# Introduction to Metabolomics

## Introduction to chemometrics
### How to analyse many metabolites

UNIVERSITY OF CAMBRIDGE

Department of Plant Sciences

# Introduction to chemometrics
## - Aims

- Align and Bin raw spectral output files

- Merge files for statistical analysis

- Identify many metabolites

- Multivariate data analysis – PCA, Cluster analysis

- ANOVA, false-positive corrections

- Metadata

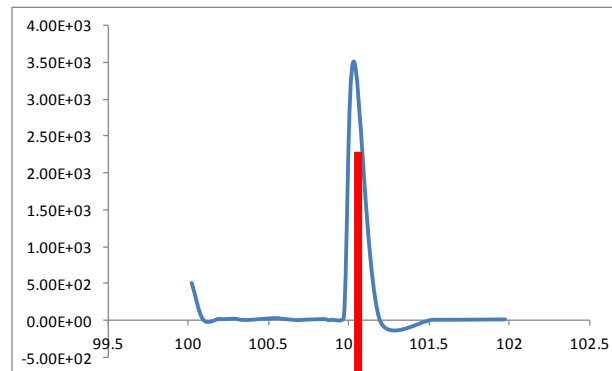# •Align and "Bin" raw spectral output files

## Technical rep 1

| | |
|---|---|
| 100.0181 | 4.27E+02 |
| 100.1587 | 8.33E+00 |
| 100.1869 | 9.78E+00 |
| 100.3746 | 3.13E+01 |
| 100.5719 | 2.41E+01 |
| 100.666 | 4.44E+00 |
| 100.713 | 4.56E+00 |
| 100.7412 | 9.67E+00 |
| 100.7694 | 5.78E+00 |
| 100.793 | 2.44E+00 |
| 100.8212 | 4.22E+00 |
| 100.8589 | 1.20E+01 |
| 100.9624 | 8.59E+01 |
| 101.0284 | 3.39E+03 |
| 101.4859 | 4.78E+00 |
| 101.8025 | 1.13E+01 |
| 101.9065 | 1.04E+01 |
| 101.9728 | 4.67E+00 |

## Technical rep 2

| | |
|---|---|
| 100.0174 | 4.36E+02 |
| 100.1534 | 7.56E+00 |
| 100.1769 | 1.18E+01 |
| 100.2472 | 1.76E+01 |
| 100.3928 | 8.00E+00 |
| 100.435 | 1.81E+01 |
| 100.482 | 6.67E+00 |
| 100.5149 | 5.78E+00 |
| 100.5525 | 1.57E+01 |
| 100.5901 | 4.00E+00 |
| 100.7782 | 2.48E+01 |
| 100.8912 | 2.78E+00 |
| 100.9194 | 1.39E+01 |
| 100.9618 | 7.48E+01 |
| 101.0277 | 3.44E+03 |
| 101.6789 | 1.70E+01 |
| 101.7025 | 9.44E+00 |
| 101.7262 | 7.56E+00 |
| 101.7782 | 1.61E+01 |
| 101.9484 | 3.67E+00 |

## Technical rep 3

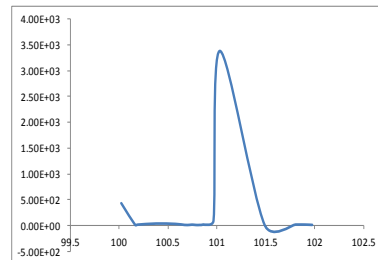| | |
|---|---|
| 100.0179 | 5.06E+02 |
| 100.0882 | 1.09E+01 |
| 100.1867 | 1.81E+01 |
| 100.2148 | 1.38E+01 |
| 100.2805 | 2.18E+01 |
| 100.3322 | 5.11E+00 |
| 100.3979 | 7.89E+00 |
| 100.5294 | 2.89E+01 |
| 100.6093 | 1.44E+01 |
| 100.6799 | 1.78E+00 |
| 100.7175 | 6.67E+00 |
| 100.8445 | 1.56E+01 |
| 100.8728 | 7.78E-01 |
| 100.8963 | 8.67E+00 |
| 100.9669 | 4.27E+01 |
| 101.0282 | 3.51E+03 |
| 101.1884 | 1.90E+01 |
| 101.5234 | 6.89E+00 |
| 101.9726 | 1.20E+01 |

