**Shotgun Data Profiling (SDP)**

# Goal for this tutorial

- To perform an exploratory and biomarker analysis on shotgun metagenomics data and visualize the results within KEGG metabolic networks along with pathway analysis.

MicrobiomeAnalyst -- comprehensive statistical, visual and meta-analysis of microbiome data

**Click here to start**

Starting from marker gene abundance data (OTU table, BIOM file, mothur output)

Marker Data Profiling (MDP)

Shotgun Data Profiling (SDP)

Starting from gene list or gene abundance data annotated by KO, EC or COG

Visually exploring your 16S rRNA data with a public data in a 3D PCoA plot

Projection with Public Data (PPD)

Taxon Set Enrichment Analysis (TSEA)

Starting with a list of taxa of interest (strains, species or higher level taxa)

# Shotgun Data Profiling (SDP)



Two types of user inputs:

❖ A list of gene IDs.

❖ Abundance table (in text or BIOM format)

Note genes need to be annotated in KO, EC, or COG for functional analysis,

# A) 1. Upload a list



Upload a list of gene IDs

3 gene ID types supported (KO, COG and EC Number).

Gene ID type — KEGG Orthology IDs (KO)

Try our example

You can try our example also

Step 1 : Choose the parameters above. Copy and paste a list of gene IDs along with their expression value

```
K01623  5
K00128  24
K00016  38.5
K00873  53
K01689  90
K01834  132.5
K00134  77
K01803  28.5
K00850  106
K01810  108
K01835  48
K01792  32
K01785  29
K00382  42
K00927  83.5
K00886  18
K01222  4
```

Submit

Step 2 : Click "Submit" to proceed.

# 2. Data Integrity Check

**Data processing summary**

Uploaded gene ID type: ko

Abundance measure **provided**

Total number of genes: 568

Mapped to our database: 563

The abundance range: [ 1.0 - 309.0 ]

By default, all genes will be used for analysis in the next stage

You can further **Filter genes** on the right panel by their abundance (if available).

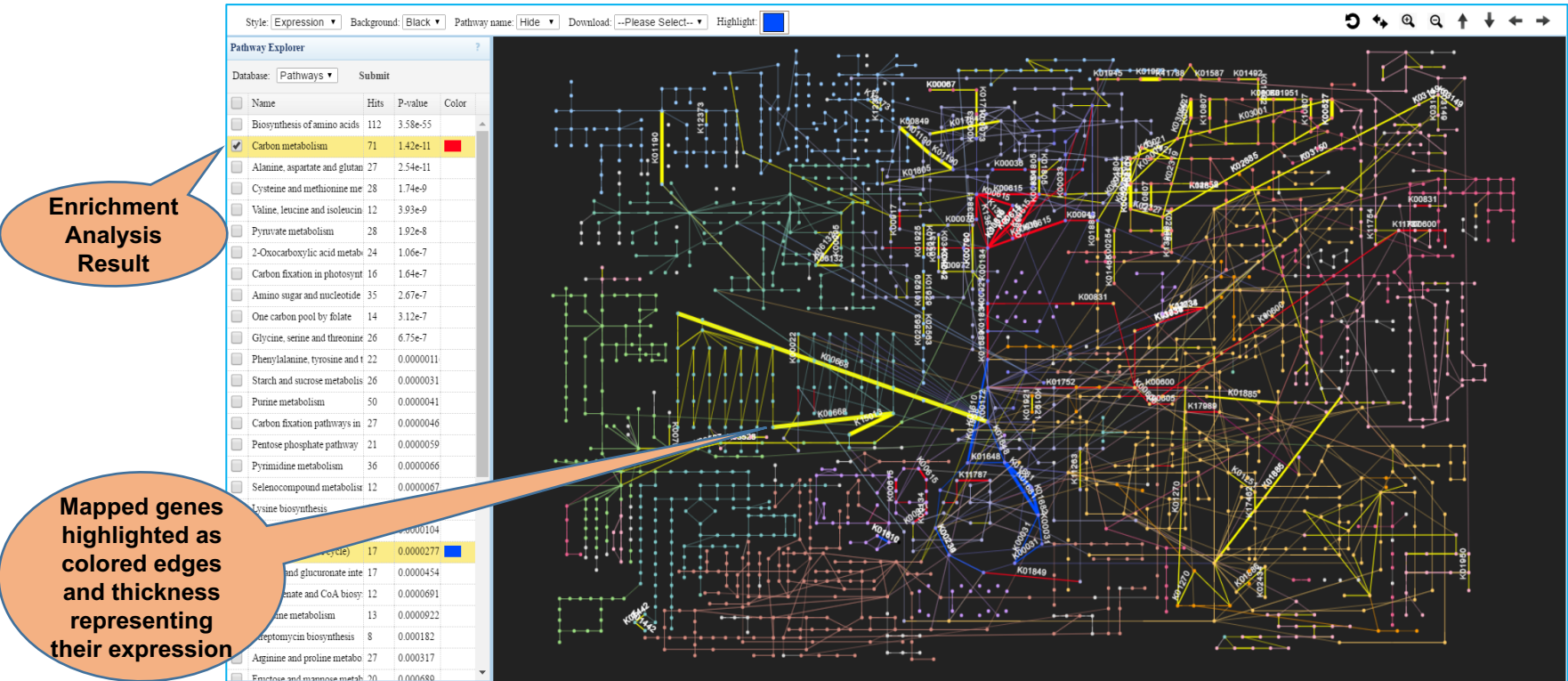Or click the **Proceed** button at bottom right to proceed.

