Data Analysis & Biomarker Discovery
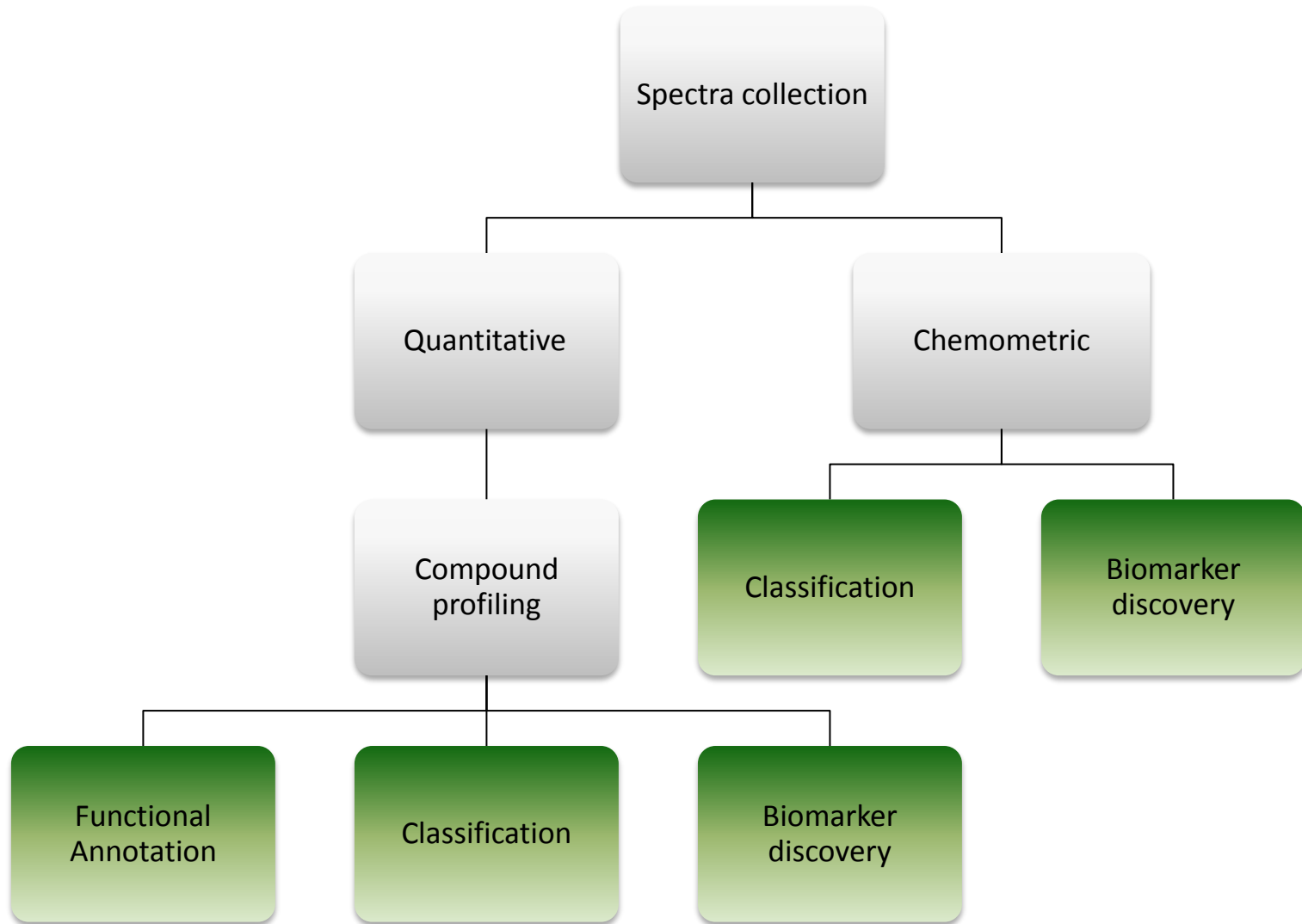
# MetaboAnalyst 2.0 & ROCCET
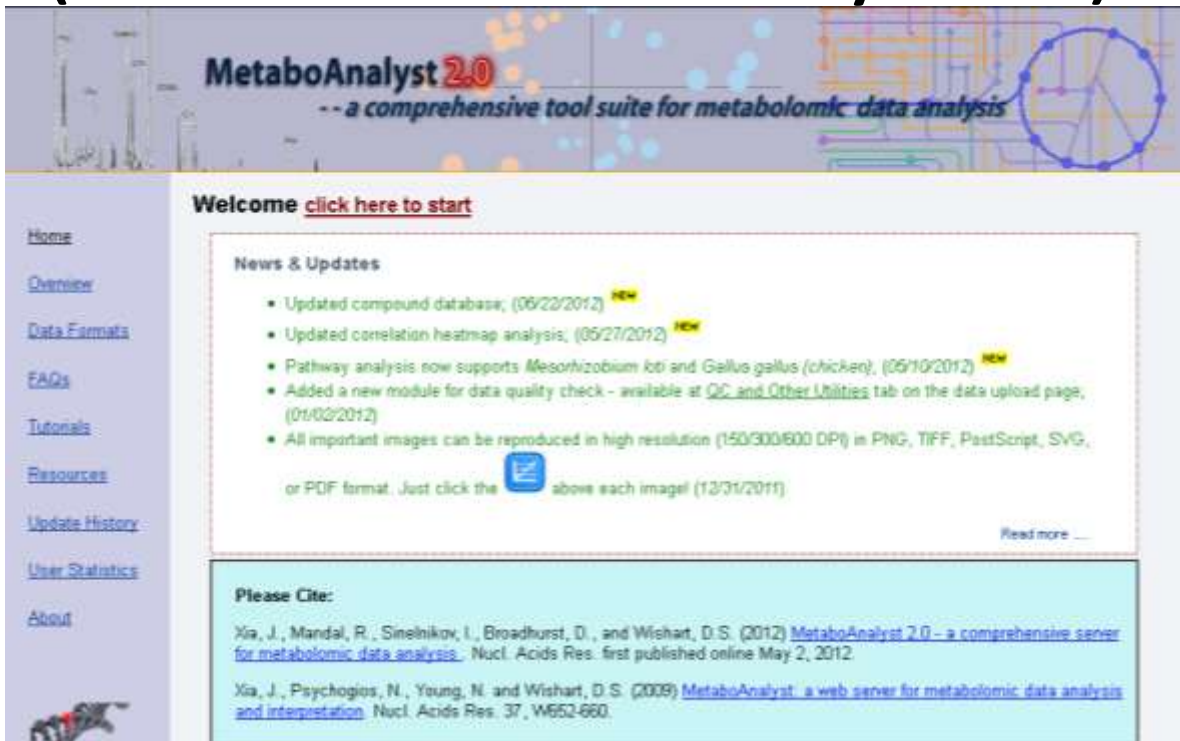
Jianguo Xia, PhD

University of Alberta, Canada

# Outline

- Introduction (updates) of two free web application
  - MetaboAnalyst 2.0
  - ROCCET
- Background & basic concepts
- Screenshot tutorials
- Live demo (if we have time)

# Metabolomic Data Analysis

# MetaboAnalyst
# (www.metaboanalyst.ca)
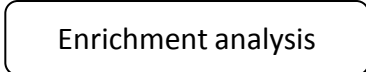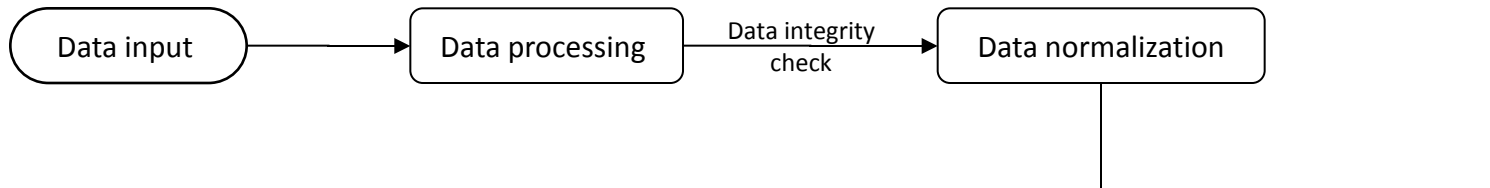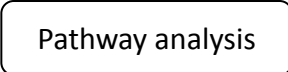
- GC/LC-MS raw spectra
- Peak lists
- Spectral bins
- Concentration table

- Spectra processing
- Peak processing
- Noise filtering
- Missing value estimation

- Row-wise normalization
- Column-wise normalization
- Combined approach

**Data input** → **Data processing** — Data integrity check → **Data normalization**

Functional Interpretation

Statistical Exploration

**Enrichment analysis**

**Pathway analysis**

**Time-series analysis**

**Two/multi-group analysis**

- Over representation analysis
- Single sample profiling
- Quantitative enrichment analysis

- Enrichment analysis
- Topology analysis
- Interactive visualization

- Data overview
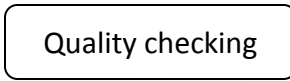- Two-way ANOVA
- ANOVA - SCA
- Time-course analysis

- Univariate analysis
- Correlation analysis
- Chemometric analysis
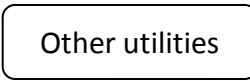- Feature selection
- Cluster analysis
- Classification

