

# BIOCRATES

LIFE SCIENCES

The Deep Phenotyping Company

# MetaboAnalyst Tutorial

## Analysis of Met/IDQ™ Data in MetaboAnalyst

# MetaboAnalyst Tutorial

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Standard Data Format</b>	<b>7</b>
<b>3</b>	<b>Data Cleaning</b>	<b>7</b>
<b>4</b>	<b>Data Pretreatment</b>	<b>9</b>
<b>5</b>	<b>Data Analysis Using MetaboAnalyst</b>	<b>10</b>
5.1	Exporting Data for MetaboAnalyst	10
5.2	MetaboAnalyst Data Analysis Workflow	14
5.2.1	Data Upload	15
5.2.2	Missing Value Estimation	17
5.2.3	Data Filtering	18
5.2.4	Data Normalization	18
5.2.5	Univariate Analysis of Data Sets with Two Groups	20
5.2.5.1	Fold Change Analysis	21
5.2.5.2	T-Test	22
5.2.5.3	Volcano Plot	23
5.2.6	Univariate Analysis of Data Sets with More Than Two Groups	24
5.2.7	Multivariate Analysis	25
5.2.7.1	Principal Component Analysis (PCA)	26
5.2.7.2	Partial Least Squares-Discriminant Analysis (PLS-DA)	30
5.2.7.3	Heatmap	34
5.2.7.4	Pattern Hunter	36
5.2.7.5	Download Results	37
5.3	MetaboAnalyst Biomarker Analysis	38
5.3.1.1	Data Upload	39
5.3.1.2	Data Normalization	39

5.3.1.3	ROC Analysis .....	41
5.3.1.4	Download Results .....	42
	Ordering and Technical Support .....	43

The information in this manual is subject to change without notice and should not be construed as a commitment by BIOCRATES® Life Sciences AG to assume responsibility for any errors that may appear. While every precaution has been taken in the preparation of this manual, BIOCRATES® Life Sciences AG shall not be liable for punitive, incidental, or consequential damage in connection with or arising from the use of this manual.

Met/DQ™ (hereinafter referred to as Met/DQ) is a trademark of BIOCRATES® Life Sciences AG. All other trademarks are the sole property of their respective owners. It is not intended to encourage use of these products in any manner that might infringe on the intellectual property rights of others.

Copyright © BIOCRATES Life Sciences AG 2020. All rights reserved. BIOCRATES® is a registered trademark of BIOCRATES® Life Sciences AG.

**For Research Only. Not for use in diagnostic procedures.**



BIOCRATES Life Sciences AG

Eduard-Bodem-Gasse 8

A-6020 Innsbruck

Austria

Phone: +43 (0)512 579 823

Fax: +43 (0)512 579 823 329

[office@biocrates.com](mailto:office@biocrates.com)

[biocrates.com](http://biocrates.com)

Document version: 1

Edition: NA

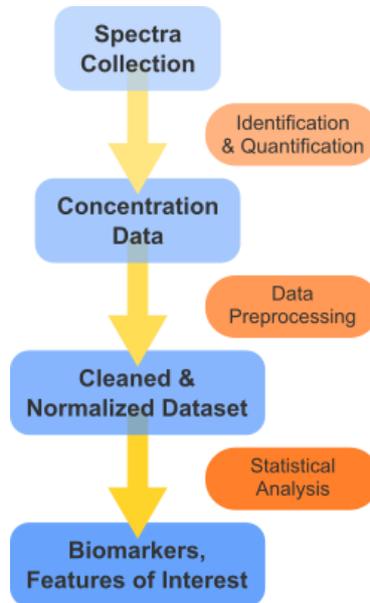
Date of the last revision: 2020-03-10

Filename: MetaboAnalyst-Guide-1.pdf

**MetaboAnalyst Tutorial**

## 1 Introduction

Obtaining biologically relevant information from large and often complex metabolomics data sets is a major challenge faced by scientists. Following metabolite quantification, the data must be cleaned before identifying biomarkers or features of interest. While data preprocessing and pretreatment are critical parts of metabolomic data analysis, there are no particularly straightforward rules that guide investigators to select the best preprocessing and pretreatment options. Figure 1 represents a data processing workflow for targeted metabolomics, starting with analyte measurement and ending with the identification of possible biomarkers. Before statistical analysis can take place, the data must first be cleaned and normalized.



**Figure 1: Data processing workflow for targeted metabolomics data.**

MetaboAnalyst (<https://www.metaboanalyst.ca/>) provides a ready-to-use framework for data cleaning, preprocessing, and basic statistical analysis for metabolomics data sets. Furthermore, Biocrates' kit workflow software, MetIDQ, allows for easy export of metabolomics data sets to the data format accepted by MetaboAnalyst. In order to utilize this functionality, a statistical group must be assigned to each sample in the exported data set. Data should then be exported as a MetaboAnalyst .csv file (refer to MetIDQ user manual).

For training purposes, two data sets are provided with this guide to demonstrate analysis involving data sets with two or more than two groups:

MetIDQ 2 group dataset.csv

MetIDQ 4 group dataset.csv

For more information on MetaboAnalyst, please refer to the following:

<https://www.metaboanalyst.ca/docs/Overview.xhtml>

<https://www.metaboanalyst.ca/docs/Faqs.xhtml>

<https://www.metaboanalyst.ca/docs/Tutorial.xhtml>

## 2 Standard Data Format

The standard data format used when uploading to MetaboAnalyst is a data frame with sample names in the first column, group labels in the second column, and the metabolite concentrations for each sample in the following columns (see Table 1).

**Table 1: Standard MetaboAnalyst data table format.**

<i>Sample</i>	<i>Group</i>	<i>Met 1</i>	<i>Met 2</i>	<i>Met 3</i>	<i>...</i>	<i>Met N</i>
S1	A	99.380	10.177	51.484	...	71.882
S2	A	101.195	10.786	50.446	...	73.318
S3	A	102.165	9.375	49.668	..	72.056
S4	B	99.481	8.291	48.111	...	73.282
S5	B	101.282	10.867	50.209	...	73.572
S6	B	99.430	9.950	47.602	...	71.983

## 3 Data Cleaning

Data preprocessing or cleaning is generally performed for analytes that are below the limit of detection (LOD), missing for a particular experimental group, or show poor reproducibility. This helps to remove clutter from the data set that may obscure interpretation.

Some suggested data cleaning rules are described below:

### **Remove metabolites with missing or < LOD values (“80% rule”)**

Metabolites that are found with valid measurements above LOD in at least 80% of the samples in at least one experimental group can be used for statistical analysis. If this criterion is met, missing values can be imputed with a fraction of the LOD (see **Missing Value Imputation** below). Metabolites that do not meet this criterion should be removed before statistical analysis.

### **Remove metabolites with poor repeatability (optional)**

Metabolites that demonstrate poor reproducibility can also be removed from analysis. Generally, as calculated for the quality control (QC) samples (or other replicated samples), metabolites with a coefficient of variation above 30% should be removed.

### **Missing value imputation (optional)**

Missing values often occur if a particular analyte is present in one sample or group, but below LOD in others. The major problem of such low reported values lies in further statistical analysis of the data. The calculation of statistical values (mean value, standard deviation ...) becomes problematic as low values are truncated.

Missing value imputation can be used to replace missing or below LOD values with non-zero values while maintaining the data structure.

The most common and easiest strategy is simple replacement, where  $< \text{LOD}$  values are replaced with zero (not recommended), with some fraction of the detection limit (usually either  $1/2$  or  $1/\sqrt{2}$ ), or with the detection limit itself. Other small value replacement techniques include constant small value, median, k-nearest neighbour, random forest, logspline, among others.

### **Batch effect removal (QC normalization)**

Experiments conducted in batches or run across multiple Kit plates can introduce systematic non-biological variations in metabolomics measurements that occur between batches. This is referred to as "batch effect". Multiple approaches for batch effect removal have been published. Ratio-based methods scale the concentration of each analyte in each sample based on a set of reference samples in each batch. As several quality control samples are recommended to be used within our kits, this is the preferred approach.

