

## Identify Significant Features Using Concentration Data

By Jianguo Xia ([jianguox@ualberta.ca](mailto:jianguox@ualberta.ca))

Last update: 4/15/2009

This tutorial shows how to identify significant features using methods provided in MetaboAnalyst. The example used is compound concentration data obtained by targeted (i.e. quantitative) metabolic profiling of  $^1\text{H}$  NMR spectra of urine samples collected from 57 cancer patients. There are two groups of patients – Cachexia (Y) refers to the group with significant skeletal muscle loss; Cachexia (N) refers to the group with no obvious skeletal muscle loss. Cachexia is defined as the loss of weight, muscle atrophy, fatigue, weakness and significant loss of appetite in someone who is not actively trying to lose weight. Cachexia is often seen in end-stage cancer, and in that context is called "cancer cachexia". The exact mechanism behind cachexia is poorly understood, but there is probably a role for inflammatory cytokines, such as tumor necrosis factor-alpha ( $\text{TNF-}\alpha$ ) - which is also nicknamed cachexin, Interferon gamma ( $\text{IFN}\gamma$ ), and Interleukin 6 (IL-6). The goal in this tutorial is to identify metabolites that are significantly different between these two groups of cancer patients (cachexic vs. non-cachexic). These metabolites could serve as potential early-stage biomarkers for detecting cachexia and for exploring its underlying metabolic basis.

**Step 1.** Go to the **Data Formats** page (found by clicking on the **Data Formats** hyperlink on MetaboAnalyst's home page), click the download link after the “Compound concentration data” option to download the compressed zip file. Unzip the file and save as “compounds.csv”.

**Comma Separated Values (.csv) format ([show details](#)), including :**

- Compound concentration data ([download](#))
- Binned NMR/MS spectra data ([download](#))
- Paired time-series concentration data ([download](#))
- Processed peak intensity table ([download](#))

**Step 2.** Go the MetaboAnalyst **Home** page and click “click here to start” to enter the data upload page.



**Step 3.** In the **Upload** page, go to the “Upload your data” panel, select the options as indicated below and click “Submit”

**1) Upload your data :** [? data format](#)

**Comma Separated Values (.csv) :**

**Data type :**  **Concentrations**  **Spectral bins**  **Peak intensity table**

**Format:**

**Data file :**

**Note:** Alternatively, you can directly select the first option in the “Try our test data” without downloading the example.

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

Data Type	Description
<input checked="" type="radio"/> <b>Concentrations</b> <a href="#">Tutorial Report</a>	Urinary compound concentrations from 57 cancer patients measured by 1H NMR (unpublished data). Group 1- no weight loss; group 2 - severe weight loss
<input type="radio"/> <b>NMR spectral bins</b> <a href="#">Tutorial Report</a>	Binned 1H NMR spectra of 50 urine samples using 0.04 ppm constant width ( <a href="#">Psihogios NG, et al.</a> ) Group 1- control; group 2 - severe kidney disease.

**Step 4.** The data integrity check will run automatically and the result is shown below. For lists of concentrations the data integrity check will assess the content (look for consistent formatting and the presence of two groups), determine whether the data is paired or determine if negative numbers exists. In this case, 869 zero values and no missing values were identified in the data. Since zero values may cause some algorithms not to work properly, MetaboAnalyst will replace these zero values with a small positive value (the half of the minimum positive number detected in the data). Click “Skip” to go to normalization step. If missing values had been detected, then the most appropriate from a variety of methods provided by MetboAnalyst could have been used to deal with this issue (for such an example, see MetaboAnalyst Tutorial 4).

**Note:** missing values are represented as NA (no quotes) or empty values.

**Data processing information**

Checking data content ...passed

Two groups were detected based on the sample labels.

Samples are not paired.

All data values are numeric.

All data values are non-negative.