**Outputs**

**Image Center**

**Quality checking**

**Other utilities**

- Processed data
- Result tables
- Analysis report
- Images

- Resolution: 150/300/600 dpi
- Format: png, tiff, pdf, svg, ps

- Methods comparision
- Temporal drift
- Batch effect
- Biolgoical checking

- Peak searching
- Pathway mapping
- Name/ID conversion
- Lipidomics

# MetaboAnalyst Overview

- Raw data processing

- Data reduction & statistical analysis

- Functional enrichment analysis

- Metabolic pathway analysis

- Quality control analysis
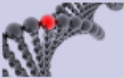
# Data processing overview

- Supported data formats
  - Concentration tables
  - Peak lists
  - Spectral bins
  - Raw spectra (* not recommended)

# Example Datasets

# Data Processing

Purpose: to convert various raw data forms into data matrices suitable for statistical analysis

# Data Upload

# Alternatively …



**2) Try our test data :** ( You can download these data here )

| Data Type | Description |
| --- | --- |
| ○ **Concentrations**<br>Tutorial\|Report | Metabolite concentrations of 77 urine samples from cancer patients measured by 1H NMR (Eisner R, et al.). Group 1- cachexic; group 2 - control |
| ○ **Concentrations** | Metabolite concentrations of 39 rumen samples measured by proton NMR from dairy cows fed with different proportions of barley grain (Ametaj BN, et al.). Group label - 0, 15, 30, or 45 - indicating the percentage of grain in diet. |
| ○ **NMR spectral bins**<br>Tutorial\|Report | Binned 1H NMR spectra of 50 urine samples using 0.04 ppm constant width (Psihogios NG, et al.) Group 1- control; group 2 - severe kidney disease. |
| ○ **NMR peak lists** | Peak lists and intensity files for 50 urine samples measured by 1H NMR (Psihogios NG, et al.). Group 1- control; group 2 - severe kidney disease. |
| ○ **Concentrations (paired)**<br>Tutorial\|Report | Compound concentrations of 14 urine samples collected from 7 cows at two time points using 1H NMR (unpublished data). Group 1- day 1, group 2- day 4. |
| ○ **MS peak intensities** | LC-MS peak intensity table for 12 mice spinal cord samples (Saghatelian et al.). Group 1- wild-type; group 2 - knock-out. |
| ○ **MS peak lists** | Three-column LC-MS peak list files for 12 mice spinal cord samples (Saghatelian et al.). Group 1- wild-type; group 2 - knock-out. |
| ○ **LC-MS spectra**<br>Tutorial\|Report | NetCDF spectra of 12 mice spinal cord samples collected by LC-MS (Saghatelian et al.). Group 1- wild-type; group 2 - knock-out. |
| ○ **GC-MS spectra** | NetCDF spectra collected by GC-MS. This is a dummy data set for testing spectra processing only. Each group is a triplicate of a single spectrum . Group 1- Sunflower oil, group 2- Olive oil. |

Submit

# Data Integrity Check

# Data Normalization

# Normalization Result

# Quality Control

- Dealing with outliers
  - Detected mainly by visual inspection
  - May be corrected by normalization
  - May be excluded
- Dealing with missing values
- Noise reduction

# Visual Inspection

- What does an outlier look like?



Finding outliers via PCA        Finding outliers via Heatmap

| Statistical Analysis | Enrichment Analysis | Pathway Analysis | Time Series | QC & Other Utilities |

## Functions for Quality Check

### Comparing the Agreement between Two Measurements

In metabolomics researches, different protocols are often explored to find to best approach. The function allows you to visually compare the agreement between two measurements and to detect outliers.
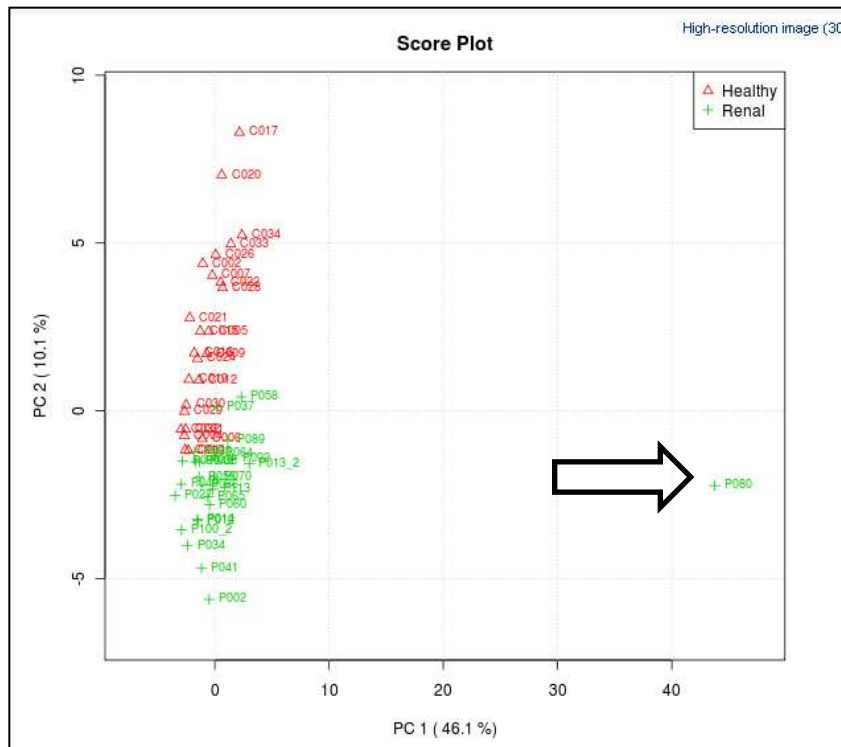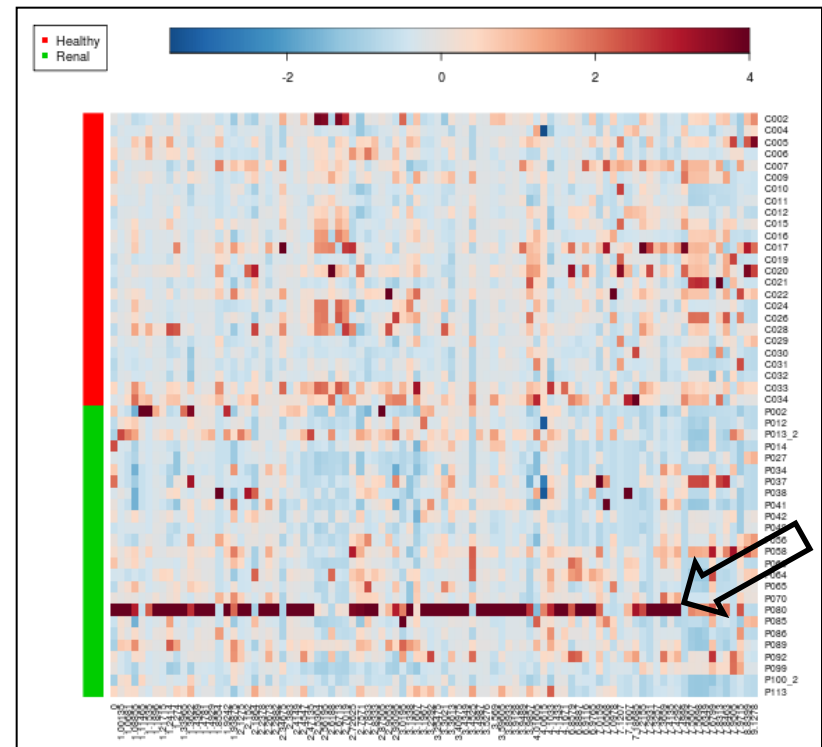
### Detecting (and Correcting) for Time Drift

The method aims to detect if temporal drift is present in the measurements collected over a long period of time. User can adjust the time window to calculate pair-wise p values between data points measured at each time frame. Finally, the method allows users to correct the drift using the LOWESS correction.
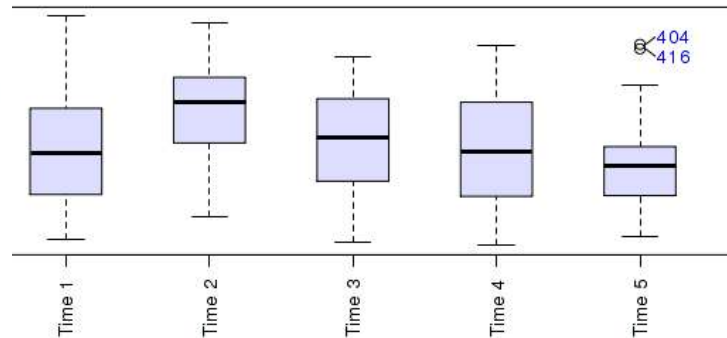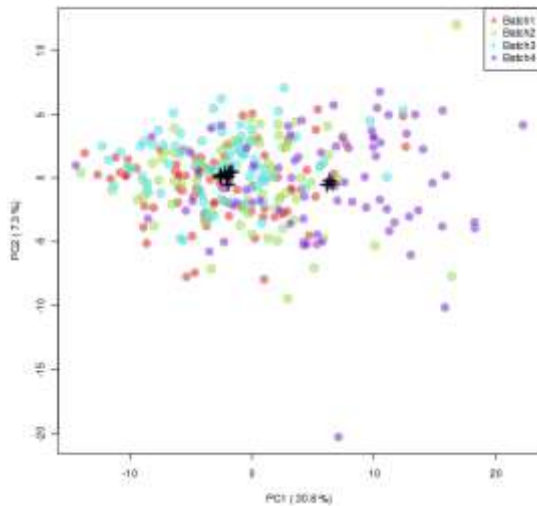
### Checking Batch Effects for Large Number of Samples

The methods aims to detect the batch effect in large scale metabolomics studies with a **randomized experiment design**. The method allows high-level visualization of samples in each batch using scatter plot, boxplot, heatmap and principal component analysis (PCA).

