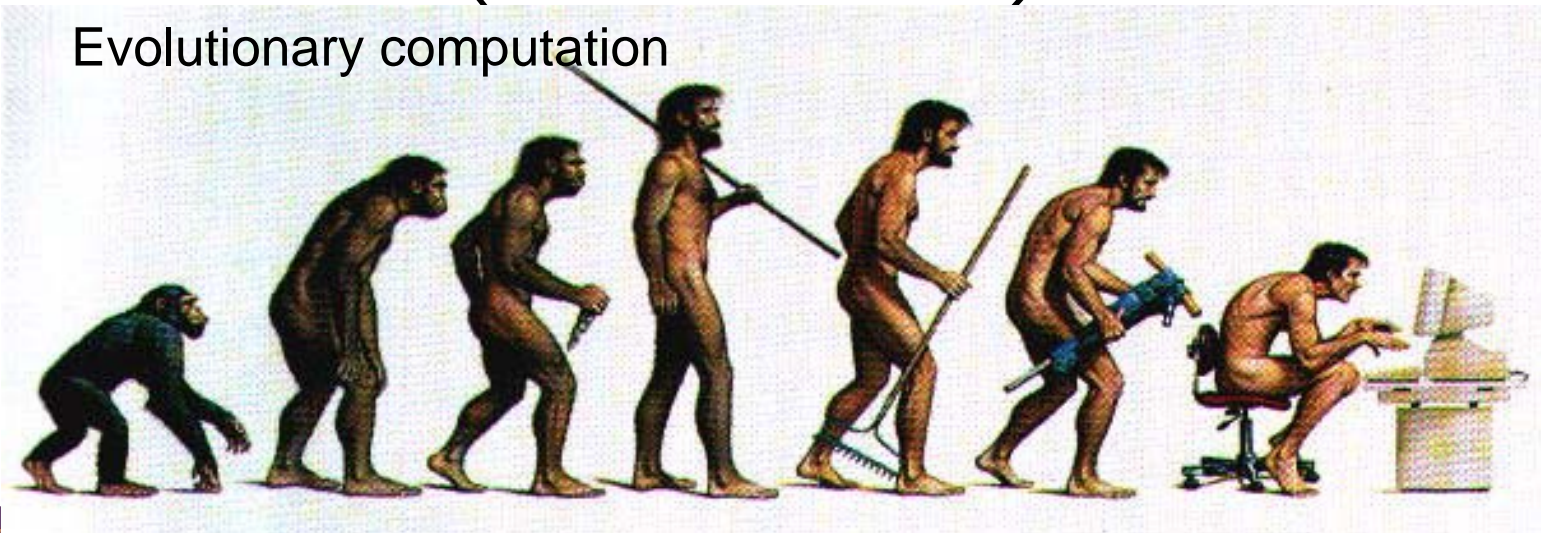


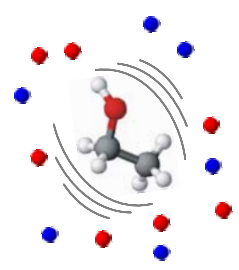
Houston we have a problem...