Filter low count
genes:                                5        Update

genes with low count can
be filtered out

Click "Proceed" to visualize
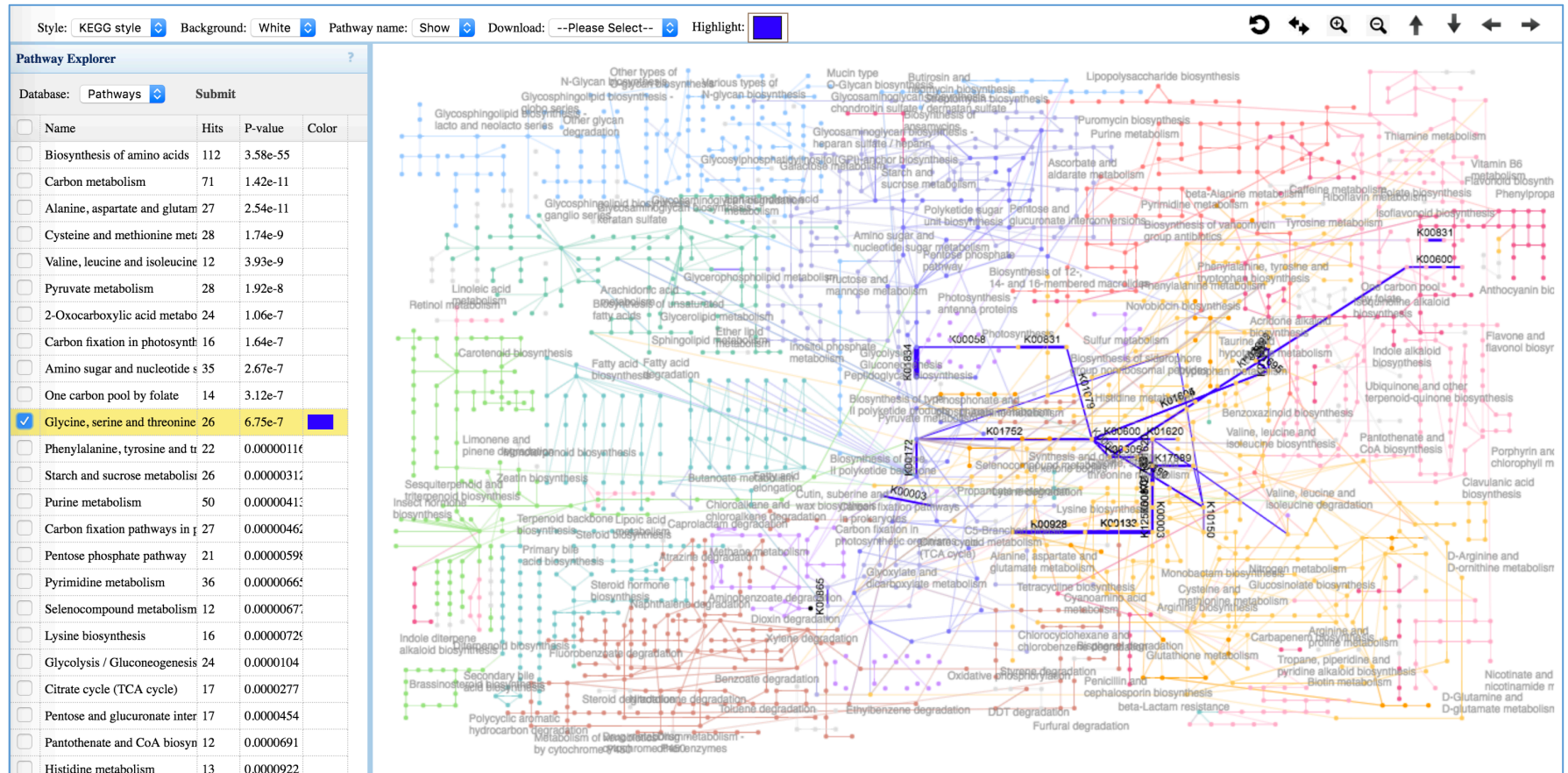the result within KEGG
metabolic network