### Checking against reference concentrations in HMDB

The methods compares the measured concentration values in user data against the normal reference values stored in HMDB. Therefore, the comparison is only meaningful for **human biofluid samples (blood/urine/CSF)**. The approach is useful to examine sample qualities, wrong labels, etc.

# Quality Check Module

# Outlier Removal

# Data Filtering

- Characteristics of noise & uninformative features
  - Low intensities
  - Low variances (default)

# Noise Reduction



**Steps**
- Home
- Upload
- ▼ Processing
  - Pre-process
  - Name check
  - Conc. check
  - Data check
  - Missing value
  - → Data filter
  - Data editor
  - Color picker
  - Normalization
- ► Statistics
- ► Enrichment
- ► Pathway
- ► Time Series
- ► Peak search
- ► Metabolites
- Download
- Log out

## Data filtering

The purpose of the data filtering is to identify and remove variables that are unlikely to be of use when modeling the data. No phenotype information are used in the filtering process, so the result can be used with any downstream analysis. This step is strongly recommended for chemometrics datasets (i.e. spectral binning data) with large number of variables, many of them are from baseline noises. Filtering can usually improve the results. For details, please see the paper by Hackstadt, et al.

Non-informative variables can be characterized in two groups:

1. Variables of very small values - these variables can be detected using **mean** or the robust estimate **median** which is not affected by extreme values or outliers;
2. Variables that are near-constant throughout the experiment conditions - these variables can be detected using **standard deviation (SD)** or the robust estimate **interquantile range (IQR)**. The **coefficient of variation (CV)** (CV=mean/SD) is another useful variance measure independent of the mean.

The following empirical rules are applied during data filtering:

| Number of Variables | Variables Filtered |
|---|---|
| < 250 | 5% |
| 250 - 500 | 10% |
| 500 - 1000 | 25% |
| > 1000 | 40 % |

Please note, in order to reduce the computational burden to the server, the maximum allowed number of variables is 5000. If over 5000 variables were left after filtering, only the top 5000 will be used in the subsequent analysis.

# Missing values

# Dimension Reduction & Statistical Analysis

# Common tasks

- To identify important features;
- To detect interesting patterns;
- To assess difference between the phenotypes
- To facilitate classification / prediction

# Home

## Steps

Upload

► Processing

▼ Statistics

- Fold change
- T-test
- Volcano plot
- ANOVA
- Correlations
- PCA
- PLSDA
- SAM
- EBAM
- Dendrogram
- Heatmap
- SOM
- K-means
- RandomForest
- SVM

► Enrichment

► Pathway

► Time Series

► Peak search

► Metabolites

Download

Log out

## Select an analysis path to explore :

**Univariate Analysis**

Fold Change Analysis, t-Tests, and Volcano plot **(two-group only)**

One-way ANOVA and Correlation Analysis

**Multivariate Analysis**

Principal Component Analysis (PCA)

Partial Least Squares - Discriminant Analysis (PLS-DA)

**Significant Feature Identification**

Significance Analysis of Microarray (and Metabolites) (SAM)

Empirical Bayesian Analysis of Microarray (and Metabolites) (EBAM)

**Cluster Analysis**

Hierarchical Clustering - Dendrogram and Heatmap

Partitional Clustering - K-Means and Self Organizing Map (SOM)

**Classification & Feature Selection**

Random Forest

Support Vector Machine (SVM) **(two-group only)**

# ANOVA

# View Individual Compounds

# Overall correlation pattern

# High resolution image

# Template Matching

- Looking for compounds showing interesting patterns of change
- Essentially a method to look for linear trends or periodic trends in the data
- Best for data that has 3 or more groups

# Template Matching (cont.)



**Top 25 compounds correlated with the 1-2-3-4**

Strong linear + correlation to grain %

Strong linear - correlation to grain %

# PCA Scores Plot

# PCA Loading Plot

# PLS-DA Score Plot

# Evaluation of PLS-DA Model

- PLS-DA Model evaluated by cross validation of $Q^2$ and $R^2$

- More components to model improves quality of fit, but try to minimize this value

- 3 Component model seems to be a good compromise here

- Good $R^2/Q^2$ (>0.7)

# Important Compounds

# Model Validation

# Heatmap Visualization

# Heatmap Visualization (cont.)

# Download Results

# Analysis Report

## 2.2 Correlation Analysis

Correlation analysis can be used to identify which features are correlated with a feature of interest. Correlation analysis can also be used to identify if certain features show particular patterns under different conditions. Users first need to define a pattern in the form of a series of hyphenated numbers. For example, in a time-series study with four time points, a pattern of of 1-2-3-4 is used to search compounds with increasing the concentration as time changes; while a pattern of 3-2-1-3 can be used to search compounds that decrease at first, then bounce back to the original level.

Figure 3 shows the important features identified by correlation analysis. Table 3 shows the details of these features.

Table 3: Important features identified by Pattern search using correlation analysis

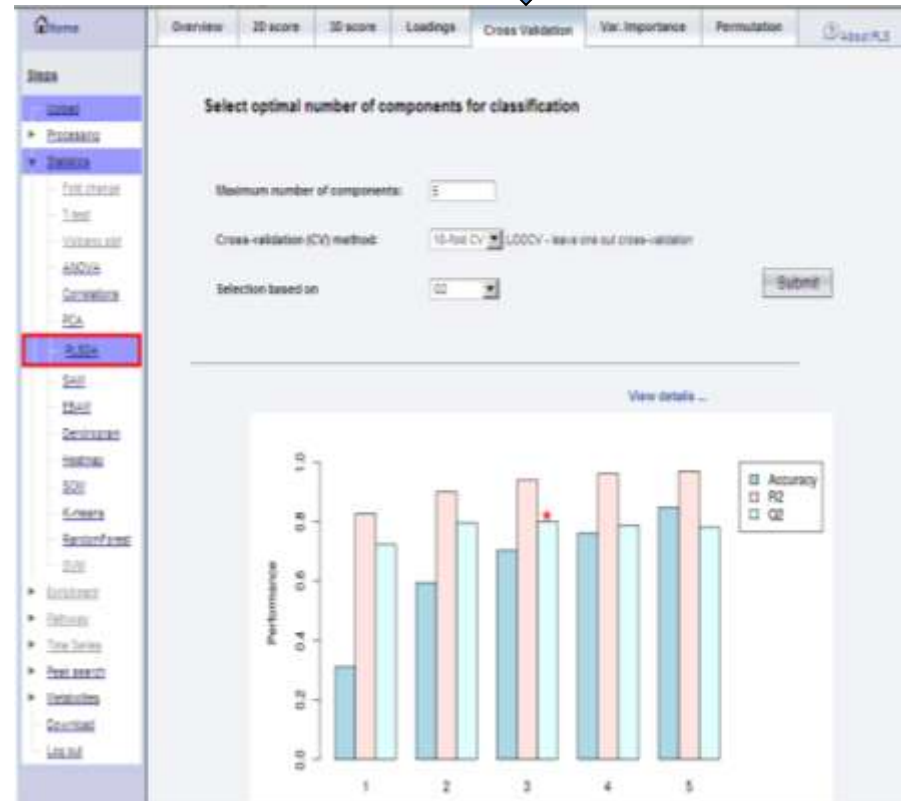| | Compounds | correlation | t-stat | p-value | FDR |
|---|---|---|---|---|---|
| 1 | Butyrate | -0.51252 | 15032 | 3.4067e-06 | 0.00080058 |
| 2 | Isobutyrate | -0.49758 | 15784 | 5.9015e-06 | 0.00092488 |
| 3 | 3-PP | -0.47935 | 15638 | 0.00014063 | 0.0016524 |
| 4 | Acetate | -0.45453 | 15359 | 0.00024011 | 0.0025416 |
| 5 | 3-HB | -0.41943 | 14024 | 0.007862 | 0.041057 |
| 6 | Isovalerate | -0.39861 | 13818 | 0.011086 | 0.086193 |
| 7 | Lysine | -0.34401 | 12291 | 0.15439 | 0.30381 |
| 8 | Methanol | -0.24287 | 12277 | 0.13978 | 0.30381 |
| 9 | Ferulate | -0.22929 | 12148 | 0.16028 | 0.32753 |
| 10 | Fumarate | -0.21966 | 12050 | 0.17906 | 0.33398 |
| 11 | Histidine | -0.2169 | 12023 | 0.18474 | 0.33398 |
| 12 | Propionate | -0.21015 | 11956 | 0.19912 | 0.34861 |
| 13 | Maltose | -0.2003 | 11859 | 0.22148 | 0.37177 |
| 14 | Acetoacetate | -0.17772 | 11638 | 0.27907 | 0.39248 |
| 15 | Choline | -0.13886 | 11054 | 0.47111 | 0.65124 |
| 16 | Tyrosine | -0.10857 | 10933 | 0.51847 | 0.67680 |
| 17 | PAG | -0.079788 | 10668 | 0.62921 | 0.79927 |
| 18 | 3-HP | -0.074918 | 10620 | 0.65058 | 0.80438 |
| 19 | Formate | -0.061347 | 10387 | 0.78623 | 0.84626 |
| 20 | Aspartate | -0.031981 | 10198 | 0.84674 | 0.86518 |
| 21 | Caffeine | 0.011541 | 9763 | 0.84297 | 0.84297 |
| 22 | Ribose | 0.039063 | 9405.1 | 0.81387 | 0.85004 |
| 23 | 1,3-D | 0.043158 | 9453.6 | 0.79419 | 0.84854 |
| 24 | Succinate | 0.04804 | 9435 | 0.78842 | 0.84854 |
| 25 | Glucose | 0.087544 | 9311.6 | 0.72787 | 0.83439 |
| 26 | Cadaverine | 0.060643 | 9280.8 | 0.71382 | 0.83439 |
| 27 | Phenylacetate | 0.063742 | 9260.2 | 0.69086 | 0.83439 |
| 28 | Hypoxanthine | 0.10911 | 8802 | 0.50847 | 0.67689 |
| 29 | Ethanol | 0.18304 | 8071.6 | 0.26471 | 0.3888 |
| 30 | NDMA | 0.18492 | 8063 | 0.25975 | 0.3888 |
| 31 | Proline | 0.18713 | 8031.2 | 0.25399 | 0.3888 |
| 32 | Glutamate | 0.19354 | 7969.8 | 0.23929 | 0.38619 |
| 33 | Benzoate | 0.21978 | 7708.6 | 0.17884 | 0.33398 |
| 34 | Valerate | 0.23936 | 7515.1 | 0.14221 | 0.30381 |
| 35 | Glycerol | 0.26991 | 7213.3 | 0.096569 | 0.23888 |
| 36 | Glycine | 0.28064 | 7107.3 | 0.083833 | 0.21812 |
| 37 | Nicotinate | 0.28612 | 7063 | 0.078511 | 0.21706 |
| 38 | Methylamine | 0.28905 | 7024.2 | 0.07451 | 0.21706 |
| 39 | Isoleucine | 0.30355 | 6881 | 0.060303 | 0.18896 |
| 40 | Xanthine | 0.30854 | 6861.3 | 0.058555 | 0.18896 |
| 41 | Dimethylamine | 0.33298 | 6600.1 | 0.038526 | 0.13856 |
| 42 | Leucine | 0.35142 | 6407.9 | 0.028264 | 0.11068 |
| 43 | Valine | 0.3809 | 6116.7 | 0.016744 | 0.071541 |
| 44 | Lactate | 0.42384 | 5692.5 | 0.0071709 | 0.041057 |
| 45 | Uracil | 0.45172 | 5417 | 0.0038928 | 0.026137 |
| 46 | Endotoxin | 0.50141 | 4926.1 | 0.0011471 | 0.0089853 |
| 47 | Alanine | 0.62026 | 3751.5 | 2.5337e-06 | 0.00080058 |