A total of 869 , ( 28.2 %) zero values were detected

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 imputation      Skip

**Step 5.** Now we arrive at the data normalization step. The internal data structure is transformed now to a table with each row representing a urine sample (from a patient) and each column representing a feature (a compound with a concentration). With the data structured in this format, two types of data normalization protocols - row-wise normalization and column-wise normalization -- may be used. These are often applied sequentially to reduce systematic variance and to improve the performance for downstream statistical analysis. Row-wise normalization aims to normalize each sample (row) so that they are comparable to each other. For row-wise normalization MetaboAnalyst supports normalization to a constant sum, normalization to a reference sample (probabilistic quotient normalization), normalization to a reference feature (creatinine or an internal standard) and sample-specific normalization (dry weight or tissue volume). In contrast to row-wise normalization, column-wise normalization aims to make each feature (column) more comparable in magnitude to each other. Four widely-used methods are offered in MetaboAnalyst - log transformation, auto-scaling, Pareto scaling, and range scaling. Urine concentrations are usually normalized by creatinine concentration to adjust for dilution effects (select option 4 - Normalization by a reference feature and choose “creatinine”). However, in this case, creatinine is the product of protein breakdown and is related to the skeletal muscle loss. Since normalizing to a metabolite that might be important for understanding this disorder might introduce biological bias, we need to choose another kind of normalization process. As a result we choose to normalize by a reference sample “NETCR4” (a general rule is to choose a sample in the control group with the fewest missing values). After deciding to normalize by reference sample for our row-wise (sample) normalization we then choose “Log normalization” for our column normalization to make the metabolite concentration values more comparable among different compounds.

**Hint:** Remember to click on Normalization by a reference sample and not only select NETCR4

The image shows a software interface for normalization settings. It is divided into two main sections: "Row-wise normalization" and "Column-wise normalization".

**Row-wise normalization**

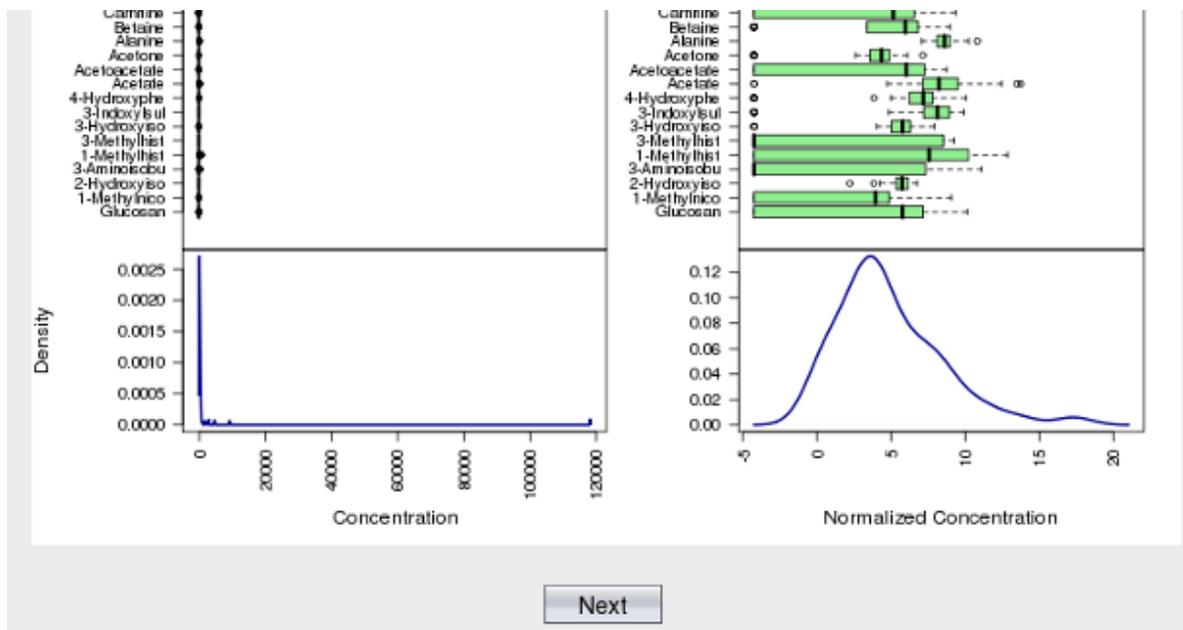
- None
- Normalization by sum
- Normalization by a reference sample NETCR4 ▾
- Normalization by a reference feature <Not set> ▾
- Sample specific normalization (i.e. dry weight, volume) [Click here to specify](#)