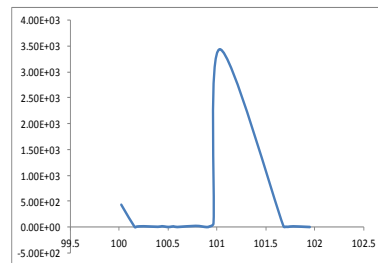Technical rep 1

Technical rep 2

Technical rep 3
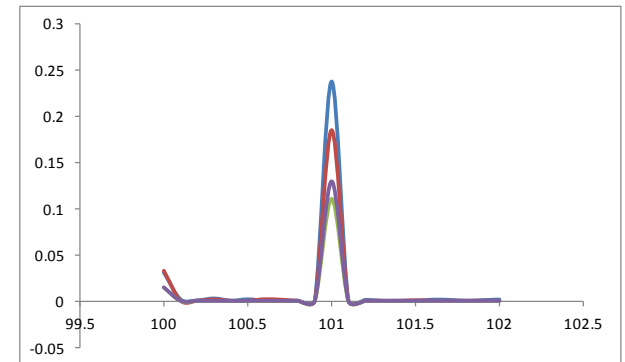
# •Align and "Bin" raw spectral output files

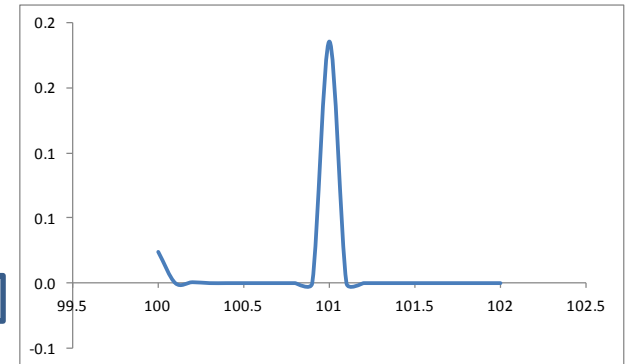"Bin" 101 = contains all data between 100.95 to 101.05

Peak alignment



| | |
|---|---|
| 100 | 0.024 |
| 100.1 | 0 |
| 100.2 | 0.0007 |
| 100.3 | 0 |
| 100.4 | 0 |
| 100.5 | 0 |
| 100.6 | 0 |
| 100.7 | 0 |
| 100.8 | 0 |
| 100.9 | 0 |
| 101 | 0.1855 |
| 101.1 | 0 |
| 101.2 | 0 |
| 101.3 | 0 |
| 101.4 | 0 |
| 101.5 | 0 |
| 101.6 | 0 |
| 101.7 | 0 |
| 101.8 | 0 |
| 101.9 | 0 |
| 102 | 0 |

Compare peak height for that bin
across many samples

Normalise height to %TIC (TIC/sum total TIC * 100)

# Why do we merge? - help with identifying many metabolites

•Identify many compounds – earlier we looked at identifying one metabolite
•Very few sites allows the searching for many metabolites