## 2.5 Hierarchical Clustering

In (agglomerative) hierarchical cluster analysis, each sample begins as a separate cluster and the algorithm proceeds to combine them until all samples belong to one cluster. Two parameters need to be considered when performing hierarchical clustering. The first one is similarity measure - Euclidean distance, Pearson's correlation, Spearman's rank correlation. The other parameter is clustering algorithms, including average linkage (clustering uses the centroids of the observations), complete linkage (clustering uses the farthest pair of observations between the two groups), single linkage (clustering uses the closest pair of observations) and Ward's linkage (clustering to minimize the sum of squares of any two clusters). Heatmap is often presented as a visual aid in addition to the dendrogram.

Hierarchical clustering is performed with the hclust function in package stat. Figure 17 shows the clustering result in the form of a dendrogram. Figure 18 shows the clustering result in the form of a heatmap.
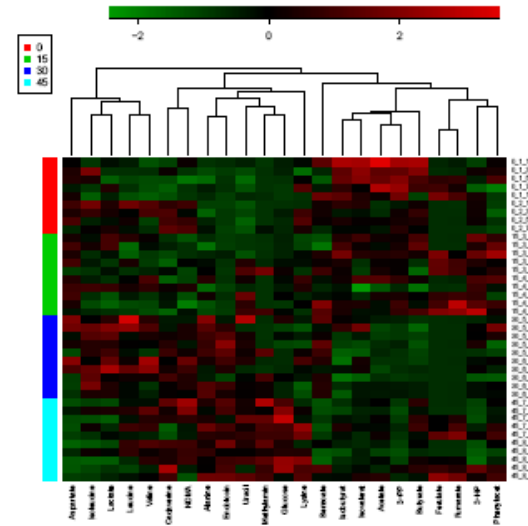
Figure 17: Clustering result shown as heatmap (distance measure using pearson, and clustering algorithm using ward).

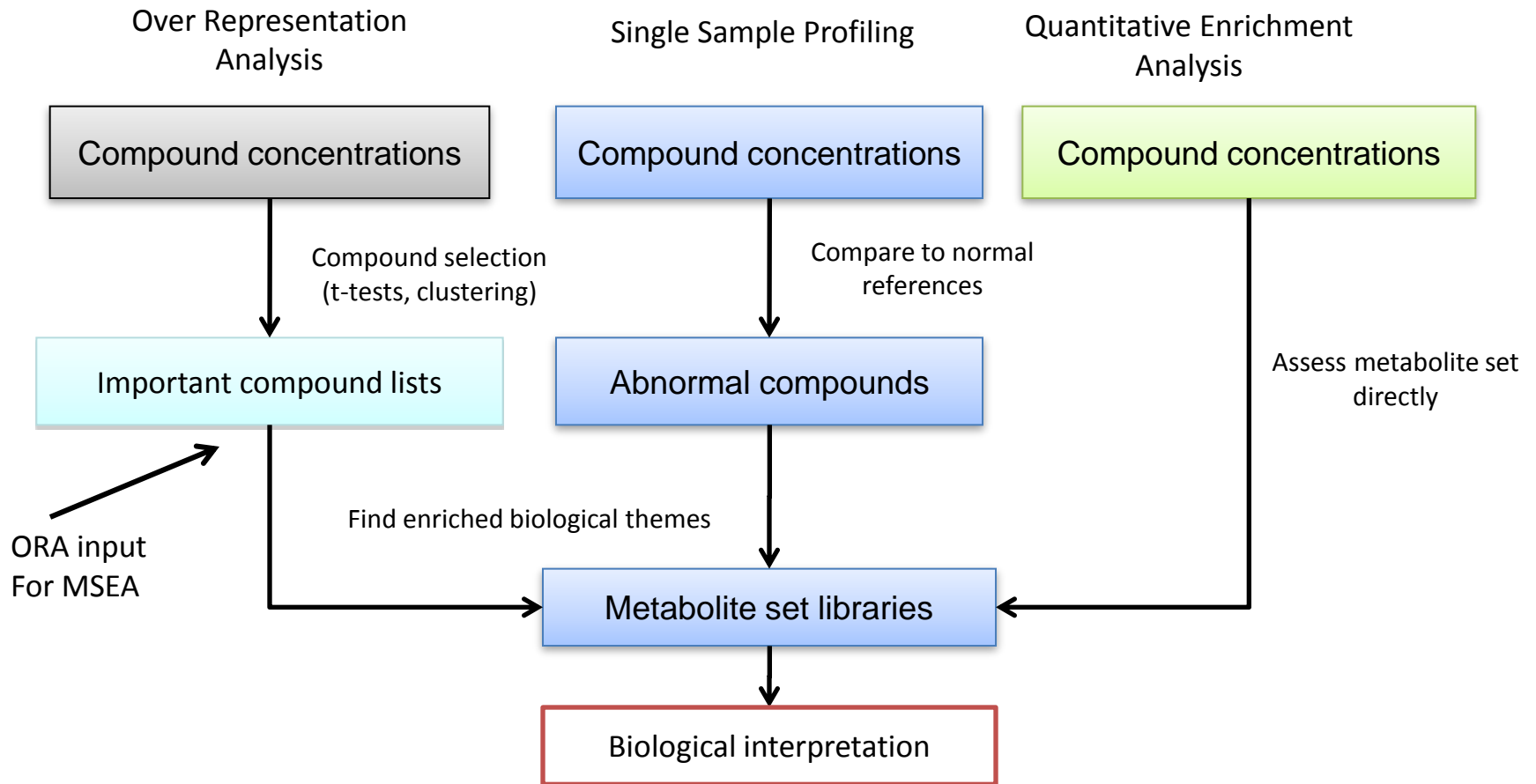# Metabolite Set Enrichment Analysis (MSEA)

# Enrichment Analysis

- Purpose: To test if there are some <span style="color:red">biologically meaningful groups</span> of metabolites that are significantly enriched in your data

- Biological meaningful groups
  - Pathways
  - Disease
  - Localization

- Currently, only supports human metabolomic data

# MSEA

- Accepts 3 kinds of input files
- 1) list of metabolite names only (ORA)
- 2) list of metabolite names + concentration data from a single sample (SSP)
- 3) a concentration table with a list of metabolite names + concentrations for multiple samples/patients (QEA)

# The MSEA approach

Over Representation Analysis

Single Sample Profiling

Quantitative Enrichment Analysis

Compound concentrations

Compound concentrations

Compound concentrations

Compound selection (t-tests, clustering)

Compare to normal references

Assess metabolite set directly

Important compound lists

Abnormal compounds

ORA input
For MSEA

Find enriched biological themes

Metabolite set libraries

Biological interpretation

# Start with a compound List

# Upload Compound List

# Compound Name Standardization

## Compound Label Standardization:

PLease note:

- Query names in normal white indicate exact match - marked by "1" in the download file;
- Query names highlighted in yellow indicate **approximate matches** (for compound name matches) - marked by "2" in the downloaded file. Users should manually select the correct match by clicking the **View** link of the corresponding compounds. Otherwise, the first match will be used;
- Query names highlighted in red indicate **no match** - marked by "0" in the downloaded file;
- Greek alphabets are not recognized, they should be replaced by English names (i.e. alpha, beta)

| Query | Match | HMDB | PubChem | KEGG | Details |
|-------|-------|------|---------|------|---------|
| Acetoacetic acid | Acetoacetic acid | HMDB00060 | 96 | C00164 | |
| Beta-Alanine | Beta-Alanine | HMDB00056 | 239 | C00099 | |
| Creatine | Creatine | HMDB00064 | 586 | C00300 | |
| Dimethylglycine | Dimethylglycine | HMDB00092 | 673 | C01026 | |
| Fumaric acid | Fumaric acid | HMDB00134 | 723 | C00122 | |
| Glycine | Glycine | HMDB00123 | 750 | C00037 | |
| Homocysteine | Homocysteine | HMDB00742 | 778 | C05330 | |
| L-Cysteine | L-Cysteine | HMDB00574 | 5862 | C00097 | |
| L-Isolucine | L-Isoleucine | HMDB00172 | 791 | C00407 | View |
| L-Phenylalanine | L-Phenylalanine | HMDB00159 | 6140 | C00079 | |
| L-Serine | L-Serine | HMDB00187 | 5951 | C00065 | |
| L-Threonine | L-Threonine | HMDB00167 | 6288 | C00188 | |
| L-Tyrosine | L-Tyrosine | HMDB00158 | 6057 | C00082 | |
| L-Valine | L-Valine | HMDB00883 | 1182 | C00183 | |
| Phenylpyruvic acid | Phenylpyruvic acid | HMDB00205 | 997 | C00166 | |
| Propionic acid | Propionic acid | HMDB00237 | 1032 | C00163 | |
| Pyruvic acid | Pyruvic acid | HMDB00243 | 1060 | C00022 | |
| Sarcosine | Sarcosine | HMDB00271 | 1088 | C00213 | |

# Name Standardization (cont.)

# Select a Metabolite Set Library

# Result

# Result (cont.)

| | Metabolite Set | Total | Hits | Expect | P value | Holm P | FDR | Details |
|---|---|---|---|---|---|---|---|---|
| | GLYCINE, SERINE AND THREONINE METABOLISM | 26 | 9 | 0.567 | 2.74E-10 | 2.19E-8 | 2.19E-8 | View |
| | PROTEIN BIOSYNTHESIS | 19 | 6 | 0.415 | 9.93E-7 | 7.85E-5 | 3.97E-5 | View |
| | PHENYLALANINE AND TYROSINE METABOLISM | 13 | 5 | 0.284 | 3.15E-6 | 2.46E-4 | 8.4E-5 | View |
| | METHIONINE METABOLISM | 24 | 5 | 0.524 | 8.98E-5 | 0.00691 | 0.0018 | View |
| | AMMONIA RECYCLING | 18 | 3 | 0.393 | 0.00581 | 0.441 | 0.0774 | View |
| | PROPANOATE METABOLISM | 18 | 3 | 0.393 | 0.00581 | 0.441 | 0.0774 | View |
| | CYSTEINE METABOLISM | 8 | 2 | 0.175 | 0.0117 | 0.863 | 0.133 | View |
| | GLUTATHIONE METABOLISM | 10 | 2 | 0.218 | 0.0183 | 1.0 | 0.162 | View |
| | BETAINE METABOLISM | 10 | 2 | 0.218 | 0.0183 | 1.0 | 0.162 | View |
| | ASPARTATE METABOLISM | 12 | 2 | 0.262 | 0.0261 | 1.0 | 0.209 | View |
| | VALINE, LEUCINE AND ISOLEUCINE DEGRADATION | 36 | 3 | 0.785 | 0.0397 | 1.0 | 0.288 | View |
| | TYROSINE METABOLISM | 38 | 3 | 0.829 | 0.0456 | 1.0 | 0.304 | View |
| | UREA CYCLE | 20 | 2 | 0.436 | 0.0677 | 1.0 | 0.417 | View |
| | CITRIC ACID CYCLE | 23 | 2 | 0.502 | 0.0868 | 1.0 | 0.496 | View |
| | CATECHOLAMINE BIOSYNTHESIS | 5 | 1 | 0.109 | 0.105 | 1.0 | 0.536 | View |
| | ARGININE AND PROLINE METABOLISM | 26 | 2 | 0.567 | 0.107 | 1.0 | 0.536 | View |
| | ALANINE METABOLISM | 6 | 1 | 0.131 | 0.124 | 1.0 | 0.585 | View |
| | TAURINE AND HYPOTAURINE METABOLISM | 7 | 1 | 0.153 | 0.144 | 1.0 | 0.638 | View |
| | BUTYRATE METABOLISM | 9 | 1 | 0.196 | 0.181 | 1.0 | 0.758 | View |
| | PANTOTHENATE AND COA BIOSYNTHESIS | 10 | 1 | 0.218 | 0.199 | 1.0 | 0.758 | View |
| | KETONE BODY METABOLISM | 10 | 1 | 0.218 | 0.199 | 1.0 | 0.758 | View |
| | GLUCOSE-ALANINE CYCLE | 12 | 1 | 0.262 | 0.234 | 1.0 | 0.851 | View |
| | BETA-ALANINE METABOLISM | 13 | 1 | 0.284 | 0.251 | 1.0 | 0.873 | View |
| | SPHINGOLIPID METABOLISM | 15 | 1 | 0.327 | 0.284 | 1.0 | 0.908 | View |
| | MITOCHONDRIAL ELECTRON TRANSPORT CHAIN | 15 | 1 | 0.327 | 0.284 | 1.0 | 0.908 | View |

Page: 1 of 2 Go

Next

# The Matched Metabolite Set

# Single Sample Profiling

# Single Sample Profiling (cont.)

# Concentration Comparison

# Concentration Comparison (cont.)



L-Threonine

| Study | Concentration | Reference | Note |
|---|---|---|---|
| Study 1 | 36.2 (10.82 – 61.58) | Shoemaker JD, Elliott WH: Automated screening of urine samples for carbohydrates, organic and amino acids after treatment with urease. J Chromatogr. 1991 Jan 2;562(1-2):125-38. (Pubmed) | |
| Study 2 | 12.7 (4.934 – 20.4) | Doctor's Data | |
| Study 3 | 1 (0.16 – 2.4) | Geigy Scientific Tables, 8th Rev edition, pp. 165-177. Edited by Cornelius Lentner | |
| Study 4 | 4.9 (2.4 – 7.4) | Geigy Scientific Tables, 8th Rev edition, pp. 165-177. Edited by Cornelius Lentner | |
| Study 5 | 16 (7 – 25) | Geigy Scientific Tables, 8th Rev edition, pp. 165-177. Edited by Cornelius Lentner | |
| Study 6 | 18 (8.4 – 27.6) | Geigy Scientific Tables, 8th Rev edition, pp. 165-177. Edited by Cornelius Lentner | |

# Quantitative Enrichment Analysis

# Result

# The Matched Metabolite Set

# Metabolic Pathway Analysis

# Pathway Analysis

- Purpose: to extend and enhance metabolite set enrichment analysis for pathways by
  - Considering the **structures of pathway**
  - Dynamic pathway visualization
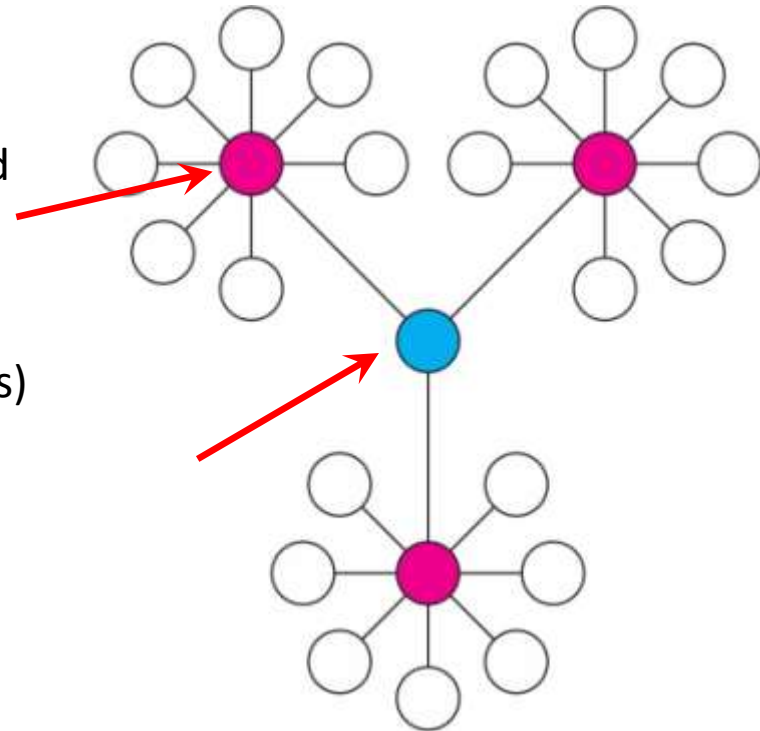- Currently supports ~1500 pathways covering 17 organisms (based on KEGG)

# Network Topology

▶ Which positions are important?