**Column-wise normalization**

- None
- Log (log<sub>2</sub> transformation)
- Autoscaling (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)

## MetaboAnalyst Tutorial 1

The normalization result is shown below. On the left is a plot (box-whisker plot on top, linear distribution plot on the bottom) of the data prior to normalization. On the right is a plot (box-whisker plot on top, linear distribution plot on the bottom) of the data after normalization. As can be seen by comparing the linear concentration curve on the left (which has an exponential decay character) to the log-transformed curve on the right (which looks reasonably Gaussian), the log normalization step along with the reference sample normalization makes the concentration data reasonably “normal”. You can also try other normalization approaches and compare their results.



**Step 6.** Now we have finished data processing and normalization. The data are now suitable for different statistical analyses. There are many feature selection methods available in MetaboAnalyst. Here we will only show results from Volcano Plot, PLS-DA and SAM methods. The screen shot below shows MetaboAnalyst's analysis view. Please note the navigation panel on the left. A color change indicates the corresponding step has been successfully performed. All the data analysis methods can be directly accessed by clicking the corresponding hyperlink.

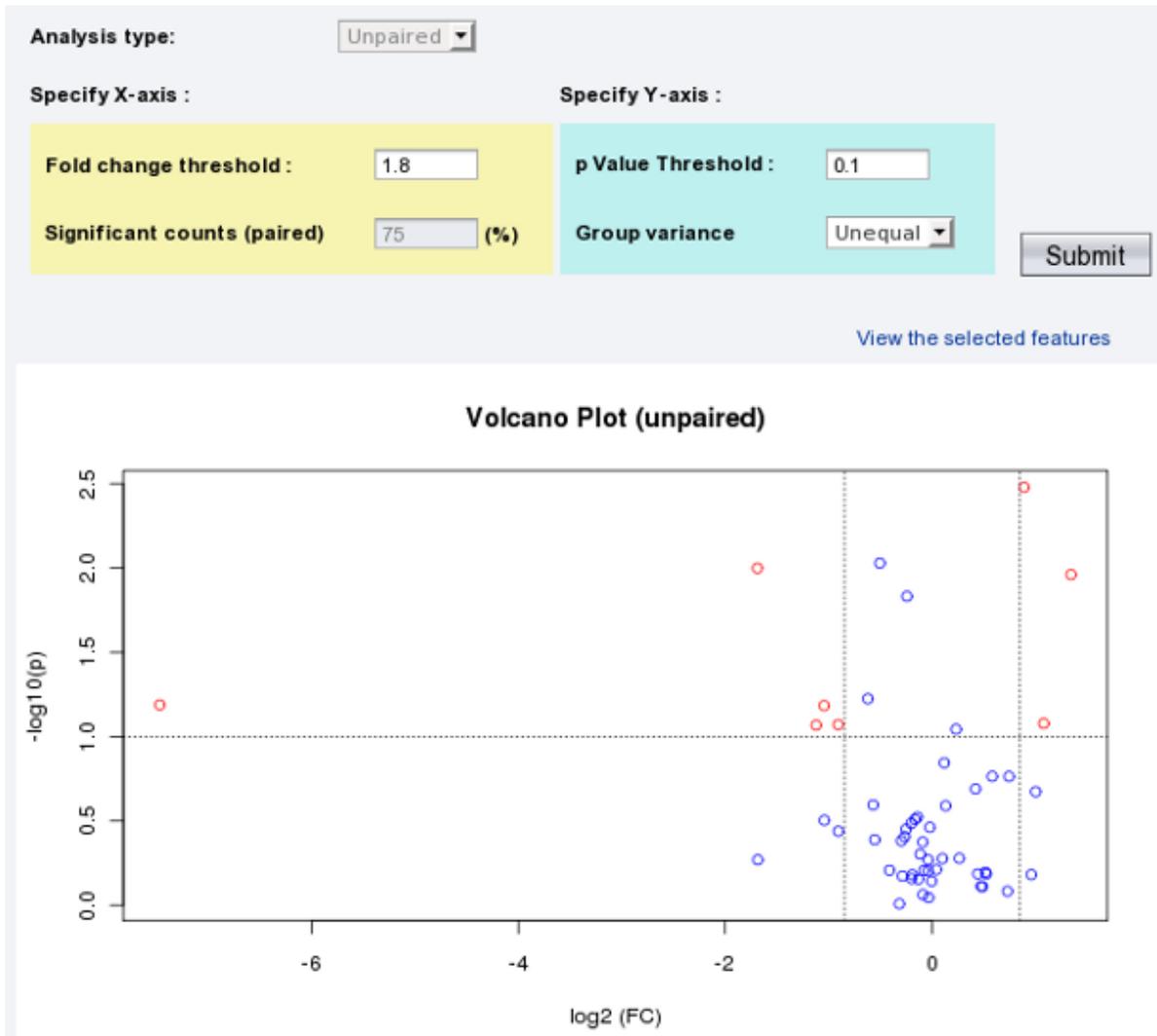
**Steps**

- [1. Upload](#)
- [2. Process](#)
- [3. Normalize](#)
- 4. Analyze**
  - [Univariate](#)
  - [PCA](#)
  - [PLSDA](#)
  - [SAM](#)
  - [EBAM](#)
  - [Tree & heatmap](#)
  - [Kmean & SOM](#)
  - [RandomForest](#)
  - [R-SVM](#)
- [5. Peak Search](#)
- [6. Pathway Mapping](#)
- [7. Download](#)
- [Log Out](#)

**Select an analysis path to explore :**

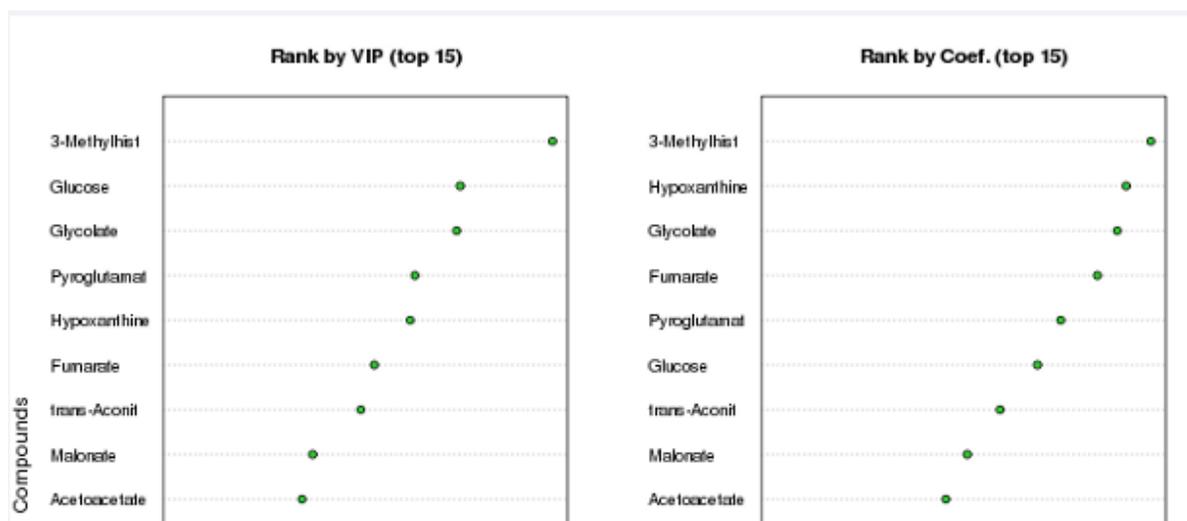
- Univariate Analysis**
  - [Fold Change, t-Tests, and Volcano plot](#)
- Chemometrics**
  - [Principal Component Analysis \(PCA\)](#)
  - [Partial-Least Square - Discriminant Analysis \(PLS-DA\)](#)
- Significant Feature Identification**
  - [Significance Analysis of Microarray \(SAM\)](#)
  - [Empirical Bayesian Analysis of Microarray \(EBAM\)](#)
- Cluster Analysis**
  - [Hierarchical Clustering - Dendrogram and Heatmap](#)
  - [Partitional Clustering - K-Means and Self Organizing Map \(SOM\)](#)
- Classification & Feature Selection**
  - [Random Forest](#)
  - [Support Vector Machine \(SVM\)](#)