Provides processing and summary information for user uploaded gene list.

# 3. KEGG Metabolic Networks (I)



1. Click "Submit" on the Pathway Explorer to perform pathways enrichment analysis.

2. Select a highlight color (default orange)

3. Click on a pathway name (a table row) to highlight the corresponding pathways

   - Gene IDs (KO) are represented as edge (reaction linking two metabolites) in the network and its thickness are based on their expression levels.

# 3. KEGG Metabolic Networks (II)



Customizing the styles using the menus on the top too bar, for example:

- Switching background from black to white;
- Showing the pathway names.

# B) Analyzing shotgun gene count data

# Data Formatting

**1. Tab-delimited text file**

- User have to upload both gene abundance table and metadata file separately.

- Manipulate data headings in a spreadsheet program like MS Excel

- Save as a **tab delimited (.txt) or comma-separated (.csv) file**

- The headings **#NAME** : (all capital letters) must be used
  - ❖ #NAME is for sample names (first column in abundance; first row in metadata file)
  - ❖ 2nd Column of metadata file is for the clinical metadata.

**2. BIOM format**

- Standard format for storing gene abundance information (metadata file separately in .txt file).

**For Example:**

| #NAME | sample1 | sample2 | sample3 | sample4 | sample5 |
|---|---|---|---|---|---|
| COG0002 | 1 | 2 | 2 | 2 | 3 |
| COG0005 | 1 | 0 | 0 | 1 | 2 |
| COG0006 | 1 | 4 | 0 | 1 | 2 |
| COG0008 | 1 | 1 | 1 | 2 | 1 |
| COG0009 | 2 | 1 | 0 | 2 | 0 |
| COG0012 | 1 | 0 | 2 | 1 | 1 |
| COG0013 | 1 | 2 | 0 | 1 | 0 |
| COG0014 | 2 | 1 | 0 | 0 | 1 |
| COG0015 | 0 | 0 | 1 | 1 | 0 |
| COG0016 | 2 | 0 | 0 | 1 | 1 |
| COG0017 | 1 | 1 | 0 | 4 | 0 |
| COG0018 | 4 | 3 | 2 | 1 | 0 |
| COG0019 | 2 | 3 | 2 | 2 | 3 |
| COG0020 | 1 | 1 | 0 | 0 | 1 |
| COG0021 | 1 | 0 | 1 | 1 | 0 |

| #NAME | Type |
|---|---|
| sample1 | lean |
| sample2 | lean |
| sample3 | lean |
| sample4 | obese |
| sample5 | obese |

**Abundance table and Metadata file in tab-delimited (.txt) format**

# 1. Data Upload

Step 1: Upload your gene abundance profile data in table or BIOM format

Step 2: Chose a gene ID type
3 IDs supported (KO, COG and EC numbers)

Upload a list of gene IDs

Upload a gene abundance table

| Gene ID type | -- Please Specify -- |
| Abundance file (.txt or .csv) | Choose File No file chosen ❓ |
| Metadata file (.txt or .csv) | Choose File No file chosen ❓ |