## 4 Data Pretreatment

The selection of a pretreatment method is an essential step in metabolomics data analysis that greatly affects the final metabolite list and results. Every method has its own merits and drawbacks, thus the best-suited method must be chosen based on the biological question to be answered, the properties of the data set, and the data analysis method selected

**Scaling methods** aim to adjust for fold differences between variables by dividing each variable by a scaling factor. This factor is specific to each variable and can be a descriptor of data dispersion such as the standard deviation (e.g. in autoscaling and Pareto scaling) or of size measure such as the mean or median.

Autoscaling brings all standard deviations to one, thus allowing the analytes to be studied on the basis of correlations rather than covariance. In contrast, Pareto scaling uses the square root of the standard deviation as scaling factor. This causes large fold changes to decrease more than small fold changes and reduces the difference between them compared to unprocessed data. While scaling can be useful to ensure highly abundant metabolites (e.g. hexoses) and metabolites at very low concentrations (e.g. biogenic amines) retain equal importance, it should be used with caution to avoid causing problems in the data set.

**Non-linear transformations** are commonly used to correct for heteroscedasticity or to make skewed distributions more symmetric. For pretreatment of kit data,  $\log_2$  transformation is recommended. This transformation is commonly applied to meet assumptions of statistical tests (e.g. symmetric distribution of data, correction for heteroscedasticity and skewness of data) and to improve the interpretability and visualization.

A detailed description of the different scaling methods, transformation types and their advantages and disadvantages was described by van der Berg, et al. (doi: [10.1186/1471-2164-7-142](https://doi.org/10.1186/1471-2164-7-142), Table 1).

## 5 Data Analysis Using MetaboAnalyst

The standard data input format for MetaboAnalyst is a data frame with sample names in the first column, group labels in the second column, and the variables in the remaining columns (Table 1). The second column (group) is mandatory, if no group information is available a dummy group label must be used.

For data analysis with MetaboAnalyst:

- Comma Separated Values (.csv)
- Both sample and group names must be unique and not contain any special characters.
- Concentration data should contain only numeric and positive values (using empty or NA for missing values).
- At least two groups are required
- At least three replicates are required in each group

Additional information about the data formats used by MetaboAnalyst can be found here: <http://www.metaboanalyst.ca/MetaboAnalyst/faces/docs/Format.xhtml>

### 5.1 Exporting Data for MetaboAnalyst

When exporting data directly from MetIDQ, make sure to have at least two groups defined for all samples and choose “MetaboAnalyst Comma Separated Values (\*.csv)” when exporting from MetSTAT. Sample barcodes will be used as the sample identification and groups will be selected from the categories defined in MetLIMS.

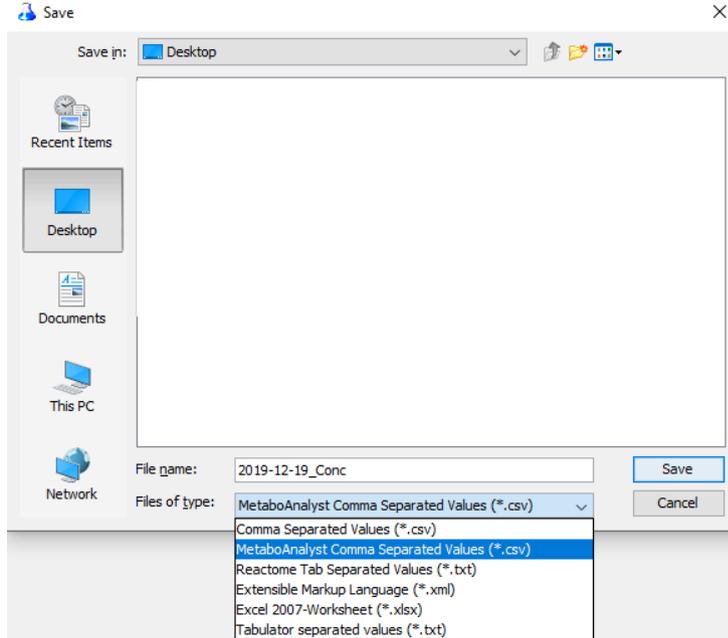
1	<input type="checkbox"/>	...	1020158015	QC Level 1	p180_QC1
2	<input type="checkbox"/>	...	1020158020	QC Level 2	p180_QC2
3	<input type="checkbox"/>	...	1020158034	QC Level 3	p180_QC3

While MetaboAnalyst contains some sample normalization and scaling functions, it is recommended to do all normalization and sample amount scaling (tissue factor, cell number, etc.) in MetIDQ prior to exporting. If QC normalization is used, it would also be required to uncheck the QC samples in MetSTAT to avoid including them in the statistical analysis. Blanks, zero samples, and calibration standards should also be unchecked or excluded from MetSTAT before exporting for MetaboAnalyst.

Under the “Save Options” in MetSTAT, values that are  $< LOD$  should also be replaced with NA, otherwise, since MetaboAnalyst has no way of knowing the LOD of the samples, it will treat these numbers as real concentration values. Other concentration values that have an undesirable status can also be replaced at this point.

Repl...	Status	with
<input type="checkbox"/>	STD/QC < Limit	STD/QC < Limit
<input checked="" type="checkbox"/>	< LOD	NA
<input type="checkbox"/>	< LLOQ	< LLOQ
<input type="checkbox"/>	> ULOQ	> ULOQ
<input type="checkbox"/>	Valid	Valid

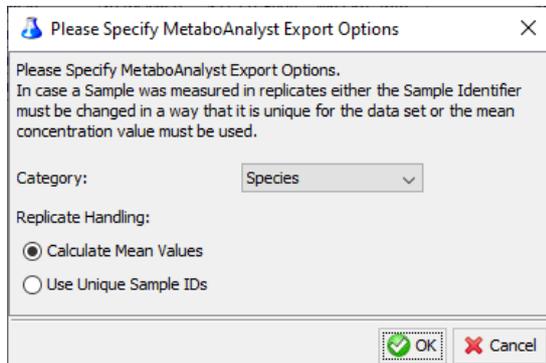
 Export



When selecting to save the file as a MetaboAnalyst CSV, MetIDQ will ask which category should be used for the statistical analysis. It is required that all samples included in the export should be assigned a group from the chosen category.

MetIDQ will also ask how replicated samples should be handled in the export. In general, it is recommended to average replicate injections, unless there are experimental design reasons to keep them separated in the analysis.

Selecting “Calculate Mean Values” will average all sample replicates containing the same barcode. Selecting “Use Unique Sample IDs” will keep each sample replicate as an individual entry.



The exported MetaboAnalyst CSVs can now be uploaded directly to MetaboAnalyst.

**Note on MetaboINDICATOR™ (formerly RatioExplorer):**

If Metabolism Indicators are activated in MetSTAT, they will be exported alongside the analyte concentrations in the MetaboAnalyst output.

It is recommended to disable the indicators when performing multivariate analyses as the combination of metabolite concentrations with ratios and sums of concentrations will bias or confound the analysis. For univariate analysis, it is up to the user whether to include or disable these indicators. Variation in concentrations as well as the changes in the sums and ratios of analytes may provide some insight into data analysis.

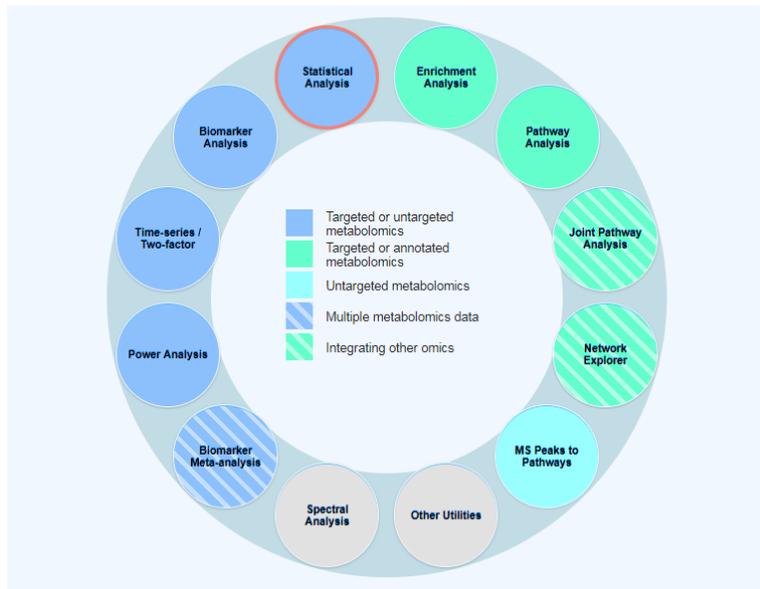
## 5.2 MetaboAnalyst Data Analysis Workflow

To begin a new statistical analysis, navigate to the MetaboAnalyst Home page (<https://www.metaboanalyst.ca/>) and use the link “>> **click here to start** <<” to enter the module selection page.