**Step 7.** Generally the simplest kind of analysis that can be performed on this type of metabolomic data is univariate data analysis. Univariate analyses are often first used to obtain an overview of the data or a rough ranking of potentially important features before applying more sophisticated data analysis tools. Univariate analysis examines each variable separately without taking into account the effect of multiple comparisons. MetaboAnalyst's univariate analysis path supports three commonly used methods - fold-change analysis, t-tests, and volcano plots. To begin the univariate analysis, click the "Univariate" link on the navigation panel to the left. From here we will perform a volcano plot analysis. Volcano plots are used to compare the size of the fold change to the statistical significance level. The X axis plots the fold change between the two groups (on a log scale), while the Y axis represents the p-value for a t-test of differences between samples (on a negative log scale). To start a volcano plot click the "Volcano" tab. Adjust the fold change (FC) threshold to 1.8 and click "Submit". As can be seen by the figure below, eight features are detected as significant and colored in red. Click the "View the selected features" link to view the names/identities of these features. The table at the bottom shows these eight features. These include glycolate, fumarate, 3-methylhistidine, glucose, *etc.*



Compounds	FC	log2(FC)	p.value	-log10(p)
Glycolate	1.856	0.893	0.0030	2.478
Fumarate	0.31	-1.688	0.01	1.999
3-Methylhistidine	2.543	1.347	0.011	1.961
Glucose	0.0060	-7.471	0.065	1.189
Pyroglutamate	0.485	-1.043	0.065	1.184
Acetoacetate	2.116	1.081	0.084	1.078
Choline	0.533	-0.908	0.085	1.071
Malonate	0.46	-1.119	0.086	1.068

**Step 8.** The Volcano plot has provided some intriguing results. We may now want to examine whether these metabolites are also detected as being significant using a slightly more sophisticated analysis tool. In particular we will use Partial-Least Squares Discriminant Analysis (PLS-DA). As a supervised method, PLS-DA can perform both classification and feature selection. The algorithm uses cross-validation to select an optimal number of components for classification. Two feature importance measures are commonly used in PLS-DA. Variable Importance in Projection or VIP score is a weighted sum of squares of the PLS loadings. The weights are based on the amount of explained Y-variance in each dimension. The other importance measure is based on the weighted sum of PLS-regression coefficients. The weights are a function of the reduction of the sums of squares across the number of PLS components. More details about these two methods can be obtained by placing your mouse over the “About PLS” link. Go back to the Analysis window and click the “PLSDA” link on the navigation panel and then click the “Var.Importance” tab. You will see the result as shown below. The graphs rank the different metabolites (the top 15) according to the VIP score on the left and according to the coefficient score on the right.



Click the “View details” link to see the data table that was used to produce the graph. Note that the VIP score is not normalized. The VIP scores tend to be 200X larger than the coefficient scores, but the relative ranking of “significant” metabolites is largely the same. 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 values >2.0 are selected and used for further data analysis, but this cut-off depends on the number of variables used. Since the number of variables in this study is less than 100, we can use a more relaxed VIP cutoff of around 1.0.

