THE UNIVERSITY OF
ALABAMA AT BIRMINGHAM
Knowledge that will change your world

GBS 724 02-08-16

# Applying MetaboAnalyst

## Xiangqin Cui, PhD

**Select MS peak list option and then load the .zip file**

## Processing MS peak list data :

Peaks need to be matched across samples in order to be compared. For two-column format (mass and intensities), peaks are grouped by their m/z values. For three column data (mass, retention time, and intensities), the program will further group peaks based on their retention time. Users need to supply tolerance values in order to proceed. Here are some suggested values: mass tolerance - 0.25 (m/z); retention time - 30 (seconds) for LC-MS peak, and 5 (seconds) for GC-MS peaks. Please note, If a sample has more than one peak in a group, they will be replaced by their sum; some groups will be excluded if none of the classes has at least half its samples represented. Finally, the program create a peak intensity table in which each sample occupies a row and each column represents a peak group identified by the median values of its position (m/z and/or retention time).

**Mass tolerance (m/z):**     0.025

**Retention time tolerance:**     30.0     Submit

---

#### MS peak processing information

The uploaded files are peak lists and intensities data.

A total of 6 samples were found.

These samples contain a total of 14304 peaks.

with an average of 2384 peaks per sample

A total of 2346 peak groups were formed.

Peaks of the same group were summed if they are from one sample.

Peaks appear in less than half of samples in each group were ignored.

**Data Integrity Check:**

1. Checking the class labels - at least three replicates are required in each class.
2. If the samples are paired, the pair labels must conform to the specified format.
3. The data (except class labels) must not contain non-numeric values.
4. The presence of missing values or features with constant values (i.e. all zeros)

**Data processing information:**

Checking data content ...passed

The uploaded files are peak lists and intensities data.

A total of 6 samples were found.

These samples contain a total of 14304 peaks.

with an average of 2384 peaks per sample

2 groups were detected in samples.

Samples are not paired.

All data values are numeric.

A total of 0 (0%) missing values were detected.

By default, these values will be replaced by a small value.

Click **Skip** button if you accept the default practice

Or click **Missing value imputation** to use other methods

| Missing value estimation | Skip |

**Note that XCMSonline filled in peaks**

---

Non-informative variables can be characterized in two groups: variables of very small values - these variables can be detected using mean or median; variables that are near-constant throughout the experiment conditions - these variables can be detected using standard deviation (SD); or the robust estimate such as interquantile range (IQR). The relative standard deviation(RSD = SD/mean) is another useful variance measure independent of the mean. The following empirical rules are applied during data filtering:

- **Less than 250 variables**: 5% will be filtered;
- **Between 250 - 500 variables**: 10% will be filtered;
- **Between 500 - 1000 variables**: 25% will be filtered;
- **Over 1000 variables**: 40% will be filtered;

Please note, in order to reduce the computational burden to the server, the **None** option is only for less than 2000 features. Over that, if you choose None, the IQR filter will still be applied. In addition, 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.

- ● Interquantile range (IQR)
- ○ Standard deviation (SD)
- ○ Median absolute deviation (MAD)
- ○ Relative standard deviation (RSD = SD/mean)
- ○ Non-parametric relative standard deviation (MAD/median)
- ○ Mean intensity value
- ○ Median intensity value
- ○ None (less than 2000 features)

| Process |

**Sample normalization**

- ○ None
- ○ Sample specific normalization (i.e. dry weight, volume)    <u>Click here to specify</u>
- ● Normalization by sum
- ○ Normalization by median
- ○ Normalization by reference sample
  - ● Specify a reference sample    | posmode_ir1 ▼ |
  - ○ Create a pooled average sample from group | Posmode_IR ▼ |
- ○ Normalization by reference feature | 50.0177/11.02 ▼ |

# Data options before stats analysis

**Data transformation**

- ● None
- ○ Log transformation        (generalized logarithm transformation or glog)
- ○ Cube root transformation (take cube root of data values)

**Data scaling**

- ○ None
- ○ Auto scaling    (mean-centered and divided by the standard deviation of each variable)
- ● Pareto scaling (mean-centered and divided by the square root of standard deviation of each variable)
- ○ Range scaling (mean-centered and divided by the range of each variable)

Now (on the right) the data are mean-centered and variance is reduced by dividing by the standard error of the mean

# Effect of normalization, mean centering and Pareto scaling



Submit

**Univariate Analysis**

Fold Change Analysis   T-tests   Volcano plot

One-way Analysis of Variance (ANOVA)

Correlation Analysis   Pattern Searching

**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   Heatmaps

Partitional Clustering:   K-means   Self Organizing Map (SOM)

**Classification & Feature Selection**

Random Forest

Support Vector Machine (SVM)

**Statistical methods available to process the data on MetaboAnalyst**

**Today we'll focus on univariate analysis (Volcano plots) and multivariate analysis (PCA and PLS-DA)**

# Fold-change  - pink = >1.5

T-Tests

Univariate analysis – Volcano plot

- **This is a useful measure**
- **The data points represent comparisons made one metabolite at a time**
- **The plot identifies metabolites with p-values <0.01 and a fold change >1.5**

# Volcano plot set up

**Volcano Plot**

The volcano plot is a combination of fold change and t-tests. Note, for unpaired samples, the x-axis is log (FC). For paired analysis, the x-axis is number of significant counts. Y-axis is -log10(p.value) for both cases.

**Analysis type:** Unpaired

**X-axis:**
Fold change threshold: 1.5
Comparison type: Grubbs_neg_1/Gr
Sig. count threshold (paired): 75.0 %
Non-parametric tests: ☐

**Y-axis:**
P value threshold: 0.01
Group variance: Equal

Submit

---

## Volcano plot with fold change=1.5 and p <0.01

Top 25 peaks (mz/rt) correlated with the 1-2



493.1546/14.2
361.2886/18.32
433.9209/11.47
309.2782/21.03
229.1536/9.46
163.0377/12.03
230.1564/9.46
492.1544/14.2
228.0514/9.18
384.2256/17.05
577.2316/7.59
439.2363/12.45
157.0843/13.98
95.0827/15.38
225.1117/17.6
384.2115/17.05
378.2112/15.8
438.2379/12.45
251.1335/12.01
345.2259/17.75
229.14155/17.61
973.4989/17.8
655.2175/6.18
309.8429/8
243.1214/17.63

-1.0     -0.5      0.0      0.5      1.0

Correlation coefficients

# Multivariate Analysis

- **In this type of analysis, principal components are identified.**
- **Contributions from all the metabolites to each principal component are calculated, reducing >1000 values to one number for each component**
- **When this done in an unsupervised way, the group information is added after principal component analysis (PCA)**
- **Partial least squares discriminant analysis (PLS-DA) is a supervised procedure and group information is included in the analysis**

# 2D- and 3D-PCA plots

# 2D- and 3D-PLS-DA plots



# From download file – plsda_vip.csv file

**Questions?**