Complex NP problem $\xrightarrow{\text{No algorithm}}$ Global optimal solution

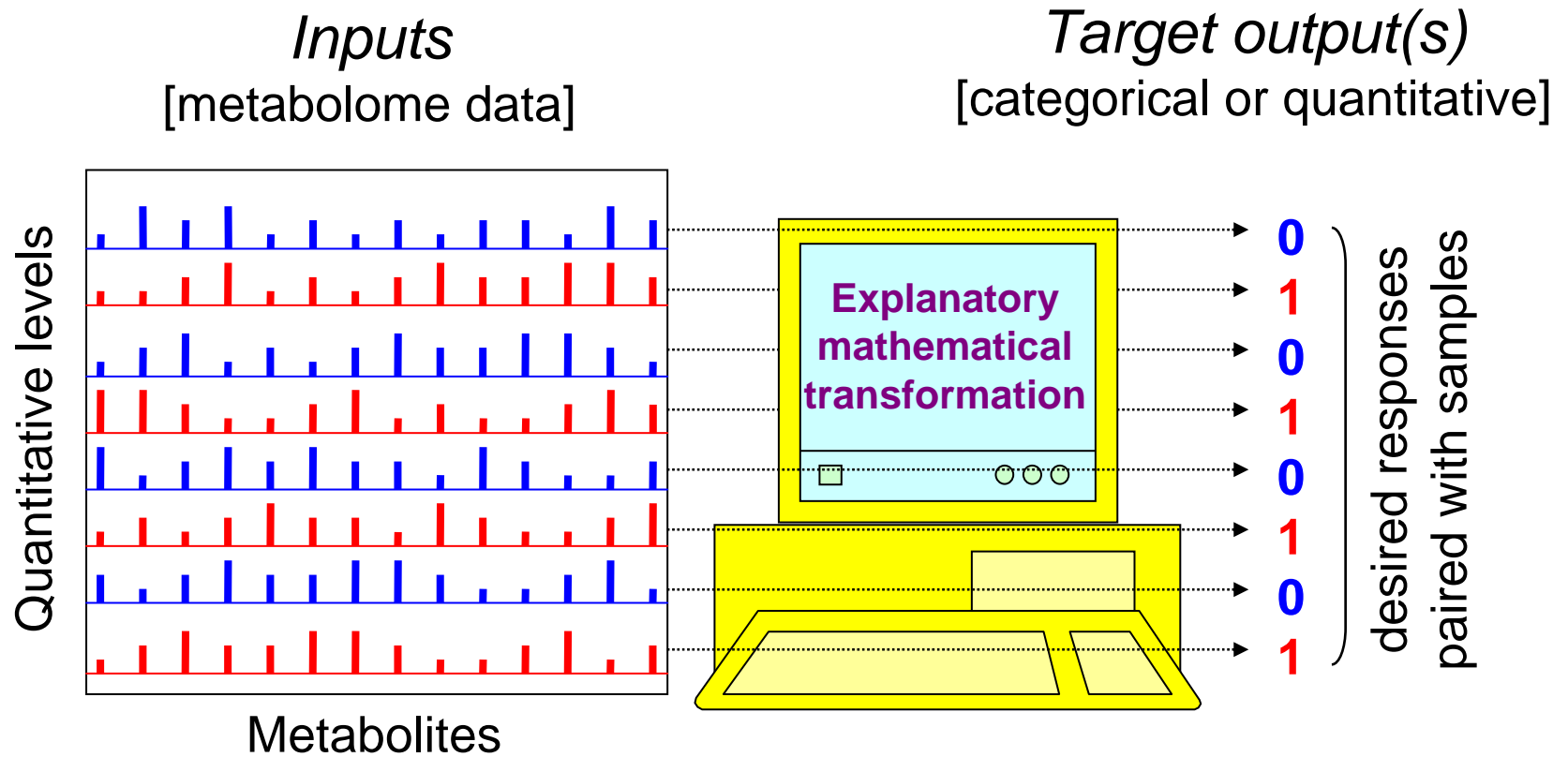
Complex NP problem $\dashrightarrow \times \dashrightarrow$ Good solution

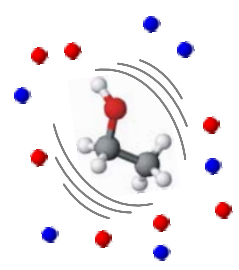
Evolutionary computation





Explanatory machine learning





Genetic algorithm

- Chromosome with n genes selects which of n inputs to use in DA, MLR, PLS or ANN

inputs: 1 2 ... i ... $n-1$ n

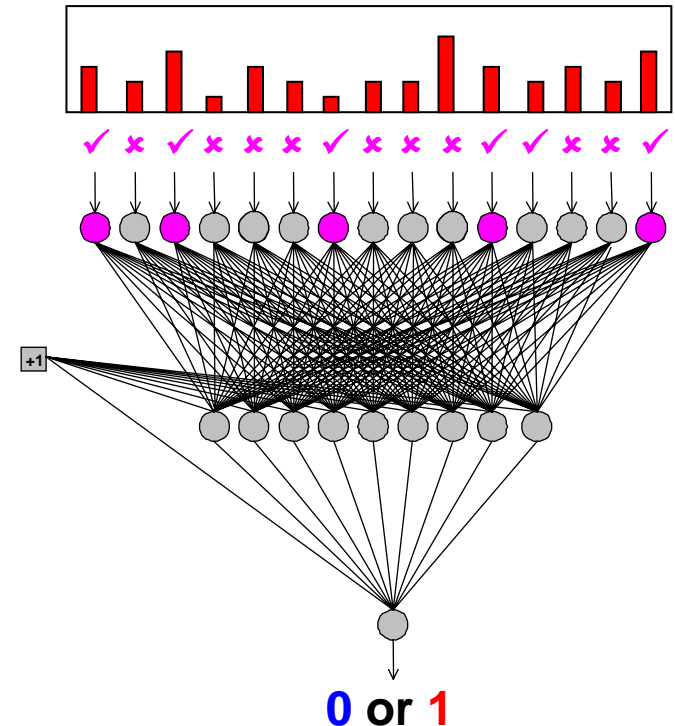
genes:

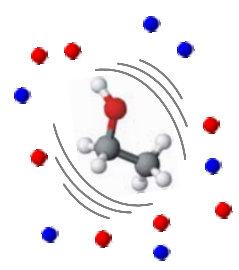
1	1	0	1	0	1	0
---	---	---	---	---	---	---



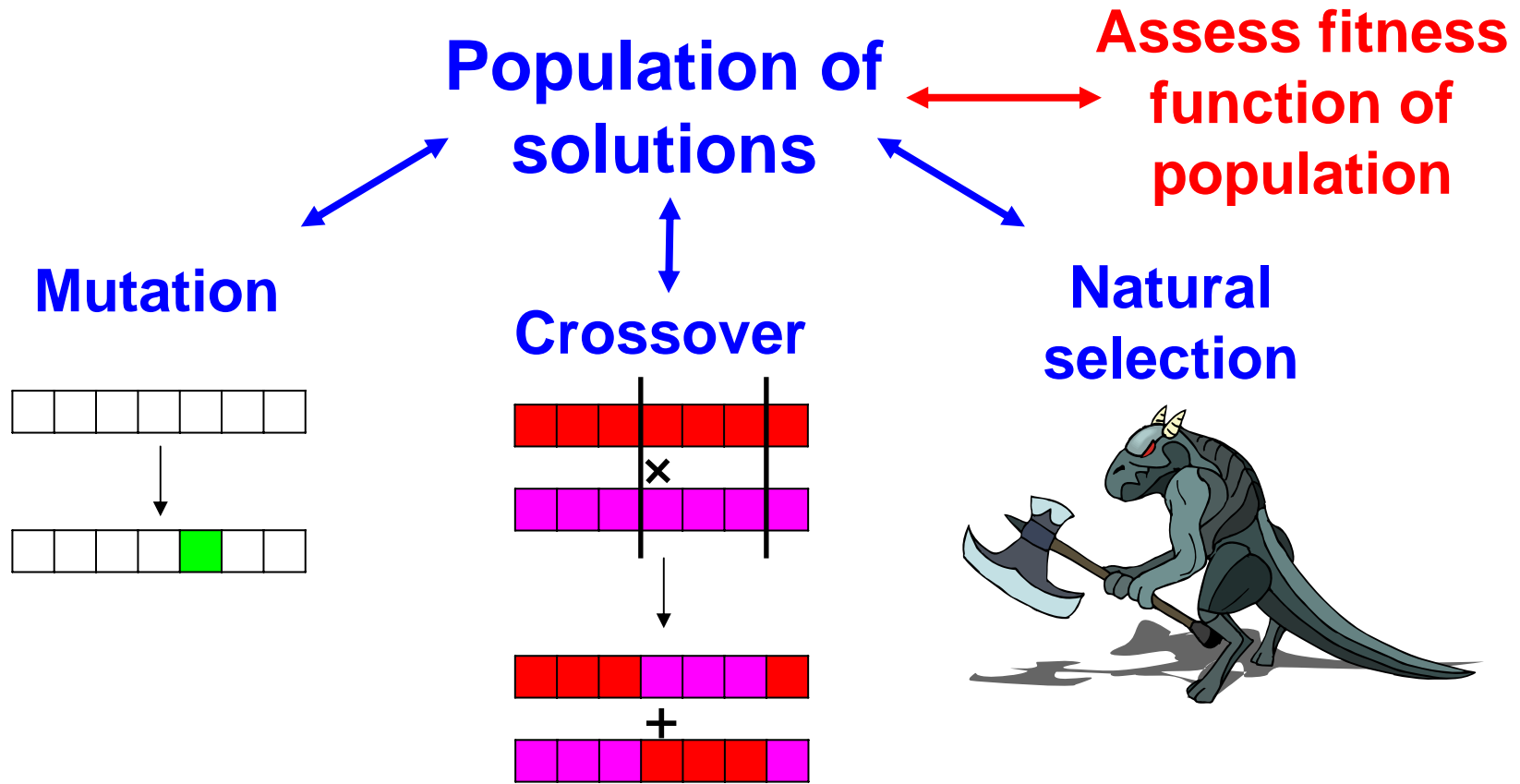
filter:

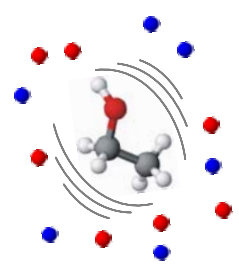
✓	✓	✗	✓	✗	✓	✗
---	---	---	---	---	---	---



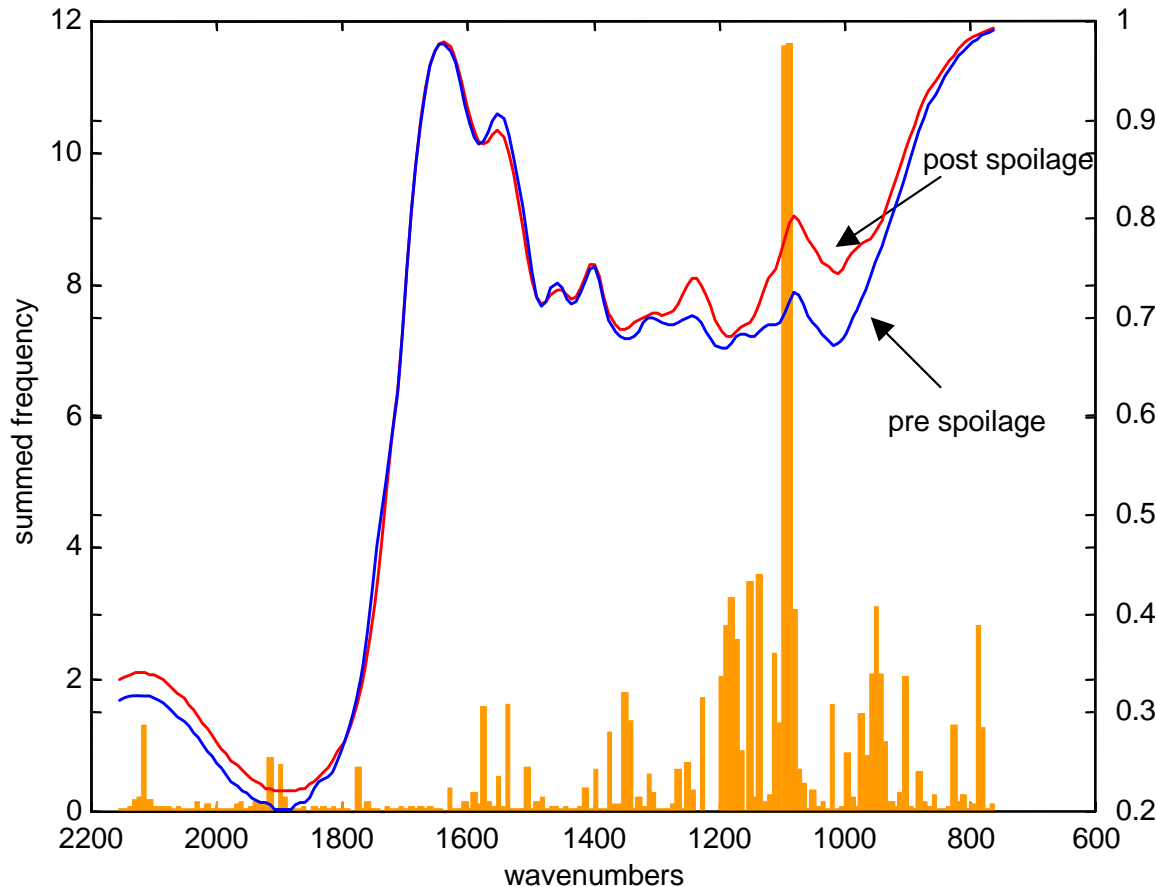


GA: *in silico* evolution





Frequency of input selections



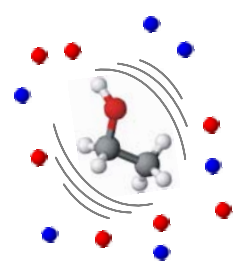
Spoilage at 10^7 =

$$\frac{1096 \text{ cm}^{-1}}{1683 \text{ cm}^{-1}} =$$

$$\frac{\text{Free amines (C-N stretch)}}{\text{Amide I (C=O vibration)}} =$$

Proteolysis of meat by
microbes detected by FT-IR

Ellis, D.I., Broadhurst, D., Kell, D.B., Rowland, J.J. & Goodacre, R. (2002) Rapid and quantitative detection of the microbial spoilage of meat using FT-IR spectroscopy and machine learning. *Applied and Environmental Microbiology* **68**, 2822-2828.