Compounds	VIP score	Compounds	Coef. score
3-Methylhistidine	2.4964	3-Methylhistidine	0.0127
Glucose	2.1209	Hypoxanthine	0.0122
Glycolate	2.1064	Glycolate	0.012
Pyroglutamate	1.9364	Fumarate	0.0116
Hypoxanthine	1.9172	Pyroglutamate	0.0108
Fumarate	1.7711	Glucose	0.0104
trans-Aconitate	1.7167	trans-Aconitate	0.0096
Malonate	1.5208	Malonate	0.0090
Acetoacetate	1.4774	Acetoacetate	0.0085
O-Acetylcarnitine	1.2359	Choline	0.0070
Choline	1.2335	O-Acetylcarnitine	0.0065
1-Methylhistidine	1.0894	Ethanol	0.0058

**Step 9.** With the completion of the PLS-DA analysis, we can try another approach to select interesting or significant features that distinguish between cachexic and non-cachexic patients. Here we'll attempt to use Significance Analysis of Microarray (SAM). SAM is designed to address False Discovery Rate (FDR) problems when running multiple tests on high-dimensional data. It first assigns a significance score to each variable based on its change relative to the standard deviation of repeated measurements. Then it chooses variables with scores greater than an adjustable threshold and compares their relative difference to the distribution estimated by random permutations of the class labels. For each threshold, a certain proportion of the variables in the permutation set will be found to be significant by chance. This number is used to calculate the FDR. To use SAM analysis, go back and click the "SAM" link on the navigational panel; you will see the following set of "Delta" plots. The Delta plots are a visualization of the table generated by SAM that contains the estimated FDR and the number of identified metabolites for a set of Delta values. Note the pop-up help balloon when you place the mouse over "About SAM". The default Delta value (0.5) has an FDR of 12% and identifies ~10 significant compounds above this threshold, as seen on the left plot. You can increase the Delta to reduce the FDR. The figure below shows that when Delta is above 0.7, FDR approaches to 0 (left panel). However, no significant compound will be identified (right panel). The default 0.5 is a compromise between the FDR and the number of compounds detected. Click 'Submit' in the bottom panel to go to Step 2 to view the result.