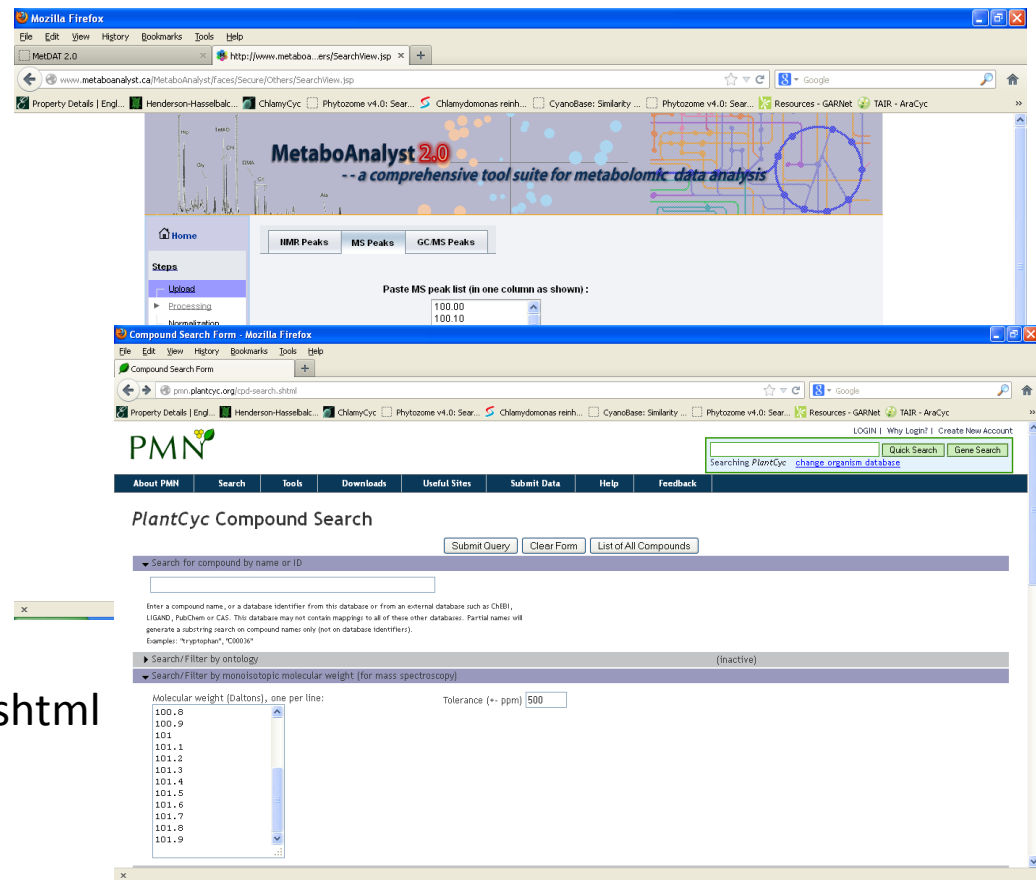Equipment software

Usually lab specific software
– excel macros

**METABOANALYST**
**http://www.metaboanalyst.ca/MetaboAnalys
t/faces/Secure/Others/SearchView.jsp**

**PlantCyc**
http://pmn.plantcyc.org/cpd-search.shtml

http://smbl.nus.edu.sg/METDAT2/

# •Align and "Bin" raw spectral output files
# •METALIGN

# Matlab – MathWorks - commercial



http://www.mathworks.co.uk/products/bioinfo/examples.html?file=/products/demos/shipping/bioinfo/mspreprodemo.html

# SIM-Stitch    University of Birmingham

# •Most people merge files for statistical analysis

# •Multivariate data analysis – PCA, Cluster analysis

# Multivariate Data Analysis

**Unsupervised**

**Principal Component Analysis (PCA)**

**Supervised**

**Partial Least Squares**

**-Discriminant Analysis (PLS-DA)**

**Hierarchical Cluster Analysis (HCA)**

Trygg et al. 2007

# PCA – Principal Component Analysis

**Objectives of PCA**

•reduce number of variables

•identify outliers

•identify any splits within the data

•discriminate between samples (eg, cold stress, GM, disease, control) or separation by other unplanned means (eg, analytical error)

How does it work?

# PCA – Principal Component Analysis



If each variable (ie. metabolite ion intensity or concentration) is thought of as a dimension,

and there are *n* variables,

every sample is at a unique position in the *n*-dimensional space defined by these *n* variables.

Very difficult for people to visualise – aim is to reduce this dimensional space by summarising the data using relatively few parameters

Ie. Make a 2D plot!

Will attempt to explain PCA…

# PCA – Principal Component Analysis

Scores out of 10

| | Banana | Apple | Melon |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

Plot these data in K-space

# PCA – Principal Component Analysis

Bendy-ess:

| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



Yellow-ness:

Bendy-ess:

| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis

| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



Yellow-ness (y)

Bendy-ess (x)

Round-ness (z)

| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis

Yellow-ness (y)

Bendy-ess (x)

Round-ness (z)

| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



Yellow-ness (y)

Bendy-ess (x)

Round-ness (z)

| | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



|  | | | |
|---|---|---|---|
| Bendy-ess: | 8 | 2 | 3 |
| Yellow-ness: | 9 | 2 | 7 |
| Round-ness: | 3 | 8 | 5 |

# PCA – Principal Component Analysis



A 2D window is inserted over this 3D space (or in real PCA datasets over all k-dimensions) that covers the most variation of the data set

A line is then placed through this window and a dot is placed in the **centre (mean centred)**

**This is what we call PC (principal component) 1 !**

# PCA – Principal Component Analysis



The second PC is calculated by looking at the variation at 90 ° to the first PC

The third PC is calculated by looking at the variation at 90 ° to the second PC

Each subsequent PC lies in an orthogonal direction of maximum variance that has not been considered by the former components.