Overall conclusions on hypothesis generation

- Many approaches can be used ranging from
 - Simple → stare and compare, ANOVA etc
 - Multivariate → loadings: PCA, DFA, PLS
 - More complex → evolutionary computation
- Need to design experiment carefully
- Hypotheses might not always be correct
 - But a good place to start when initial knowledge is non-existent.



PyChem v2.0.0 Beta

(<http://pychem.org.uk>)

- Standalone WinXP **graphical** application for MVA
- Preprocessing algorithms incorporated
- Many standard chemometrics algorithms
 - PCA, HCA, DFA, PLSR, PLS-DA.
 - Feature selection using a variety of genetic algorithm tools
- Coded in Python (<http://python.org/>) the software is both **FREE** and **OPEN SOURCE**. Source code can be downloaded at <http://sourceforge.net/projects/pychem/>.
- Can be used to analyse continuous or discrete data such as transcriptomic, metabolomic (vibrational spectra, GC-MS, LC-MS etc...)

Written by Dr Roger Jarvis

Use: Raw data

Processed data

Principal Component scores

Apply principal components

Extract discriminant functions

Apply full cross validation

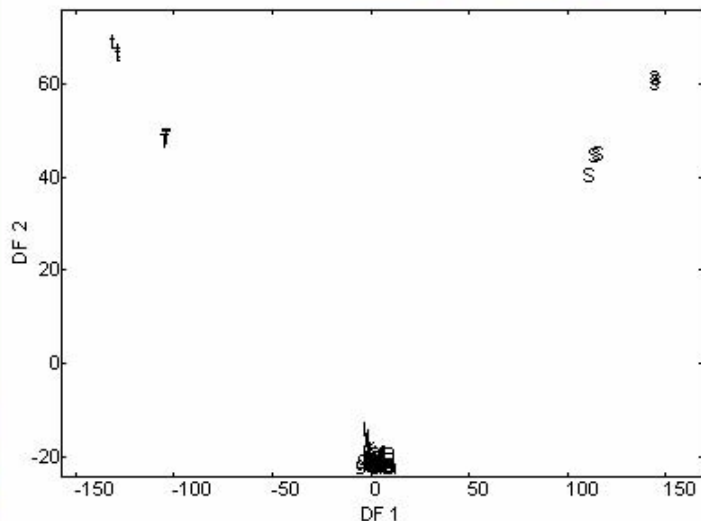
Run DFA

Export

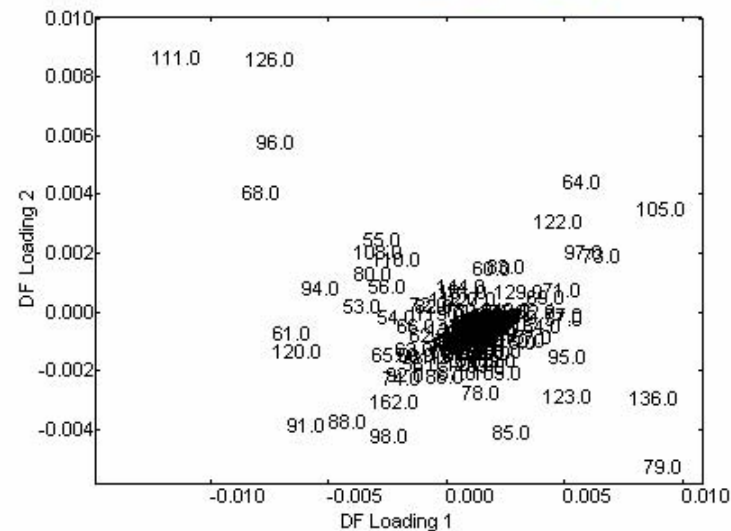
Goto PCA

Goto GA-DFA

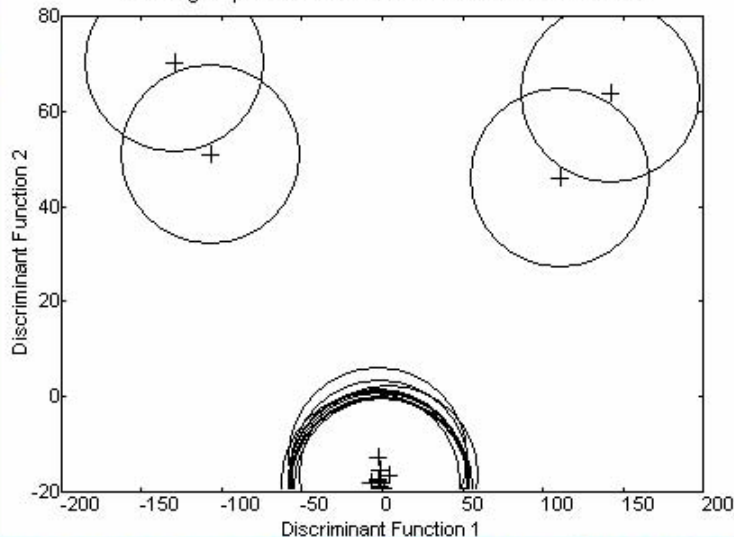
Discriminant Function vs.



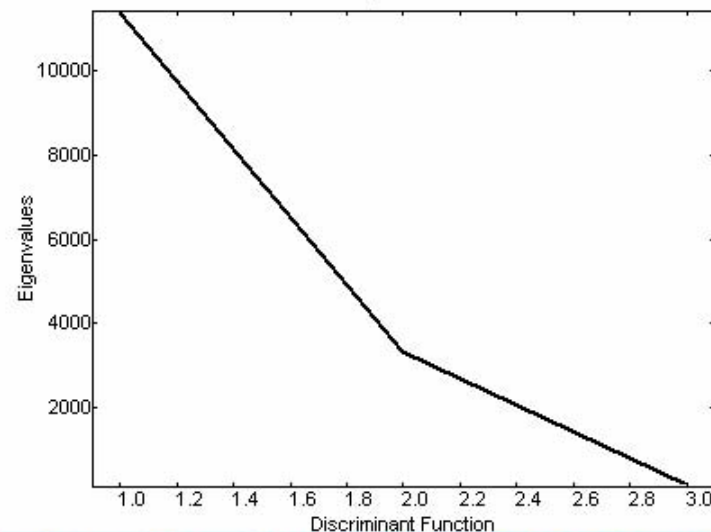
Discriminant Function Loading vs.