**Step 2** ? About SAM

Please note: the Group variance will be i

**Method**

**Samples :**

**Group variance**

Choose a delta value to control false

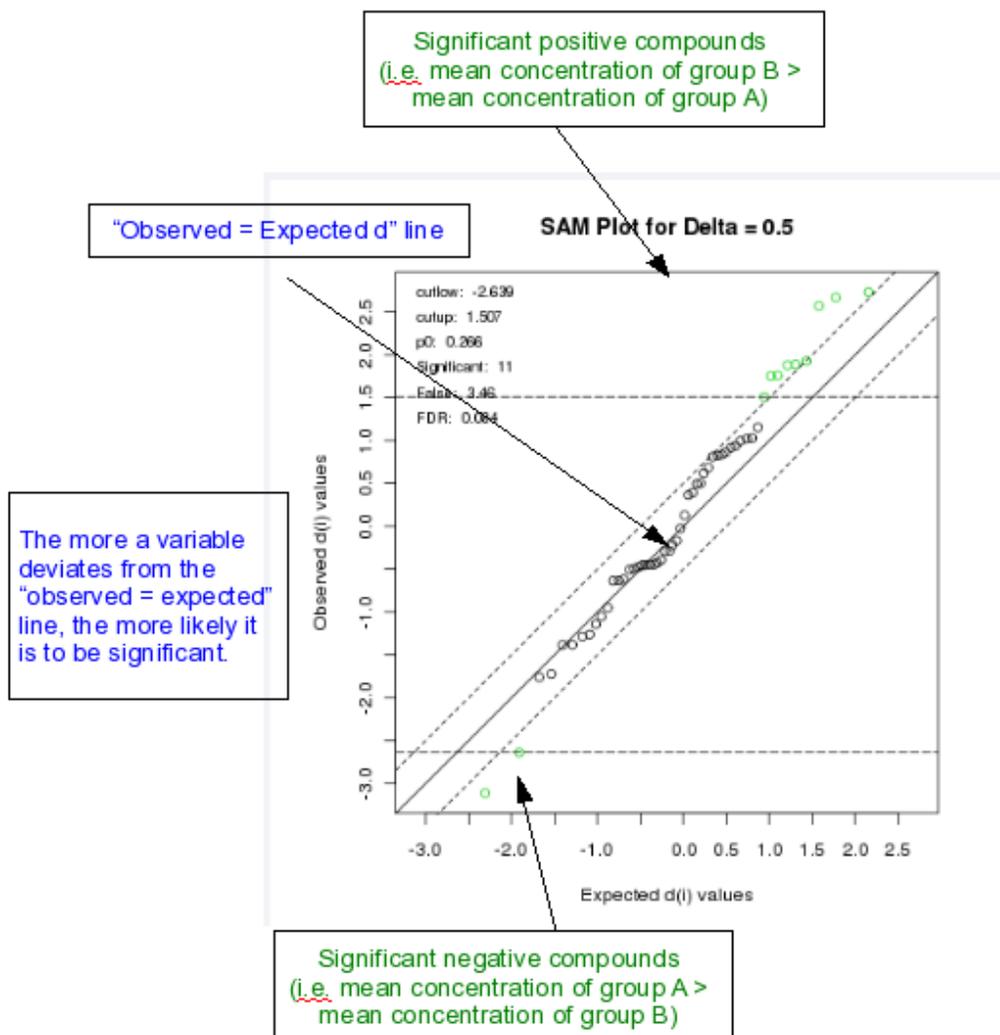
SAM is a well-established statistical method for identification of differentially expressed genes in microarray data analysis. It is designed to address false discovery rate (FDR) when running multiple tests on high-dimensional data. SAM assigns a significance score to each variable based on its change relative to the standard deviation of repeated measurements. For variable with scores greater than an adjustable threshold, its relative difference is compared to the distribution estimated by random permutations of the class labels. For each threshold, a certain proportion of the variables in the permutation set will be found to be significant by chance. The number is used to calculate the FDR. Users need to specify the **Delta** to control FDR in order to proceed.

Delta	FDR (in %)
0.1	22
0.2	18
0.3	18
0.4	14
0.5	8
0.6	6
0.7	6
0.8	2
0.9	2
1.0	0

Delta	No. of Significant Compds
0.1	45
0.2	27
0.3	27
0.4	21
0.5	11
0.6	10
0.7	10
0.8	4
0.9	3
1.0	0

Delta =

The Step 2 tab shows a typical SAM plot with Delta = 0.5. Click the “View details of the ...” button to see the SAM results table (next page). A SAM plot displays a “positive” metabolite set and a “negative” metabolite set. In the positive metabolite set, higher levels of these metabolites correlate with higher values for the cachexia phenotype. In the negative metabolite set, lower levels of these metabolites correlate with higher values for the cachexia phenotype. A total of 11 compounds were identified above the chosen threshold.



Significant compounds identified by SAM with  $\delta = 0.5$ . Note that the term “rawp” refers to the raw p-values from regular t-tests.

Compounds	d.value	stdev	rawp	q.value	R.fold
Glycolate	-3.114	0.986	0.0020	0.027	0.119
trans-Aconitate	2.729	0.919	0.0050	0.027	5.69
Fumarate	2.667	0.803	0.0070	0.027	4.414
3-Methylhistidine	-2.639	1.35	0.0080	0.027	0.085
Hypoxanthine	2.569	0.93	0.01	0.028	5.235
Trigonelline	1.925	0.361	0.056	0.102	1.618
Glucose	1.887	1.541	0.061	0.102	7.505
Pyroglutamate	1.879	1.418	0.063	0.102	6.343
Choline	1.756	1.026	0.081	0.102	3.487
Malonate	1.751	1.246	0.082	0.102	4.536
3-Hydroxyisovalerate	1.507	0.618	0.13	0.144	1.908