  ▶ Hubs

    ▶ Nodes that are highly connected (red ones)

  ▶ Bottlenecks

    ▶ Nodes on many shortest paths between other nodes (blue ones)

▶ Graph theory

  ▶ Degree centrality

  ▶ Betweenness centrality

Junker *et al. BMC Bioinformatics* 2006

# Data Upload

# Pathway Libraries

# Network Topology Analysis

# Pathway Visualization

# Pathway Visualization (cont.)

# Result



| Pathway Name | Total | Hits | p | -log(p) | Holm p | FDR | Impact | Details |
|---|---|---|---|---|---|---|---|---|
| Glycine, serine and threonine metabolism | 48 | 9 | 1.7267E-4 | 3.7628 | 0.0088061 | 0.0044709 | 0.48394 | KEGG SMP |
| Valine, leucine and isoleucine biosynthesis | 27 | 5 | 3.637E-4 | 3.4393 | 0.018185 | 0.0044709 | 0.06148 | KEGG SMP |
| Methane metabolism | 34 | 6 | 3.8485E-4 | 3.4147 | 0.018858 | 0.0044709 | 0.16466 | KEGG |
| Sulfur metabolism | 18 | 2 | 4.755E-4 | 3.3229 | 0.022824 | 0.0044709 | 0.03307 | KEGG SMP |
| Valine, leucine and isoleucine degradation | 40 | 3 | 6.5285E-4 | 3.1852 | 0.030684 | 0.0044709 | 0.02232 | KEGG SMP |
| Arginine and proline metabolism | 77 | 6 | 6.578E-4 | 3.1819 | 0.030684 | 0.0044709 | 0.06203 | KEGG SMP |
| Aminoacyl-tRNA biosynthesis | 75 | 12 | 6.9157E-4 | 3.1602 | 0.031121 | 0.0044709 | 0.11268 | KEGG |
| Nicotinate and nicotinamide metabolism | 44 | 5 | 7.0133E-4 | 3.1541 | 0.031121 | 0.0044709 | 0.04113 | KEGG SMP |
| Glutathione metabolism | 38 | 2 | 0.0011587 | 2.936 | 0.049823 | 0.0059801 | 0.0019 | KEGG SMP |
| Propanoate metabolism | 35 | 4 | 0.0013934 | 2.8559 | 0.058523 | 0.0059801 | 0.01603 | KEGG SMP |
| Nitrogen metabolism | 39 | 8 | 0.0013997 | 2.854 | 0.058523 | 0.0059801 | 0.00763 | KEGG SMP |
| Galactose metabolism | 41 | 3 | 0.001486 | 2.828 | 0.059441 | 0.0059801 | 0.01992 | KEGG SMP |
| Taurine and hypotaurine metabolism | 20 | 3 | 0.0015243 | 2.8169 | 0.059449 | 0.0059801 | 0.35252 | KEGG SMP |
| Cyanoamino acid metabolism | 16 | 4 | 0.0016826 | 2.774 | 0.06394 | 0.0061295 | 0.0 | KEGG |
| Tryptophan metabolism | 79 | 1 | 0.0018984 | 2.7216 | 0.070241 | 0.0064103 | 0.10853 | KEGG SMP |
| Phenylalanine, tyrosine and tryptophan biosynthesis | 27 | 2 | 0.0021242 | 2.6728 | 0.076472 | 0.0064103 | 0.00738 | KEGG SMP |
| Inositol phosphate metabolism | 39 | 1 | 0.002215 | 2.6546 | 0.077526 | 0.0064103 | 0.13703 | KEGG SMP |
| Pyruvate metabolism | 32 | 4 | 0.0022624 | 2.6454 | 0.077526 | 0.0064103 | 0.41957 | KEGG SMP |
| Cysteine and methionine metabolism | 56 | 2 | 0.0026796 | 2.5719 | 0.088426 | 0.0071926 | 0.02846 | KEGG SMP SMP |
| Alanine, aspartate and glutamate metabolism | 24 | 6 | 0.0029727 | 2.5268 | 0.095127 | 0.0075805 | 0.25546 | KEGG SMP SMP S |

# Not Everything Was Covered

- Clustering (K-means, SOM)
- Classification (SVM, randomForests)
- Time-series data analysis
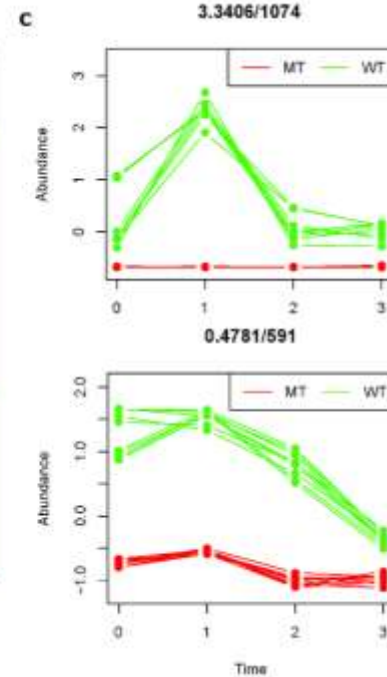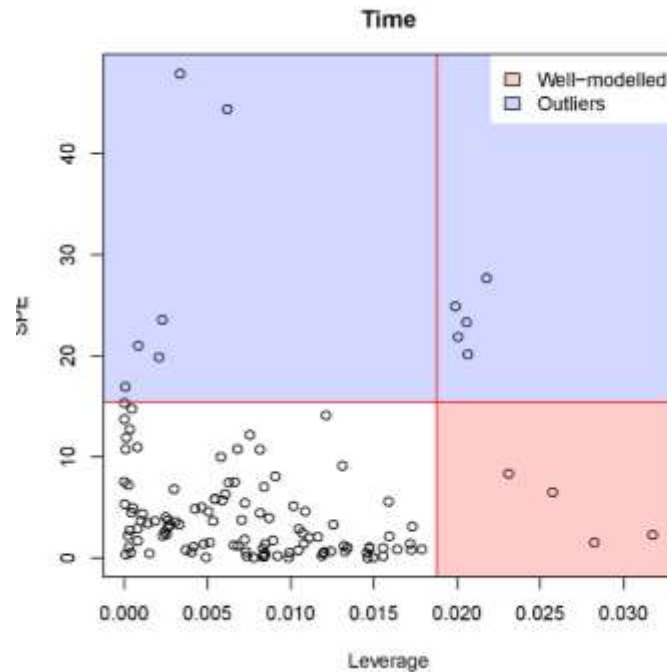- Two -factor data analysis
- Peak searching
- ….