MetaboAnalyst - statistical, functional and integrative analysis of metabolomics data

Welcome >> [click here to start](#) <<

Select the “**Statistical Analysis**” module from the roulette:



## 5.2.1 Data Upload

In the Statistical Analysis module, use the “**1) Upload your data**” panel, and select the options as indicated below (default parameters). Choose the MetaboAnalyst CSV file created by MetIDQ and then click the “Submit” button.

Tab-delimited text (.txt) or comma-separated values (.csv) file:

Data Type:  Concentrations  Spectral bins  Peak intensity table

Format:

Data File:  No file chosen

Zipped Files (.zip) :

Data Type:  NMR peak list  MS peak list  MS spectra

Data File:  No file chosen

Pair File:  No file chosen

After uploading the data, an integrity check is performed automatically and the result is shown below. If missing values have been detected, select “Missing value estimation” otherwise click “Skip” to go to the normalization step.

**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

Samples are in rows and features in columns

The uploaded file is in comma separated values (.csv) format.

The uploaded data file contains 40 (samples) by 263 (compounds) data matrix.

Samples are not paired.

4 groups were detected in samples.

Only English letters, numbers, underscore, hyphen and forward slash (/) are allowed.

Other special characters or punctuations (if any) will be stripped off.

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

## 5.2.2 Missing Value Estimation

For removal of analytes with too many missing values, set the cut-off to 20% (80% rule). For statistical analysis missing values need to be replaced by a number to avoid skewing the results. It is recommended to select “Estimate missing values using KNN”. K-Nearest Neighbor (KNN) attempts to replace the missing values with logical numbers based on the other surrounding values within the experimental group. Using this method, other missing value statuses such as “no intercept” are not considered below LOD. Press “Process” for the next step.

**Step 1. Remove features with too many missing values**

Remove features with >  % missing values

**Step 2. Estimate the remaining missing values**

Replace by a small value (half of the minimum positive value in the original data)

Exclude variables with missing values

Replace by column (feature)  ▼

Estimate missing values using  ▼

### 5.2.3 Data Filtering

The data sets generated with Biocrates' kits do not contain more than 5000 features (i.e. metabolites). For Data filtering, select “None” and click “Proceed”.

Filtering features if their RSDs are >  25 % in QC samples

None (less than 5000 features)

Interquartile 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

### 5.2.4 Data Normalization

In general, most large metabolomics data sets are skewed. In order to correct for heteroscedasticity, select the Log transformation option. This will move the data set closer to normal distribution, which is a prerequisite for parametric statistical analysis.

Scaling is a more complex topic that is outside the scope of this document. In general, the choice of scaling depends on the hypothesis and setup of the experiment. While scaling may be beneficial for some data sets, it also has the possibility to cause problems when applied incorrectly. For more information on scaling techniques and their applications, it is recommended to review the earlier mentioned paper by van den Berg, et al. When in doubt, the best option would be to select “None” under the Data Scaling option.

**Sample Normalization**

None

Sample-specific normalization (i.e. weight, volume) [Specify](#)

Normalization by sum

Normalization by median

Normalization by reference sample (PQN) [Specify](#)

Normalization by a pooled sample from group [Specify](#)

Normalization by reference feature [Specify](#)

Quantile normalization

**Data transformation**

None

Log transformation (generalized logarithm transformation or glog)

Cube root transformation (takes the cube root of data values)

**Data scaling**

None

Mean centering (mean-centered only)

Auto scaling (mean-centered and divided by the standard deviation of each variable)

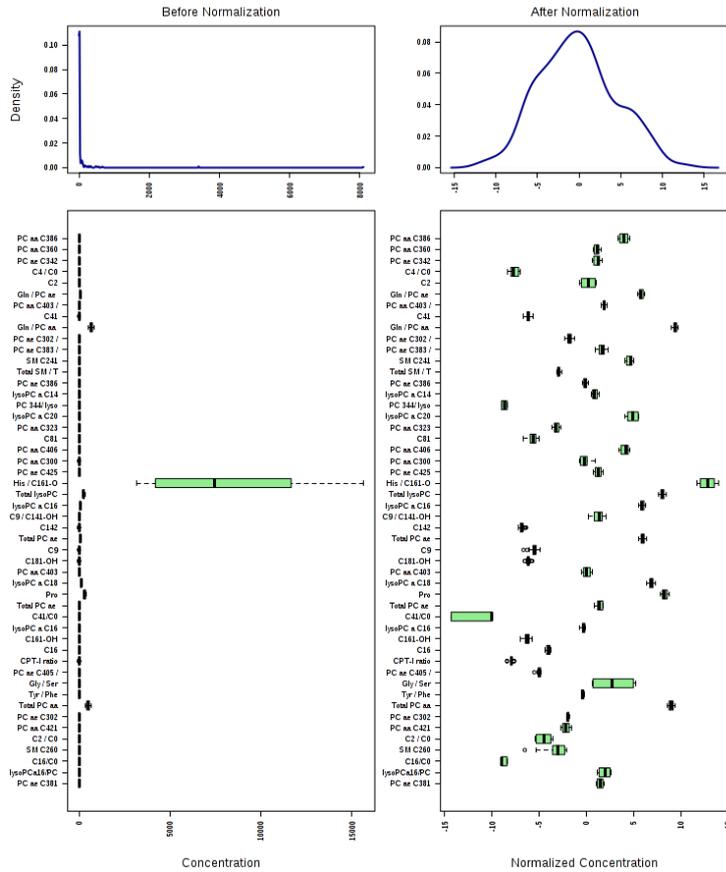
Pareto scaling (mean-centered and divided by the square root of the standard deviation of each variable)

Range scaling (mean-centered and divided by the range of each variable)

[Normalize](#) [View Result](#) [Proceed](#)

The result of the normalization process can be seen by clicking “View Result”. An example data set is shown below. The left side shows the skewed data before normalization and transformed and scaled data on the right side.

Now that data processing and normalization are finished data set is suitable for statistical analysis. Close the normalization results popup and press “Submit” for the next step.



### 5.2.5 Univariate Analysis of Data Sets with Two Groups

Depending on the number of groups in a data set, different methods for data analysis are available and highlighted in blue. For two groups, univariate analysis can be applied such as

fold change analysis, t-test, and volcano plot. If there are more than two groups, the more complex analysis of variance (ANOVA) must be used.

**Univariate Analysis**

[Fold Change Analysis](#) [T-tests](#) [Volcano plot](#)

One-way Analysis of Variance (ANOVA)

[Correlation Analysis](#) [Pattern Searching](#)

**Chemometrics Analysis**

[Principal Component Analysis \(PCA\)](#)

[Partial Least Squares - Discriminant Analysis \(PLS-DA\)](#)

[Sparse Partial Least Squares - Discriminant Analysis \(sPLS-DA\)](#)

[Orthogonal Partial Least Squares - Discriminant Analysis \(orthoPLS-DA\)](#)