**Step 10.** Based on the result from the Volcano Plot, PLS-DA and SAM, several compounds are consistently identified as being significant by different approaches. Using 3-methylhistine and glycolate as examples, let us further check which pathways they are involved in. To do so, we can use MetaboAnalyst’s data annotation tools. Click the “Pathway mapping” link on the left navigation panel, enter the two compound names separated by a semicolon (i.e. ;) and then click the “Search” button. The result shows only the pathway for Glycolate. No entry was found for 3-Methylhistidine in the pathway library of the Human Metabolome Database (HMDB). By clicking on the relative links in the resulting table, you will access the corresponding pathway as well as detailed information about the metabolite (Metabocard).

**Pathway mapping :**

Please enter compound names (separated by semi-colon):

3-methylhistidine; glycolate

---

[View all library hits](#)

Pathway Name	Members
<a href="#">Glyoxylate and Dicarboxylate Metabolism a.html</a>	<a href="#">Glycolic acid</a>

**Step 11.** Now, we want to find out about the biological function of 3-methylhistidine. Go to the HMDB ([www.hmdb.ca](http://www.hmdb.ca)), and enter “3-methylhistidine” and click “Search” button. The result is shown below.

## Human Metabolome Database

Version 2.0 | [Version 1.0](#)

Search:   [\[Advanced\]](#)

### Search Results

Search for "3-methylhistidine" returned 1 results

Showing 1-1 out of 1 hits

HMDB ID	Name	Formula	Weight
HMDB00479 <a href="#">MetaboCard</a>	<b>3-Methylhistidine</b> ... N(pi)-methyl-L-histidine; Tau-methylhistidine; I- <b>3-methylhistidine</b> ; 3-N-Methyl-L-histidine; N3-Methyl-L-histidine; pi-Methyl-L-histidine; 3-methyl-L-Histidine; N(pai) ...	C <sub>7</sub> H <sub>11</sub> N <sub>3</sub> O <sub>2</sub>	169.1811

Click the “MetaboCard” on the left panel; the result is shown below. As indicated, this compound can be used as an index of muscle protein breakdown which is relevant to the cachexia patients with significant skeletal muscle loss.

### Showing metabocard for 3-Methylhistidine (HMDB00479)

Legend: metabolite field enzyme field

St

Version	2.0
Creation Date	2005-11-16 15:48:42
Update Date	2005-11-16 15:48:42
Accession Number	HMDB00479
Common Name	<b>3-Methylhistidine</b>
Description	3-Methylhistidine is a product of peptide bond synthesis and methylation of actin and myosin. The measurement provides an index of the rate of muscle protein breakdown.

**Step 12.** Now, assume we have finished the analysis. Click the “Download” link on the navigation panel. A detailed analysis report will be generated (MetaboAnalystReport.pdf) containing introductions and results from every steps you have performed. Now, you can directly click “Download.zip” file to download all the processed data, images and the PDF report. Alternatively, you can ask MetaboAnalyst to send you the result via Email by entering your email address. The data will remain on the server for 72 hours before being automatically deleted.

Email address :

Files in your home directory	
<a href="#">Download.zip</a>	<a href="#">MetaboAnalystReport.pdf</a>
<a href="#">compounds.csv</a>	<a href="#">conc-norm.png</a>
<a href="#">data_normalized.csv</a>	<a href="#">data_original.csv</a>
<a href="#">data_processed.csv</a>	<a href="#">pls_Class.png</a>
<a href="#">pls_loading.png</a>	<a href="#">pls_pair.png</a>
<a href="#">pls_permut.png</a>	<a href="#">pls_score2d.png</a>
<a href="#">pls_score3d.png</a>	<a href="#">pls_VIP.png</a>
<a href="#">Rhistory.R</a>	<a href="#">sam-cmpd.png</a>
<a href="#">sam_fdr.png</a>	<a href="#">univar_fc.png</a>
<a href="#">univar_t.png</a>	<a href="#">univar-volcano.png</a>

-----End of tutorial-----