# Two Factor Analysis

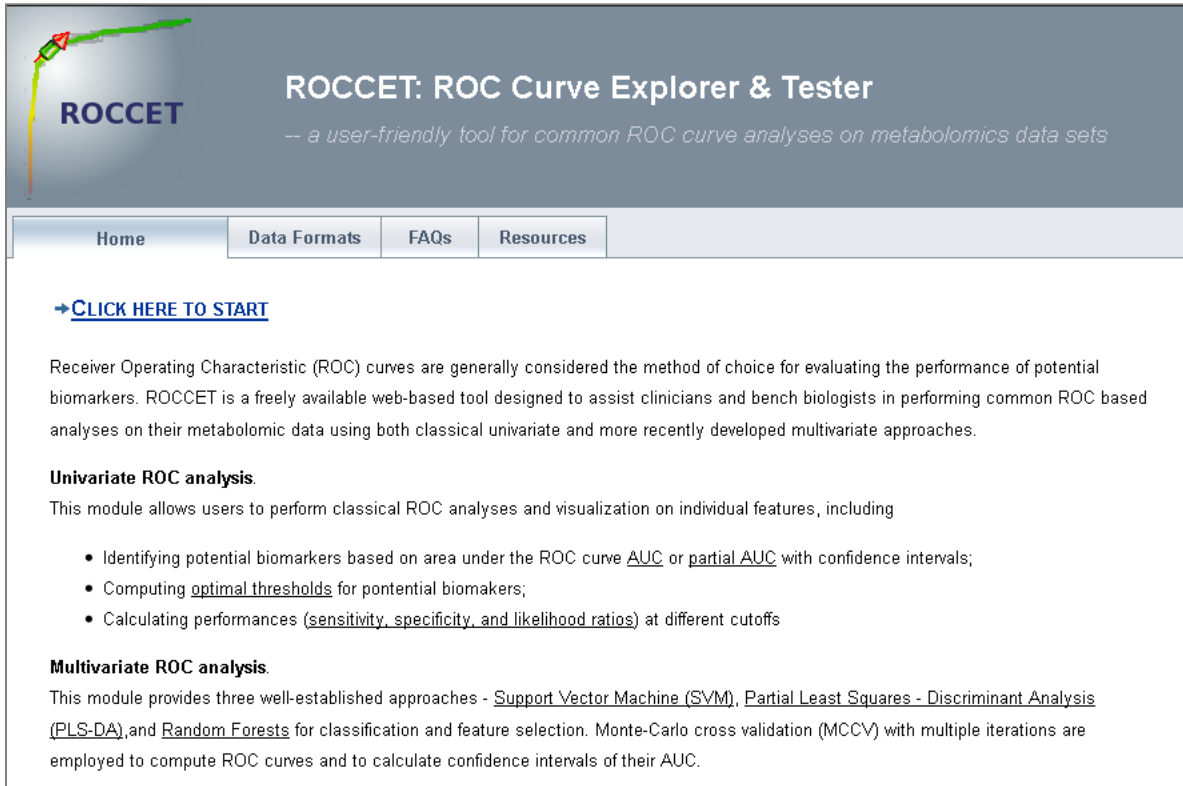- Two – way ANOVA
- Two – way heatmap

# Time series data analysis

- ANOVA-SCA
- Multivariate Empirical Bayes

# Biomarker Discovery & Performance Evaluation

# ROCCET (www.roccet.ca)



Classical ROC → ROC Explorer → ROC Tester
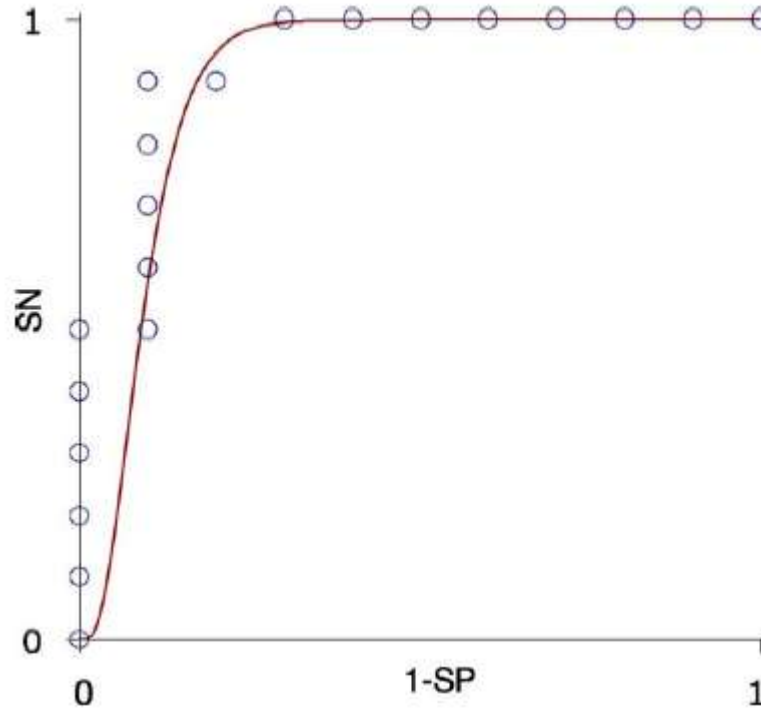
# Sensitivity, Specificity & ROC curve

- Two important performance measures in a diagnostic tests
  - Sensitivity (true positive rate)
  - Specificity (true negative rate)
- Cutoff dependent
  - Increase cutoff, will improve specificity, decease sensitivity
- ROC curves integrate these two measures

# How to construct ROC curves

- Input: a score on a univariate scale
  - A test gives continuous value (i.e. blood *Glucose* level)
  - A classifier that produces a continuous score (i.e. likelihood, probabilities)

# Classical ROC curve

2-h plasma glucose (mmol/L)

| Healthy | Diseased |
|---------|----------|
| 4.86 | |
| 5.69 | |
| 6.01 | |
| 6.06 | |
| 6.27 | |
| 6.37 | |
| 6.55 | |
| 7.29 | 7.29 |
| 7.82 | |
| | 9.22 |
| | 9.79 |
| | 11.28 |
| | 11.83 |
| 12.06 | |
| | 18.48 |
| | 18.50 |
| | 20.49 |
| | 22.66 |
| | 26.01 |



| $1 - SP$ | $SN$ |
|----------|------|
| 1.00 | 1.00 |
| 0.90 | 1.00 |
| 0.80 | 1.00 |
| 0.70 | 1.00 |
| 0.60 | 1.00 |
| 0.50 | 1.00 |
| 0.40 | 1.00 |
| 0.30 | 1.00 |
| 0.20 | 0.90 |
| 0.10 | 0.90 |
| 0.10 | 0.80 |
| 0.10 | 0.70 |
| 0.10 | 0.60 |
| 0.10 | 0.50 |
| 0.00 | 0.50 |
| 0.00 | 0.40 |
| 0.00 | 0.30 |
| 0.00 | 0.20 |
| 0.00 | 0.10 |
| 0.00 | 0.00 |

*TA Lasko, et al (2005)*

# Explore ROC space

- The ROC curve itself (visualization)
- Compare different ROC curves
  - Area under the curve
    - AUC
  - When two curves cross
    - Partial AUC (pAUC)
  - Confidence Intervals
    - Empirical ROC curves are based on samples

# Understand AUC

- Area under an ROC curve (AUC)

  a. The probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one

  a. The average specificity across all values of sensitivity

  b. The average sensitivity across all values of specificity

# ROCCET

- ROC curves – based biomarker discovery and performance evaluation
  - Classical ROC Curve Analysis for individual biomarker
  - Multivariate biomarker model creation & assessment (automatic / manual mode)
    - PLSDA, Linear SVM, Random Forests
  - Calculate AUC & partial AUC with confidence intervals
  - Other supporting utilities

# ROCCET: ROC Curve Explorer & Tester

*-- a user-friendly tool for common ROC curve analyses on metabolomics data sets*

| Home | Data Formats | FAQs | Resources |

→ **CLICK HERE TO START**

Receiver Operating Characteristic (ROC) curves are generally considered the method of choice for evaluating the performance of potential biomarkers. ROCCET is a freely available web-based tool designed to assist clinicians and bench biologists in performing common ROC based analyses on their metabolomic data using both classical univariate and more recently developed multivariate approaches.

**Univariate ROC analysis.**
This module allows users to perform classical ROC analyses and visualization on individual features, including

- Identifying potential biomarkers based on area under the ROC curve AUC or partial AUC with confidence intervals;
- Computing optimal thresholds for pontential biomakers;
- Calculating performances (sensitivity, specificity, and likelihood ratios) at different cutoffs

**Multivariate ROC analysis.**
This module provides three well-established approaches - Support Vector Machine (SVM), Partial Least Squares - Discriminant Analysis (PLS-DA),and Random Forests for classification and feature selection. Monte-Carlo cross validation (MCCV) with multiple iterations are employed to compute ROC curves and to calculate confidence intervals of their AUC.

- **ROC Explorer**
  This purpose of this module is to create and identify robust predictive models using multiple biomarkers. We have integrated feature selection and classfication procedures for the three algorithms mentioned above. The procedures are repeated multiple times in order to identify the best model as well as the most stable features. Various graphical presentations such as ROC Curve View, Probability View, Significant Feature View, etc. are provided to facilitate improved understanding of the results.
- **ROC Tester**
  This module offers flexible interface which allows users to manually construct a biomarker model and to evaluate its performance. It also allows users to allocate a subset of samples as hold-out data for validation (that outside the CV). Other features permutations tests are also available for further model assessments.

# ROCCET: ROC Curve Explorer & Tester
-- *a user-friendly tool for common ROC curve analyses on metabolomics data sets*

**Home**

- Upload
- Data check
- Data Processing
- **Analysis**
  - Univ. ROC
  - ROC Explorer
  - ROC Tester
    - Builder
    - Evaluator
- Download
- Log out

## Data Analysis Options

**Choose two target groups of interest** (for group number > 2)

Select the two groups you want to compare    0 vs. 1

**Choose an analysis path:**

⦿ **To perform classical univariate ROC curve analyses**

Perform classical univariate ROC curve analyses, such as to generate ROC curve, to calculate AUC or partial AUC as well as their 95% confidence intervals, to compute optimal cutoffs for any given feature, as well as to generate performance tables for sensitivity, specificity, and confidence intervals at different cutoffs.

○ **To perform automated biomarker selection and model evaluation (ROC Explorer)**

Perform automated biomarker selection and classification using one of the three multivariate algorithms - support vector machines (SVM), partial least squares discriminant analysis (PLS-DA), and random rorests.

○ **To create and evaluate custom biomarker models (ROC Tester)**

Manually select potential biomarker(s) and then test their performance using any of the three algorithms mentioned above. The module also allows users to hold out a subset of samples for validation purpose (i.e. outside the buildin cross validation). Users can also assess the importance of a model using permutation-based approaches.

**Submit**

# AUC, pAUC & CI



| Var. | AUC | CI |
|------|-----|-----|
| 2 | 0.978 | 0.975-0.981 |
| 3 | 0.978 | 0.974-0.981 |
| 5 | 0.975 | 0.967-0.983 |
| 10 | 0.948 | 0.936-0.96 |
| 20 | 0.916 | 0.902-0.931 |
| 39 | 0.908 | 0.894-0.923 |