**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**

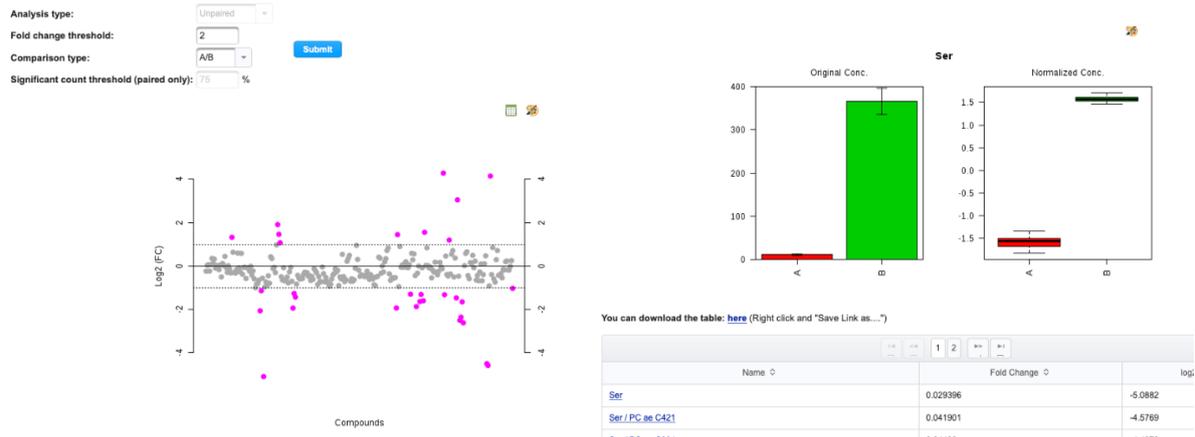
[Random Forest](#)

[Support Vector Machine \(SVM\)](#)

### 5.2.5.1 Fold Change Analysis

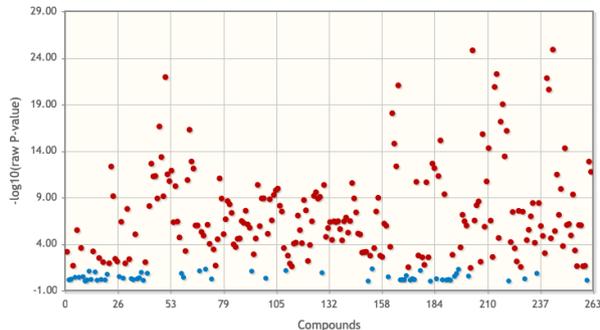
The fold change (FC) analysis compares the  $\log_2$  transformed ratio between each metabolite in each group. If the fold change value exceeds a defined threshold (in this example: 2), the measurement will be flagged as significant (pink colored).

A data table with detailed information is available by clicking on the green table symbol (📄) above the plot. The data can be sorted by name, fold change, and  $\log_2(\text{FC})$ . Clicking on any metabolite name will display the box plot for this metabolite.



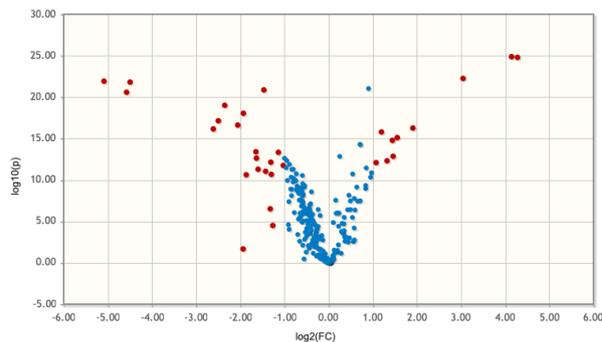
### 5.2.5.2 T-Test

A t-test will determine if the changes between metabolites in each group are considered statistically significant. Before the test is performed, a  $p$ -value threshold (significance level) is chosen. Traditionally 5% (= 0.05; by default) or 1 % (= 0.01). By clicking on the dots, the box plot for this metabolite is shown. A data table with detailed information is available by clicking on the table symbol (📄). Dots colored red indicate a  $p$ -value below the defined threshold (significant).



### 5.2.5.3 Volcano Plot

A volcano plot combines a measure of statistical significance from a statistical test (t-test, y-axis) with the magnitude of the change (fold change, x-axis) enabling quick visual identification of those data-points that display large-magnitude changes that are also statistically significant. Metabolites that fulfil the criteria chosen for fold change and t-test are shown in red. As before, a detailed table can be displayed () and dots can be clicked to show box plots.

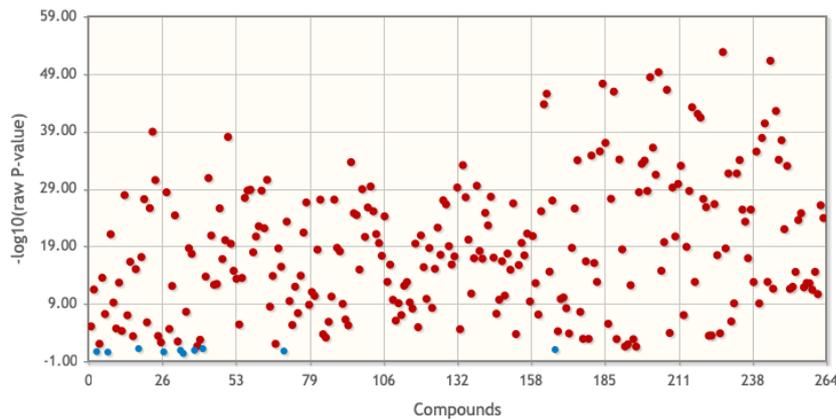


### 5.2.6 Univariate Analysis of Data Sets with More Than Two Groups

If a data set consists of more than two groups, fold change analysis, t-test, and volcano plot cannot be applied. In this case only one-way ANOVA is available.

The ANOVA compares multiple groups and determines if there is a statistically significant difference ( $p$ -value) between each group pair. As before, the  $p$ -value threshold can be set (default = 0.05). Metabolites with a  $p$ -value below the threshold are shown in red (significant).

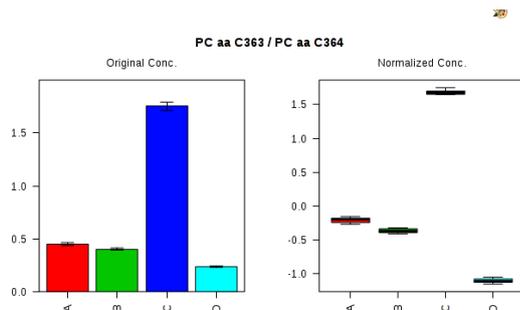
The ANOVA result for a data set with four groups is shown below.



A data table with detailed information is available by clicking on the table symbol . By clicking on the metabolite name, the box plot for this metabolite is shown.

The corresponding table will also provide the false discovery rate (FDR), also referred to as the  $q$ -value. In a statistical evaluation with many comparisons, controlling the FDR becomes important to minimize the risk of a statistically significant result being obtained by pure chance. A metabolite concentration fold change should only be considered statistically

significant if both the  $p$ -value and the FDR are below the specified thresholds. A commonly used FDR threshold is 20% (= 0.2).



You can download the table: [here](#) (Right click and "Save Link as...")