Submit

Step 3: Upload your abundance data file

Step 4: Upload your metadata file

Upload a BIOM file

Example data sets for testing

Step 4 : Click "Submit" to proceed

Example data sets for testing

| Data Type | Format | Description |
| --- | --- | --- |
| ⦿ KO Dataset | Plain text | A test example containing KO annotated read counts from 20 samples. **Class:** Diseased (10 samples), Normal (10 samples). |

Submit

You can try our example also

# 2. a) Data Integrity Check

| | |
|---|---|
| **Text Summary** | Graphic Summary |

| | |
|---|---|
| Data type: | Gene abundance table |
| File format: | text |
| Gene annotation: | ko |
| Total gene number: | 1000 |
| Genes with ≥ 2 counts: | 1000 |
| Sample number: | 20 |
| Number of experimental factors: | 1 |
| Total read counts: | 146868319 |
| Average counts per sample: | 7343415 |
| Maximum counts per sample: | 8016826 |
| Minimum counts per sample: | 6757015 |

Provides processing and summary information for user uploaded data.

# 2. b) Graphic Summary



- Provides user the information about library size or total number of reads present in of each sample and help in identifying the potential outliers due to undersampling or sequencing errors.

# 3. a) Data Filtering (Features)



- Identifying and removing variables or features that are unlikely to be of use when modeling the data. (e.g., features containing all zeros or constant across all the samples)

- **6** different approaches: on the basis of count (**abundance**) or using **statistical** approaches such as **mean, median, IQR, standard deviation or C.V.**

# 3. b) Sample Filtering (Editor)



User can select samples to remove from downstream analysis

- Users can remove samples that are detected as outlier via graphical summary result or downstream analysis. (e.g. Beta-diversity analysis)

# 4. Data Normalization



- Normalizing is required to account for **uneven sequencing depth**, **under-sampling** and **sparsity** present in such data. (useful before any meaningful comparison)

- Several normalization methods which have been commonly used in the field are present. (2 categories: **data scaling and data transformation** )

# 5. Data analysis

User can get an overview along with comparative and functional analysis of shotgun data.

Shotgun Data Analysis

Overall functional profiling
- Functional diversity overview
- Functional association analysis

Clustering analysis
- Heatmap clustering
- Dendrogram clustering
- PCA visualization
- Correlation analysis
- Pattern search

Differential abundance analysis
- Univariate analysis
- metagenomeSeq
- RNAseq methods

Biomarker analysis
- LEfSe
- Random Forests

# A. Functional Profiling



**Functional Diversity Profiling**

Functional category — KEGG metabolism

Calculate category abundance by — Total hits

Group samples based on — Phenotype

Color scheme — Palette_21

Submit

**User can select from different categories based on input gene id type :**

- **KEGG metabolism, pathways, modules or COG or EC functional category**

**The abundance of functional categories can be estimated by 3 different method to account for one to many gene mapping issue**

**Samples colored on the basis of selected experimental factor**

KEGG metabolism
- Amino acid metabolism
- Biosynthesis of other secondary metabolites
- Carbohydrate metabolism
- Energy metabolism
- Glycan biosynthesis and metabolism
- Lipid metabolism
- Metabolism of cofactors and vitamins
- Metabolism of other amino acids
- Metabolism of terpenoids and polyketides
- Nucleotide metabolism
- Xenobiotics biodegradation and metabolism

## 1. Functional Diversity Profiling

- Samples have been compared to provide a coarser view of the data by collapsing related genes (KO, COG or EC) to observations of functions. (rather than observations of specific genes)
- **5** main functional categories present to collapse within based on **gene ID type** : **KEGG metabolism, pathways, modules and COG or EC functions**.