# PCA – Principal Component Analysis



Yellow-ness (y)  PC1

PC3

Bendy-ess (x)

PC2

Round-ness (z)

PC2

PC1

Rotating the window converts the multidimension data to a 2D PCA plot
(this is called a score scatter plot)

**How do we know what causes each point to be in each position? – loadings plot**

# PCA – Principal Component Analysis

How do we see what causes each point to be in each position? – loadings plot
Displays the relationships among the variables

Loading score for each measurement is given between 0-1 (essentially an R2)

Are inversely related measurement – what is high in bendy is low in opposite measurement

PCA Score Scatter Plot

Aqueous phase Negative ionisation

Direct Injection Mass Spectrometry (DIMS)

30% of ALL the variation in the masses detected ALL samples is explained in PC1

Davey et al. 2008/2009

PCA Score Loadings Plot

0.2Da bins

Overlay PCA Score Scatter and Loadings Plot

PC2 (22%)

PC1 (30%)

SIMCA-P 11 - 09/04/2008 15:31:02

SIMCA-P 11 - 09/04/2008 15:26:59

# Simca-P software - PCA



**Data mining** – search masses in
our metabolite database

# SMILE STRUCTURE REPRESENTATION



M1.t[1] = 82.5556
M1.t[2] = -4.89724
Primary ID = 2431
**MASS 122 – Malonate**

M1.t[1] = 45.7915
M1.t[2] = -8.28212
Primary ID = 448
**MASS 115 – Fumaric acid**

M1.t[1] = 25.4171
M1.t[2] = 12.544
Primary ID = 1786

**MASS 133 – Malic acid**

M1.t[1] = 9.94449
M1.t[2] = -3.22985
Primary ID = 1730

**MASS 74 - Glycine**

# Scaling data

## Reducing bias against very intense peaks
## -need to highlight fold/relative intensities among samples of the same peak

# Scaling data

**Reducing bias against very intense peaks**
**-need to highlight fold/relative intensities among samples of the same peak**



**Unit variance scaling**
Most objective – takes data as face value

**Pareto Scaling**
Best for MS, NMR
Decreases importance of high intensity peaks
Depends on question…

# How many components should I look at?

How many components should be included in the model?

Degree of fit and the predictive ability

**Fit** = how well we are able to mathematically reproduce the data of a training set (goodness of fit) $R^2X$ = the explained variation (0-1)

**Predictive ability** = how accurately can we predict the raw X data? (goodness of predictability) $Q^2X$

Take out 25% of samples, does the PCA still look the same? **(a good reason to have many samples!)**

Does adding more components make the model better?

Poor $R^2X$ and $Q^2X$ if very noisy dataset or too few replicates

# •PLS-DA – projections to latent structures by partial least squares – discriminate analysis

•Similar to PCA

•Supervised method – good for discovering biomarkers

•Can connect other data (Y) with your MS/NMR dataset (X)
– you tell the model that some samples are from control, others treated etc.



R2X[1] = 0.365161       SIMCA-P 11.5 - 15/11/2007 17:44:56

# HCA – Hierarchical Cluster Analysis

# •ANOVA, false-positive corrections
## Which masses are statistically significantly different between our samples?
## T-test
## ANOVA

Many samples (10's or 100's)

Many bins (MW)
1000's

Many intensities
10,000's

# •ANOVA, false-positive corrections

Testing so many samples runs the risk of significant values occurring by chance
Change the usual p <0.05 to take into account the number of samples

**p<0.05** is usually considered a significant result - error rate of 5 %
(**ie: 1 in every 20 tests gives a false result**)

If a dataset has 1000 metabolites (so 1000 univariate tests…)
**50 metabolites are significantly different when they are not**

# •ANOVA, false-positive corrections

## Testing so many samples runs the risk of significant values occurring by chance
## Change the usual p <0.05 to take into account the number of samples

**p<0.05** is usually considered a significant result - error rate of 5 %
(**ie: 1 in every 20 tests gives a false result**)

If a dataset has 1000 metabolites (so 1000 univariate tests…)
**50 metabolites are significantly different when they are not**