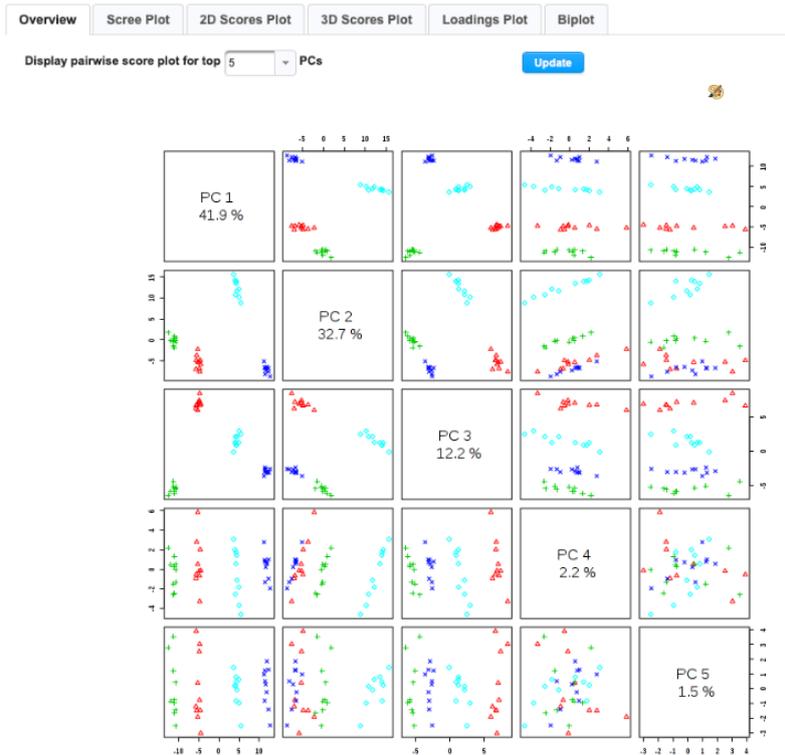
Name	f.value	p.value	$-\log_{10}(p)$	FDR	Post-hoc tests
<a href="#">PC aa C363 / PC aa C364</a>	11227.0	1.5886E-53	52.799	4.1779E-51	A - B; C - A; A - D; C - B; B - D; C - D
<a href="#">Thr / Ser</a>	9252.2	5.1459E-52	51.289	6.7669E-50	A - B; A - C; A - D; B - C; B - D; C - D
<a href="#">Gly / Ser</a>	7159.8	5.159E-50	49.287	4.5227E-48	A - B; A - C; A - D; C - B; D - B; D - C
<a href="#">Gly / Gln</a>	6406.0	3.8079E-49	48.419	2.5037E-47	A - B; A - C; A - D; B - C; B - D; D - C
<a href="#">C181/C0</a>	5550.2	5.0044E-48	47.301	2.6323E-46	D - A; D - B; D - C
<a href="#">lysoPC a C204 / lysoPC a C203</a>	4833.3	5.997E-47	46.222	2.6287E-45	B - A; A - C; D - A; B - C; D - B; D - C
<a href="#">C2 / C0</a>	4645.0	1.2239E-46	45.912	4.5984E-45	A - B; C - A; D - A; C - B; D - B; D - C
<a href="#">C2C3 / C0</a>	4430.3	2.8624E-46	45.543	9.4102E-45	A - B; C - A; D - A; C - B; D - B; D - C
<a href="#">H1</a>	3503.1	1.9349E-44	43.713	5.6543E-43	B - A; C - A; D - A; C - B; D - B; D - C
<a href="#">Om / Ser</a>	3277.5	6.3847E-44	43.195	1.6792E-42	A - B; A - C; A - D; C - B; D - B; C - D

## 5.2.7 Multivariate Analysis

Multivariate analysis can be performed with data sets containing both two or more than two groups.

### 5.2.7.1 Principal Component Analysis (PCA)

Principal components analysis can reveal hidden structure in a complex data set such as groups of observations, trends, multivariate outliers (samples). The data set is transformed to reduce the number of dimensions into principle components. In the first tab, **Overview**, score plots of the first five principal components are shown. The first two principle components will show the greatest amount of differentiation between samples.



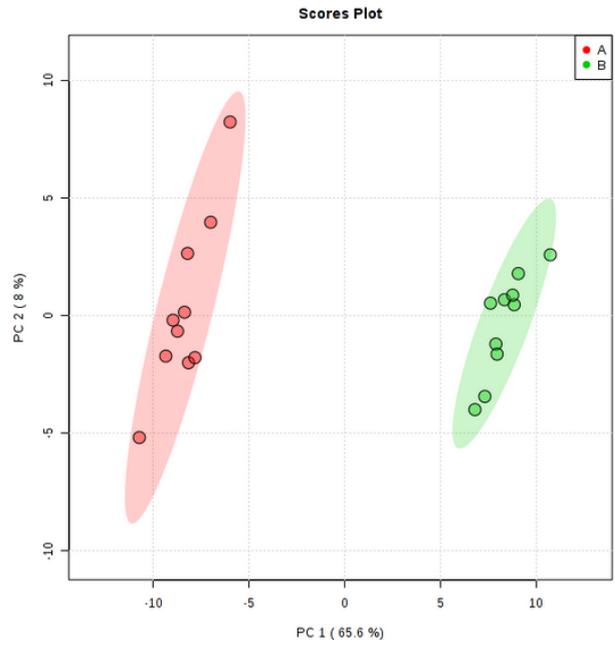
The **Scree Plot** tab displays the contribution of each of the principal components to the total variation in a data set. The green line shows the cumulative variance explained by the first five components, while the red line shows the individual variance for each principal component. Such a plot when read left-to-right can often show a clear separation between the 'most important' components and the 'least important' components. The point of separation is often called the 'elbow' (PC2 in our example).



The **2D Scores Plot** is a powerful tool that enables to visualize how the samples are related to one another. Samples are clustered according to their similarity with more distinct sample groups showing more separation and vice versa. By default, PC1 and PC2 will be plotted with their 95% confidence intervals drawn around each group.

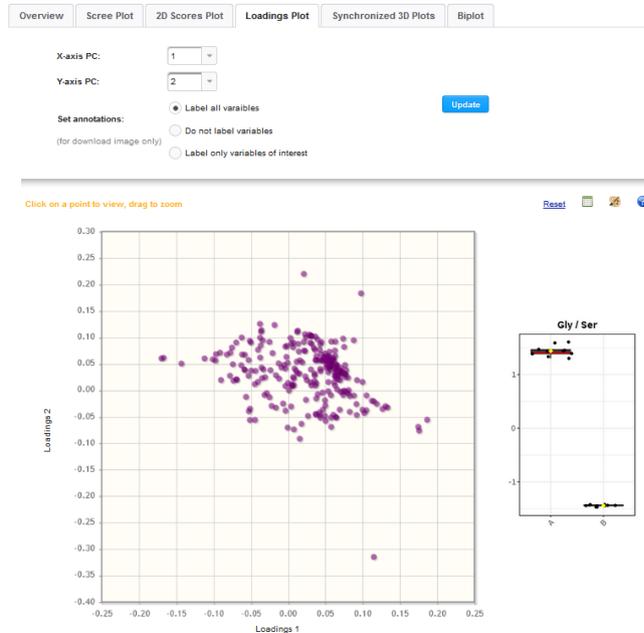
Overview | **Score Plot** | **2D Scores Plot** | Loadings Plot | Synchronized 3D Plots | Biplot

Specify PC on X-axis: 1  
Specify PC on Y-axis: 2  
Display 95% confidence regions:   
Display sample names:   
Use grey-scale colors:   
[Flip Image](#)  X axis  Y axis  All [Update](#)



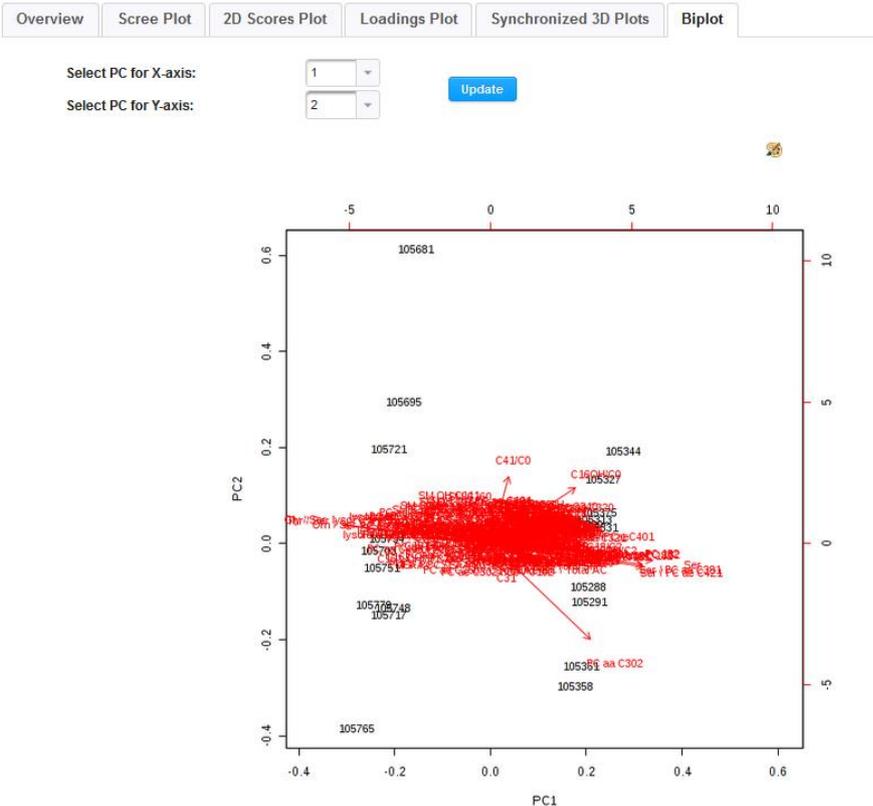
The **Loadings Plot** enables identification of metabolites that are most responsible for driving the patterns seen in the score plot. Metabolites contributing similar information are grouped together indicating correlation. When metabolites are inversely correlated, they are positioned on opposite sides of the plot origin, in diagonally opposed quadrants. The distance to the origin also conveys information. The further away from the plot origin a metabolite is, the stronger its impact.