Mean group centres with 95% confidence intervals



Eigenvalues



Use: Raw data
 Processed data
 Principal Component scores

Apply 10 principal components

Extract 3 discriminant functions

Apply full cross validation

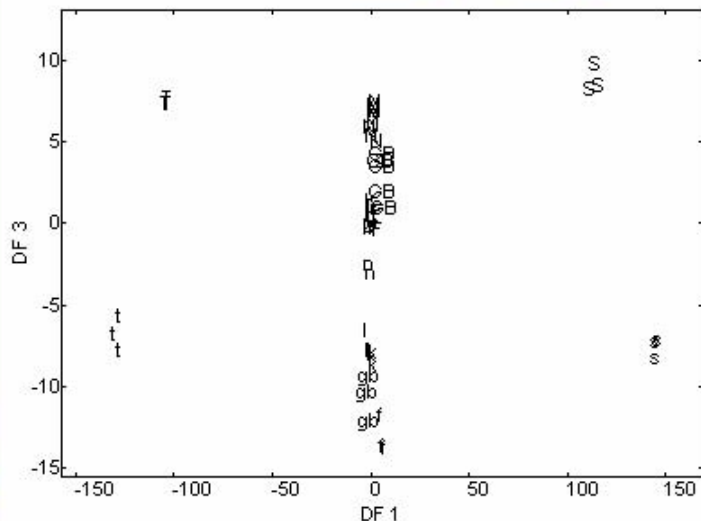
Run DFA

Export

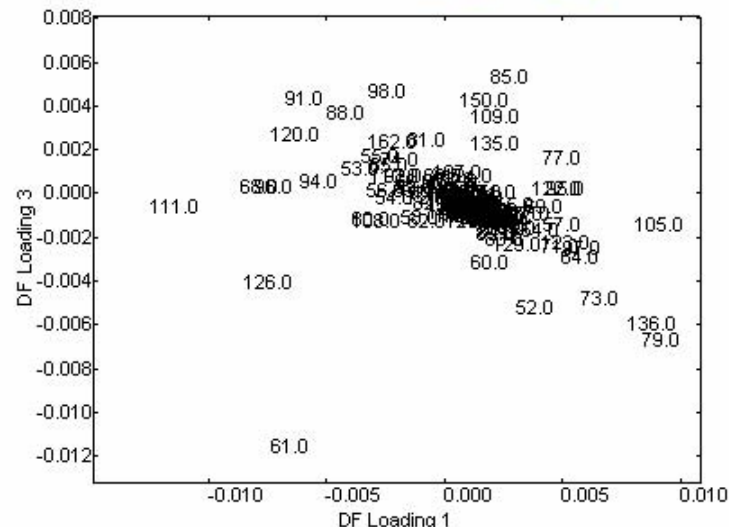
Goto PCA

Goto GA-DFA

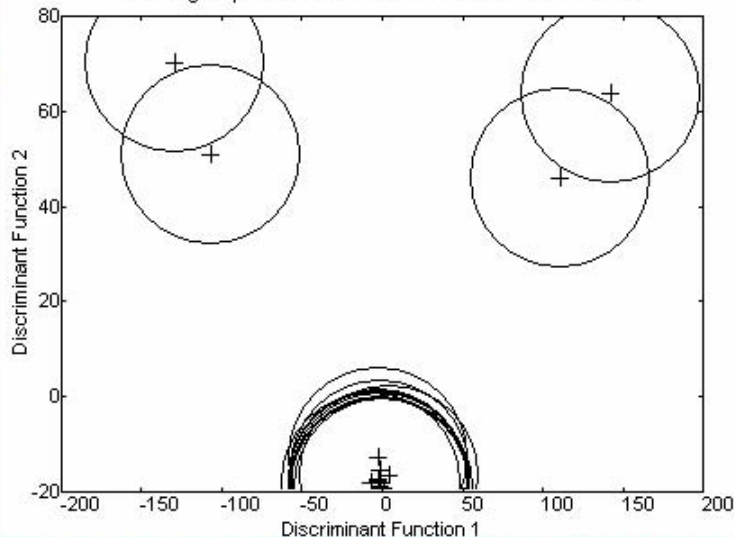
Discriminant Function 1 vs. 3



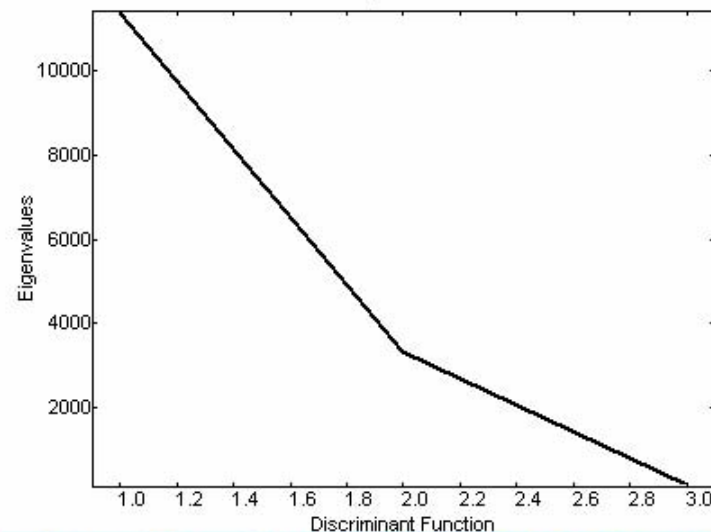
Discriminant Function Loading 1 vs. 3



Mean group centres with 95% confidence intervals



Eigenvalues



GA Parameters

No. vars. from	3	No. vars. to	3
No. inds.	5	No. runs	100
<input checked="" type="checkbox"/> Xover rate	0.8	<input checked="" type="checkbox"/> Mut. rate	0.4
Insert. rate	0.8	Max. DFs	2
Max. gen.	500	<input type="checkbox"/> Rep. until	20