Partial area under the curve (pAUC) = 0.169
95% CI: 0.16-0.178

# Posterior probabilities

# Accuracies

Metabolomics 2012

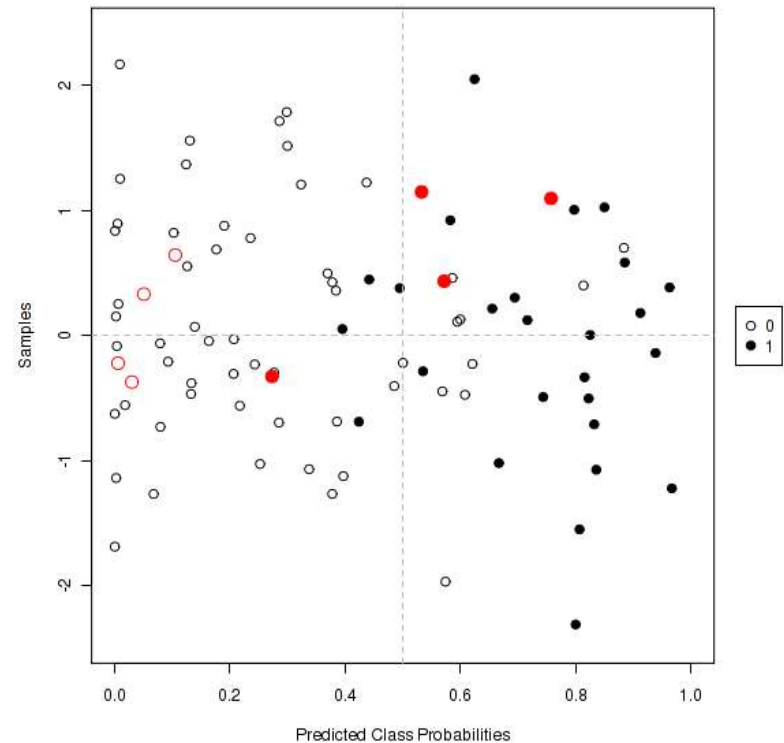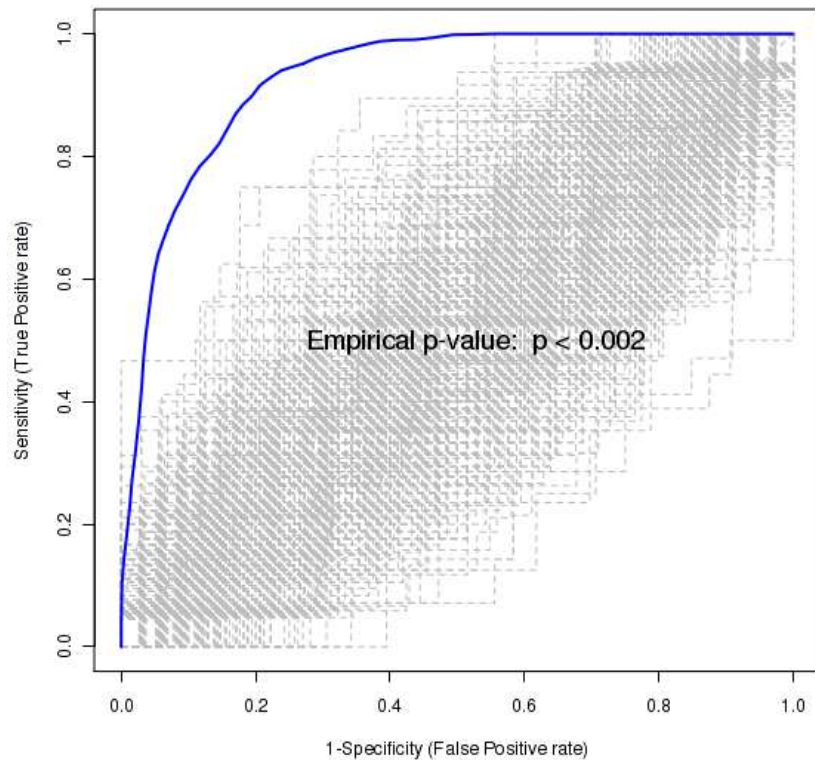# ROC & Posterior probabilities (with hold-out)
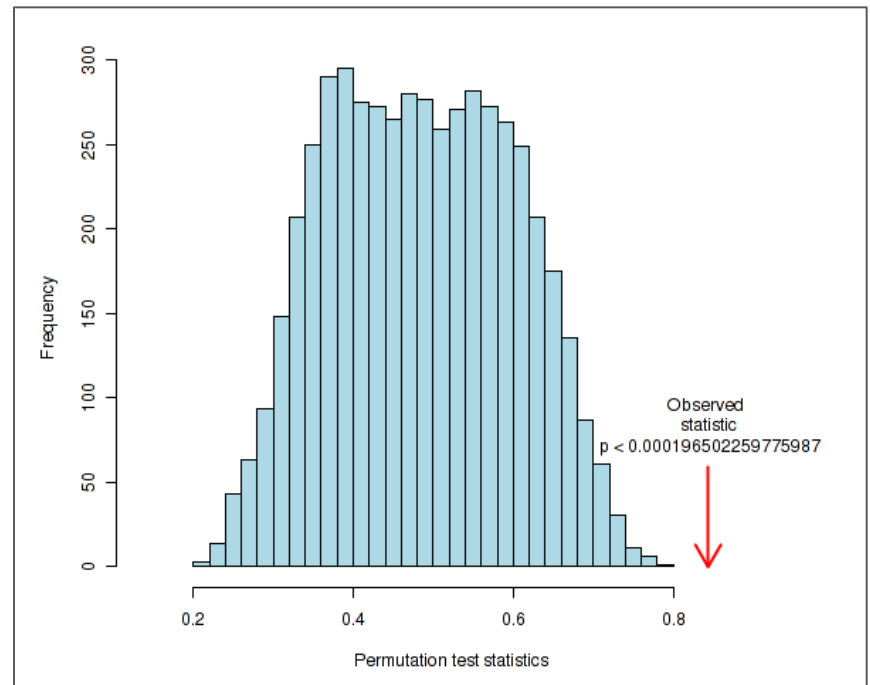
- AUC = 1

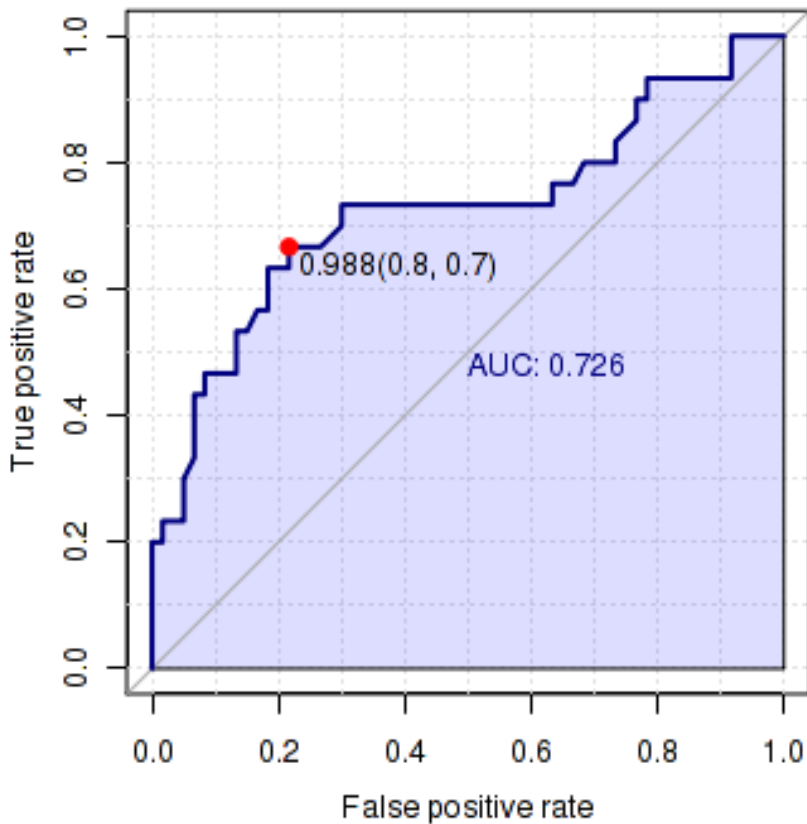- Accuracy = 7/8

# Permutations

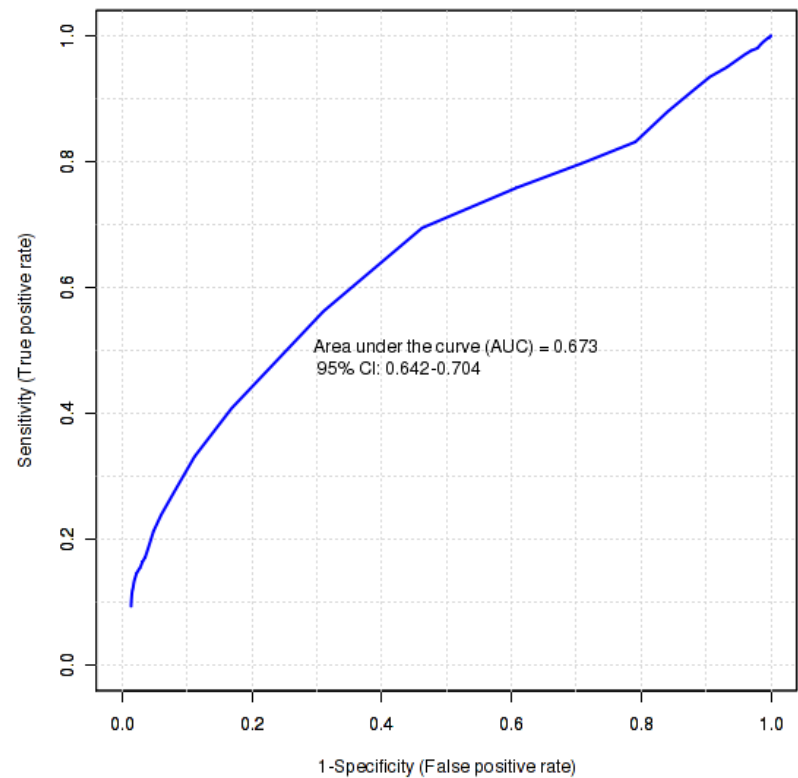**Based on AUC**

**Based on accuracy**

# Over-estimation

**Classical : 0.726**



**CV-based: 0.673**

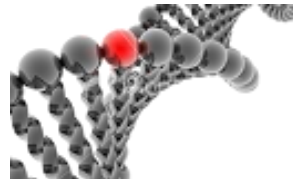# Some Technical Details (1)

- Calculate AUC
  - Empirical or non-parametric method
    - Connecting data points with straight lines
    - Trapezoid rules
- Calculate CI
  - Bootstrapping  (classical univariate)
  - Repeated random sampling & cross validation

# Some Technical Details (2)

- Biomarker selection
  - Classical univariate
    - AUC/pAUC
  - Multivariate MCCV-based
    1. Feature selection
       - PLSDA (VIP score)
       - RandomForest (mean decrease accuracy)
       - Linear SVM (feature weight)
    2. Model Selection
       - AUC/pAUC

# Acknowledgements

- Dr. David Wishart

- Dr. David Broadhurst

- Dr. Rupa Mandal

- The Metabolomic Innovation Center (TMIC)

- University of Alberta, Canada