The figure is interactive. Clicking and dragging in the figure will zoom. Clicking on individual dots will show which metabolites they correspond to in a boxplot to the right.



By comparing the score and loadings plots, the relationships between samples and metabolites can be identified. In most cases, a biplot eases this comparison by combining a score plot and a loadings plot. Samples are displayed as points while metabolites are displayed as vectors.

In most cases, the **Biplot** may become too cluttered to read properly. The data can be individually extracted from the tables in the loadings and score plots.

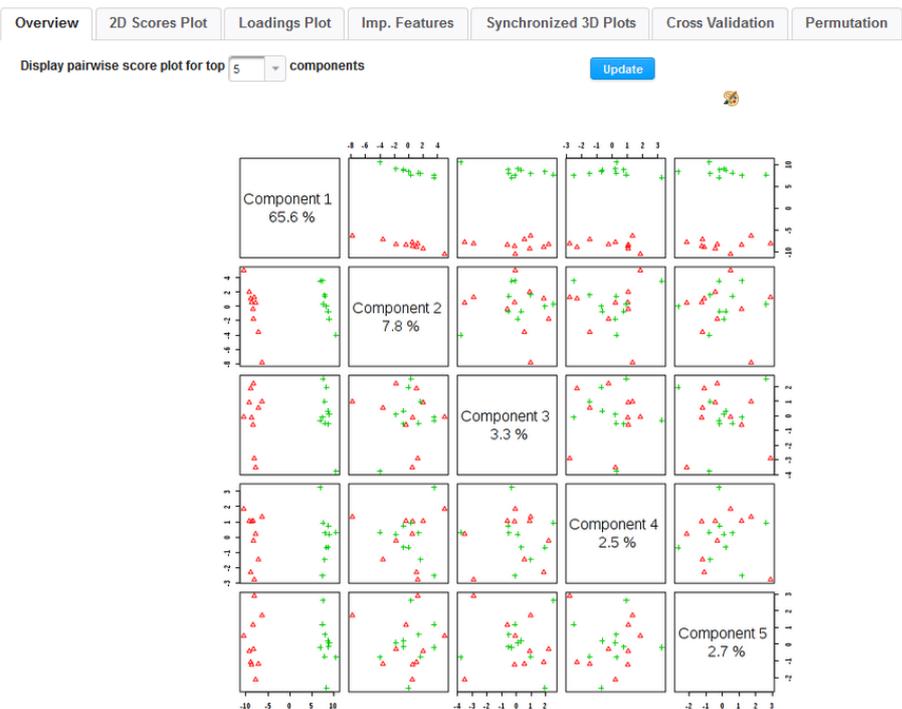


### 5.2.7.2 Partial Least Squares-Discriminant Analysis (PLS-DA)

Partial least square-discriminant analysis (PLS-DA) is performed in order to sharpen the separation between groups of samples. PLS-DA can perform both classification and feature selection. The algorithm uses cross validation to select an optimal number of components for classification. As a supervised technique, PLS-DA is prone to overfitting (it will try to separate classes even when there is no real difference between them), therefore it is

important to validate the outcomes of PLS-DA analyses. Cross validation procedures or permutation testing are commonly used for this purpose. Generally, it is best to perform PLS-DA when more than two experimental groups are present.

In the first tab, **Overview**, the first five components are shown. Scores plots and loadings plots are available and are similar to the plots described for the PCA.



**Cross Validation** of Q2 and R2 is used to select the optimal number of components that are used in the PLS-DA model for classification.

Generally speaking, an R2 value above 0.7 would be considered a substantial predictive ability. Similarly, the Q2 should be close to the R2 value indicating that the data is fitting the model.

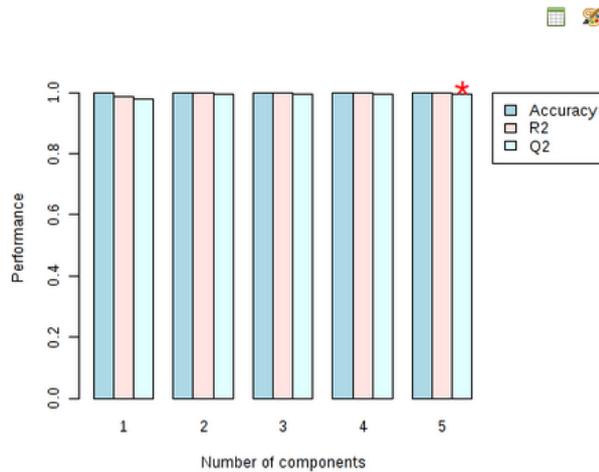
Overview | 2D Scores Plot | Loadings Plot | Imp. Features | Synchronized 3D Plots | **Cross Validation** | Permutation

Select optimal number of components for classification

Maximum components to search:

Cross validation (CV) method:

Performance measure:



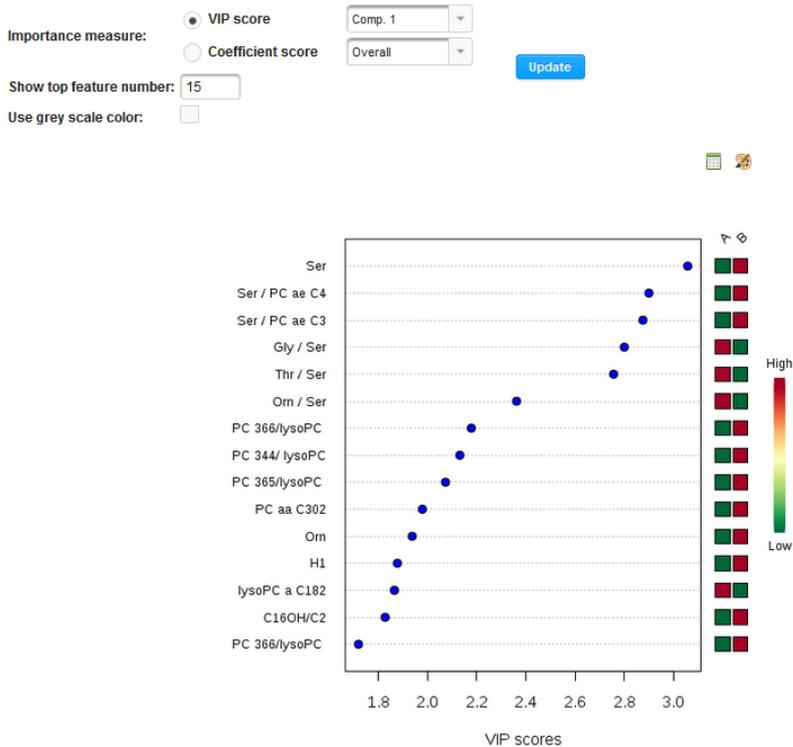
Variable importance in projection (VIP) is a weighted sum of squares of the PLS weight, which indicates the importance of the variable to the whole model. In many studies, VIP

## MetaboAnalyst Tutorial

values  $>2.0$  are selected and used for further data analysis, but this cut-off depends on the number of variables used.