Raw data Processed data

Run GA-DFA

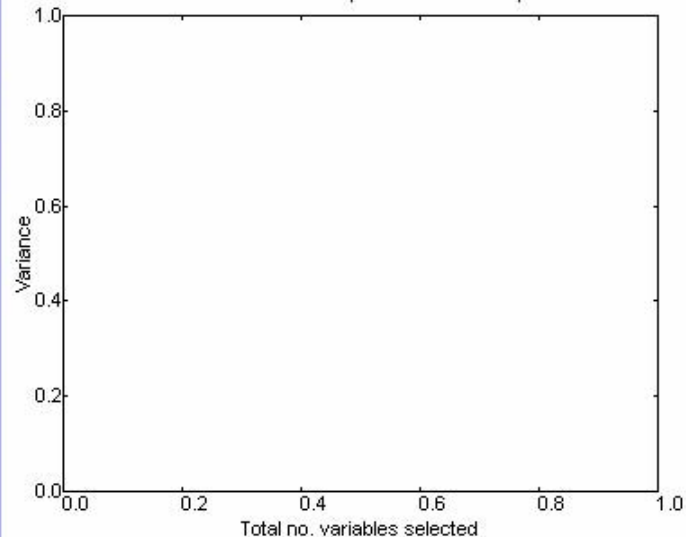
Export Results

GA Results

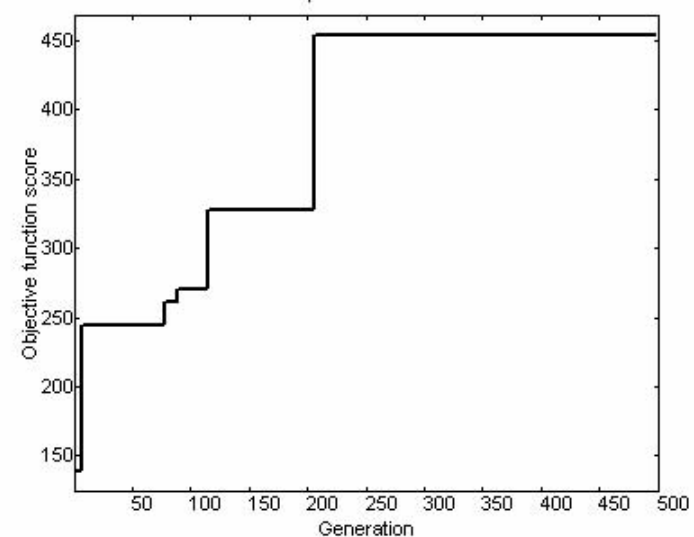
3 variables

#1	[58. 94. 109.]	352.29
#2	[58. 94. 109.]	352.29
#3	[58. 94. 109.]	352.29
#4	[109. 117. 136.]	454.11
#5	[58. 94. 109.]	352.29
#6	[57. 109. 120.]	324.76
#7	[58. 94. 109.]	352.29
#8	[69. 109. 136.]	327.45
#9	[109. 117. 136.]	454.11
#10	[58. 67. 109.]	343.35
#11	[58. 73. 109.]	336.46
#12	[109. 117. 136.]	454.11
#13	[89. 109. 117.]	392.02
#14	[58. 94. 109.]	352.29
#15	[58. 94. 109.]	352.29
#16	[58. 73. 109.]	336.46
#17	[67. 117. 136.]	331.19
#18	[109. 117. 136.]	454.11
#19	[82. 117. 136.]	338.26
#20	[58. 94. 109.]	352.29
#21	[58. 67. 109.]	343.35
#22	[69. 109. 136.]	327.45
#23	[109. 117. 136.]	454.11

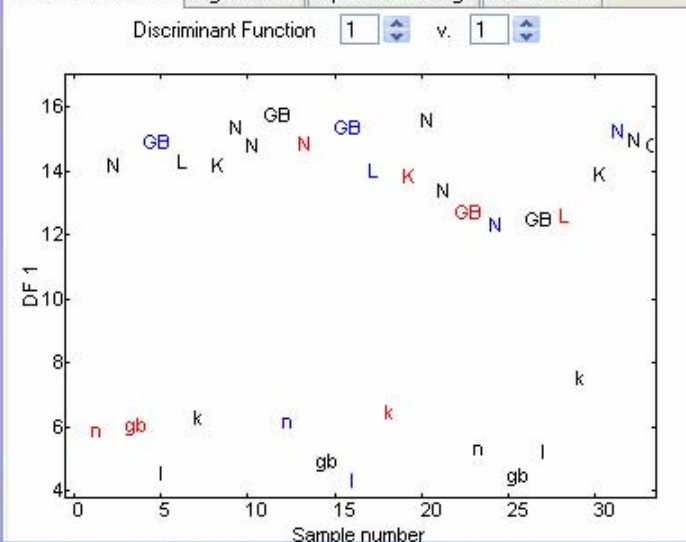
Between/Within Group Variance Comparison



GA Optimisation Curve



2D Ordination Plot



Frequency of Variable Selection