Two main False-positive corrections

**Bonferroni**   adjusted P-value is: 0.05/number of samples (very strict)

=0.05/4254        $P_{adj} = 0.00001175$

**Benjamini and Hochberg** – rank the p-values highest to lowest
$P_{adj}$ = (number of samples/p-value)*position in ranked p-value table
Keep going until you get to $P_{adj}$ of 0.05
Original p-value = 0.011718
$P_{adj} = 0.051755$
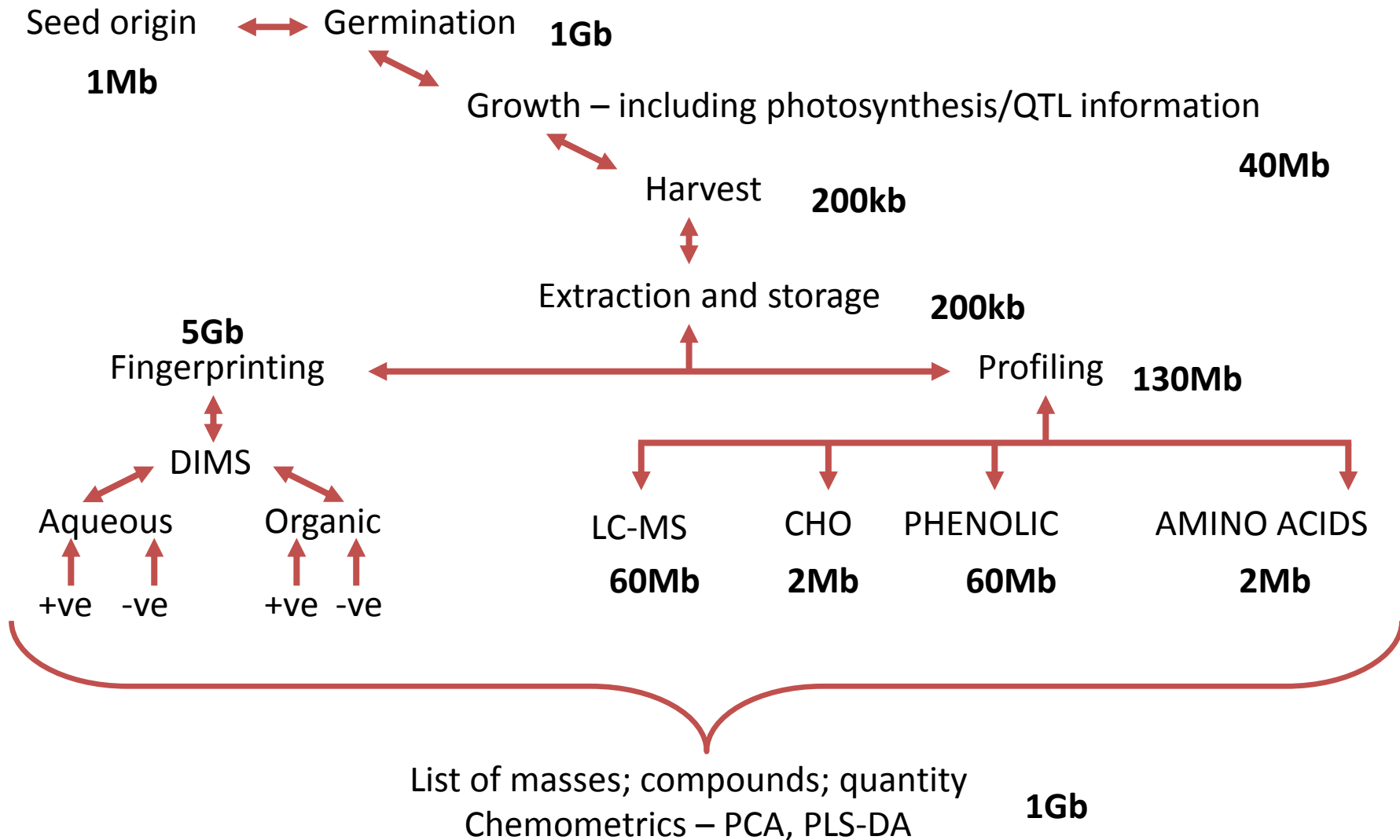
Benjamini and Hochberg.  1995. Journal of the Royal Statistical Society, Series B 57: 289–300.

# Introduction to chemometrics

• Metadata – useful for interpreting PCA, Y-var for PLS-DA

|  | **META-DATA** |
| --- | --- |
| Experiment set-up | Seed origin |
|  | Germination conditions – soil, light, humidity, CO2, day length, temp |
|  | Growth conditions – randonisation etc |
|  | Harvest conditions, time, weight of plant |
| Gas exchange | Time, growth conditions, IRGA conditions saved per raw sample run |
| Growth | Plant and cabinet data |

# Metadata – data storage

# Online practical involving multivariate statistics – PCA etc