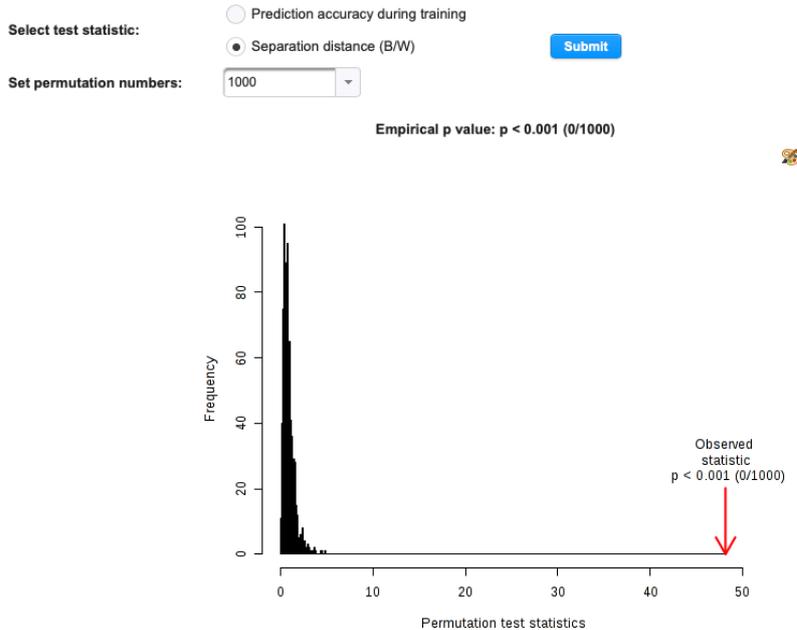


There are two importance measures in PLS-DA: one is variable importance in projection (VIP) and the other is the weighted sum of absolute regression coefficients (coef.). The colored boxes on the right indicate the relative concentrations of the corresponding metabolite in each group under study.



As pointed out before it is important to validate the results of PLS-DA analysis. This is done using permutation tests. The graph below is used to evaluate whether a class assignment is good or bad. This histogram shows the distribution formed by the permuted samples. The

arrow represents the original sample. The further away to the right of the distribution, the more significant the separation between the groups is.



### 5.2.7.3 Heatmap

Heatmap provides a visualization of the data table. The cells correspond to the concentration values in the data table. It is often useful to combine heatmaps with hierarchical clustering, which is a way of arranging items in a hierarchy based on the distance or similarity between them. The result of a hierarchical clustering calculation is displayed in a heat map as a dendrogram, which is a tree-structure of the hierarchy. Row dendrograms show the distance (or similarity) between rows and which nodes each row belongs to as a result of the clustering calculation. Column dendrograms show the distance (or similarity) between the

variables (the selected cell value columns). Select the parameters shown below and press “Submit”.

Distance Measure:

Clustering Algorithm:

Color Contrast:

Data Source:

Standardization:

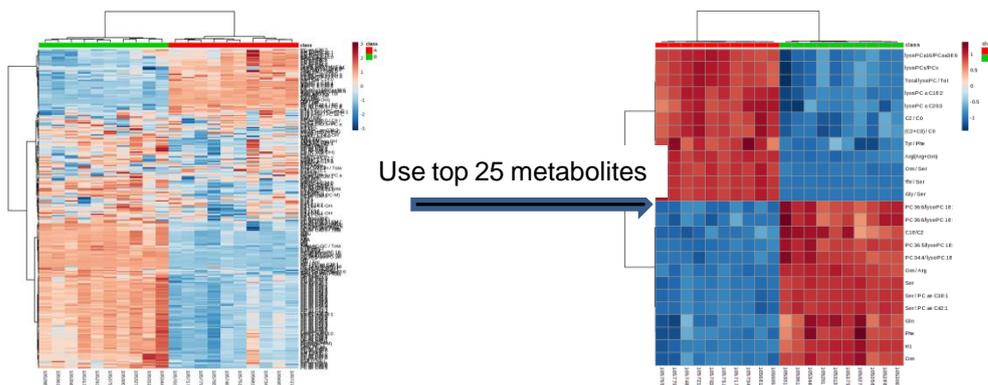
View Mode:  Overview  Detail View (< 2000 features)

Do not reorganize:

Use top:

View Options:  Show cell borders  Show only group averages

Displaying only the top x metabolites ranked by t-test can be used to show the most contrasting pattern (see below).



### 5.2.7.4 Pattern Hunter

The Pattern Hunter analysis provides a method for monitoring the changes that occur between sample groups. This can be particularly useful when comparing data collected as a time series. To monitor a pattern across a group of samples, a predefined pattern can be set for the series such as 1-2-3-4. A distance measurement such as Spearman's rank correlation coefficient can be used to rank metabolites in concentration order instead of the concentration magnitude.

Define a pattern using:

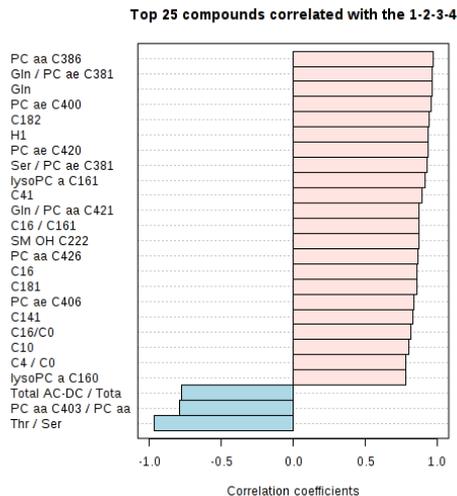
a feature of interest:

a predefined profile: 1-2-3-4

a custom profile:

Choose a distance measure: Spearman rank correlation

Submitting the settings will display the metabolites that correlate the best with the defined series.



In this way, metabolites that increase linearly with the time series or pattern can easily be identified. This works particularly well when evaluating GFR, disease progression, or response over time.

### 5.2.7.5 Download Results

After the analysis is completed the results can be downloaded. All the data (tables, figures) that have been created during the analysis can be downloaded separately or as a whole in a zip-file (Download.zip).

#### Result Download

Please download the results (tables and images) below. The **Download.zip** contains all the files in your home directory. You can also generate a **PDF analysis report** using the button below.

[Generate Report](#)

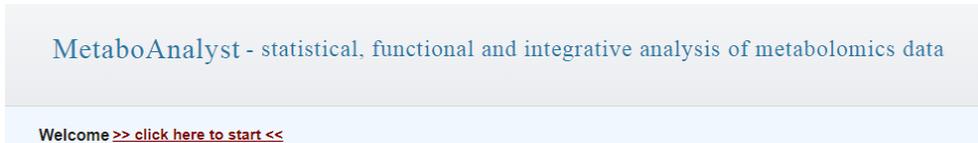
<a href="#">Download.zip</a>	<a href="#">t_test.csv</a>
<a href="#">Rhistory.R</a>	<a href="#">pls_pair_0_dpi72.png</a>
<a href="#">plsda_coef.csv</a>	<a href="#">fc_0_dpi72.png</a>
<a href="#">fold_change.csv</a>	<a href="#">2015-09-03_Conc test data 2 groups.csv</a>
<a href="#">volcano_0_dpi72.png</a>	<a href="#">pls_imp_0_dpi72.png</a>
<a href="#">heatmap_1_dpi72.png</a>	<a href="#">pls_loading_0_dpi72.png</a>
<a href="#">data_normalized.csv</a>	<a href="#">pca_score2d_0_dpi72.png</a>
<a href="#">plsda_vip.csv</a>	<a href="#">data_processed.csv</a>
<a href="#">pls_score2d_0_dpi72.png</a>	<a href="#">pca_loading_0_dpi72.png</a>
<a href="#">pls_cv_0_dpi72.png</a>	<a href="#">pca_loadings.csv</a>
<a href="#">plsda_score.csv</a>	<a href="#">pca_pair_0_dpi72.png</a>
<a href="#">tt_0_dpi72.png</a>	<a href="#">pca_scee_0_dpi72.png</a>
<a href="#">pca_biplot_0_dpi72.png</a>	<a href="#">pca_score.csv</a>
<a href="#">data_original.csv</a>	<a href="#">volcano.csv</a>
<a href="#">snorm_0_dpi72.png</a>	<a href="#">norm_0_dpi72.png</a>
<a href="#">heatmap_0_dpi72.png</a>	<a href="#">plsda_loadings.csv</a>

[Logout](#)

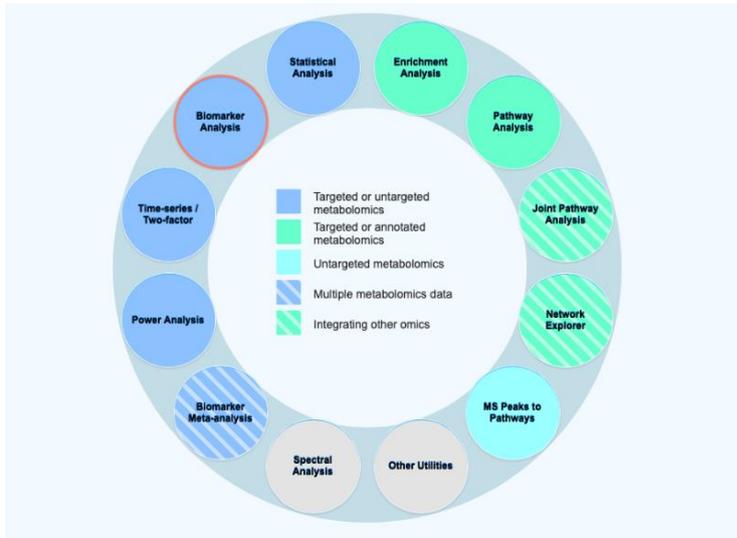
### 5.3 MetaboAnalyst Biomarker Analysis

Receiver operating characteristic (ROC) curves are generally considered the method of choice for evaluating the performance of potential biomarkers. MetaboAnalyst offers the possibility to perform ROC curve-based biomarker analysis for a single or multiple biomarkers. ROC curves can only be calculated from data sets containing two groups.

To begin a new biomarker analysis, navigate to the MetaboAnalyst Home page (<https://www.metaboanalyst.ca/>) and use the link “>> **click here to start** <<” to enter the module selection page.



Select the “Biomarker Analysis” module from the wheel.



### 5.3.1.1 Data Upload

In the Biomarker Analysis module, use the “Upload your data” panel, and select the options as indicated below. Choose the MetaboAnalyst CSV file created by MetIDQ and then click “Submit”.

#### Upload your data table (.csv or .txt):

Data Type:  Concentrations  Spectral bins  Peak intensity table

Format:

Data File:  No file selected.

After uploading the data, an integrity check is performed automatically, and the result is shown below. Missing value estimation and data filtering can be performed as described in the Statistical Analysis section.

### 5.3.1.2 Data Normalization

In addition to the normalization methods described for statistical analysis, it is possible to calculate metabolite ratios based on  $p$ -values. Ratios between metabolite concentrations may carry more information than the corresponding metabolite concentrations alone. This can improve the chances of biomarker discovery but also increases the potential for overfitting.

Otherwise, normalization can be performed as described earlier. Select the appropriate parameters for normalization and press “Submit” for the next step.

**Sample normalization**

- None
- Sample-specific normalization (i.e. weight, volume) [Specify](#)
- Normalization by sum
- Normalization by median
- Normalization by reference sample (PQN) [Specify](#)
- Normalization by a pooled sample from group [Specify](#)
- Normalization by reference feature [Specify](#)
- Quantile normalization

---

Compute and include metabolite ratios: Top 20 ▼

Ratios between two metabolite concentrations may carry more information than the two corresponding metabolite concentrations alone. MetaboAnalyst will compute ratios of all possible metabolite pairs and then choose top ranked ratios (based on p values) to be included in the data for further biomarker analysis. Note, there is a potential overfitting issue associated with the procedure. The main purpose here is to improve the chance of biomarker discovery. You need to validate the performance in independent studies. Log normalization will be performed during the process. You can only perform Data scaling in the next step.

---

**Data transformation**

- None
- Log transformation (generalized logarithm transformation or glog)
- Cube root transformation (takes the cube root of data values)

**Data scaling**

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

Normalize

View Result

Proceed

### 5.3.1.3 ROC Analysis

Different methods for ROC analysis are available. Select the method you want to perform. The classical univariate analysis is shown.

#### ROC Analysis Options :

##### Classical univariate ROC curve analyses

Perform classical univariate ROC curve analysis, 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.

##### Multivariate ROC curve based exploratory analysis (Explorer)

Perform automated important feature identification and performance evaluation. ROC curve analyses are performed based on three multivariate algorithms - support vector machines (SVM), partial least squares discriminant analysis (PLS-DA), and random forests.

##### ROC curve based model evaluation (Tester)

Users can manually select any combination of features to create biomarker models using any of the three algorithms mentioned above. The module also allows users to **hold out** a subset of samples for extra validation purpose, as well as to **predict class for new samples** (samples without class labels).

The results of the classical ROC curve analysis are shown below. By default, the analytes are sorted according to area under the curve (AUC). To view the ROC of an individual analyte press "View".

AUC values can be between 0.5 to 1.0. Where 0.5 would be a bad classifier and 1.0 indicates a good classifier.

### ROC curve analysis for individual biomarkers

The features displayed in the table below are ranked based on area under ROC curve (AUROC), T-statistics or Log2 fold change (FC). Click a header to sort the table accordingly. Click "View" in the last column to visualize its ROC curve and performance details. Note, the 95% confidence interval is calculated using 500 bootstrappings. You can use the **ROC Detail Analysis** button at the bottom for detailed sensitivity, specificity analysis at a given cutoff for the compound under current selection.

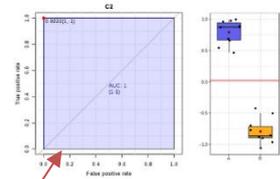
Calculate partial ROC curve limited by:

Parameter:  X-axis (max FPR)  Y-axis (min TPR)

Threshold:  range [0, 1] Update

Show optimal cutoff using:

Compute 95% confidence band (warning: time consuming)



Name	AUC	T-tests	Log2 FC	ROC Curve	Details
C2	1.0	4.8444E-13	1.3242	<a href="#">View</a>	<a href="#">→</a>
C3	1.0	7.3297E-10	0.6333	<a href="#">View</a>	<a href="#">→</a>
C4	1.0	4.2311E-7	0.59745	<a href="#">View</a>	<a href="#">→</a>
C5	1.0	1.7227E-8	0.55377	<a href="#">View</a>	<a href="#">→</a>
Arg	1.0	8.5083E-9	-0.59217	<a href="#">View</a>	<a href="#">→</a>
Gln	1.0	2.3712E-13	-0.99825	<a href="#">View</a>	<a href="#">→</a>
Gly	1.0	5.5448E-12	-0.82917	<a href="#">View</a>	<a href="#">→</a>
His	1.0	5.0034E-12	-0.80121	<a href="#">View</a>	<a href="#">→</a>
Met	1.0	1.2984E-9	-0.64768	<a href="#">View</a>	<a href="#">→</a>
Orn	1.0	2.3824E-17	-2.0599	<a href="#">View</a>	<a href="#">→</a>
Phe	1.0	4.6283E-14	-1.1362	<a href="#">View</a>	<a href="#">→</a>
Pro	1.0	7.5642E-10	-0.55842	<a href="#">View</a>	<a href="#">→</a>
Ser	1.0	1.1878E-22	-5.0994	<a href="#">View</a>	<a href="#">→</a>
Thr	1.0	3.2926E-12	-0.95996	<a href="#">View</a>	<a href="#">→</a>
Trp	1.0	1.8009E-11	-0.73101	<a href="#">View</a>	<a href="#">→</a>

#### 5.3.1.4 Download Results

After the analysis is completed the results can be downloaded. All the data (tables, figures) that have been created during the analysis can be downloaded separately or as a whole in a zip-file (Download.zip).

## MetaboAnalyst Tutorial

## Ordering and Technical Support

### **Order Absolute/IDQ® or Biocrates® Kits**

*E-Mail:*  
[sales@biocrates.com](mailto:sales@biocrates.com)

*Phone:*  
+43 512 579 823

### **Technical assistance**

*E-Mail:*  
[support@biocrates.com](mailto:support@biocrates.com)

### **Web Site**



<http://www.biocrates.com>

### **Video Tutorials**



[https://www.youtube.com/playlist?list=PLGETE8vMY-Plp\\_gSz4eMaSLG1QKB\\_mdFpk](https://www.youtube.com/playlist?list=PLGETE8vMY-Plp_gSz4eMaSLG1QKB_mdFpk)

### **Frequently Asked Questions (FAQ)**



<https://biocrates.com/support>

BIOCRATES Life Sciences AG  
Eduard-Bodem-Gasse 8, 6020 Innsbruck, Austria

tel: +43 512 579 823  
fax: +43 512 579 823 329

office@biocrates.com  
[biocrates.com](http://biocrates.com)