# A. Functional Profiling

**2. Functional Association analysis and Metabolic Network Exploration:** associations between any functional categories with the experimental factor or sample groups is calculated by integrating the abundance changes of all members within each functional group to evaluate the strength of association
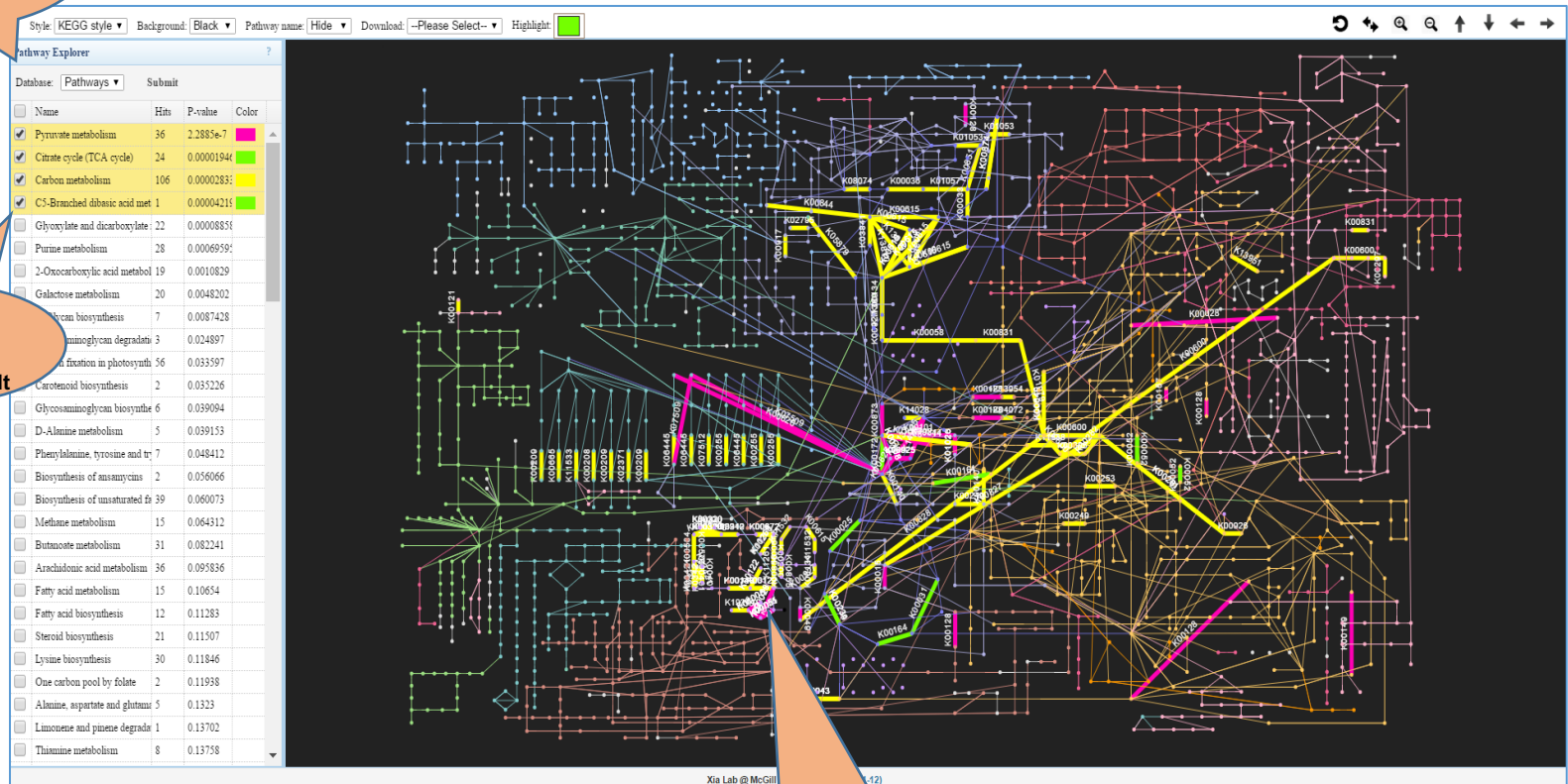
- It is based on the globaltest algorithm. For details:
    - "A global test for groups of genes: testing association with a clinical outcome". Bioinformatics 2004 Jan 1;20(1):93-9.

- Significant functional categories (pathways and modules) can be visualized within Metabolic networks.

# A. Functional Profiling

## 2. Functional Association Testing and Metabolic Network Exploration:



User can chose from 2 functional categories : pathways or modules

Functional categories association analysis Result

Significant functional categories (pathways or modules) can be highlighted with different colors

# B. Clustering Analysis

## 1. Principal Component Analysis (PCA)

- Data reduction technique that can be used to visualize the high-dimensional and complex metagenomic data into 2-3D.
- It emphasizes on variation and shows strong patterns in a dataset. (w.r.t experimental factors)
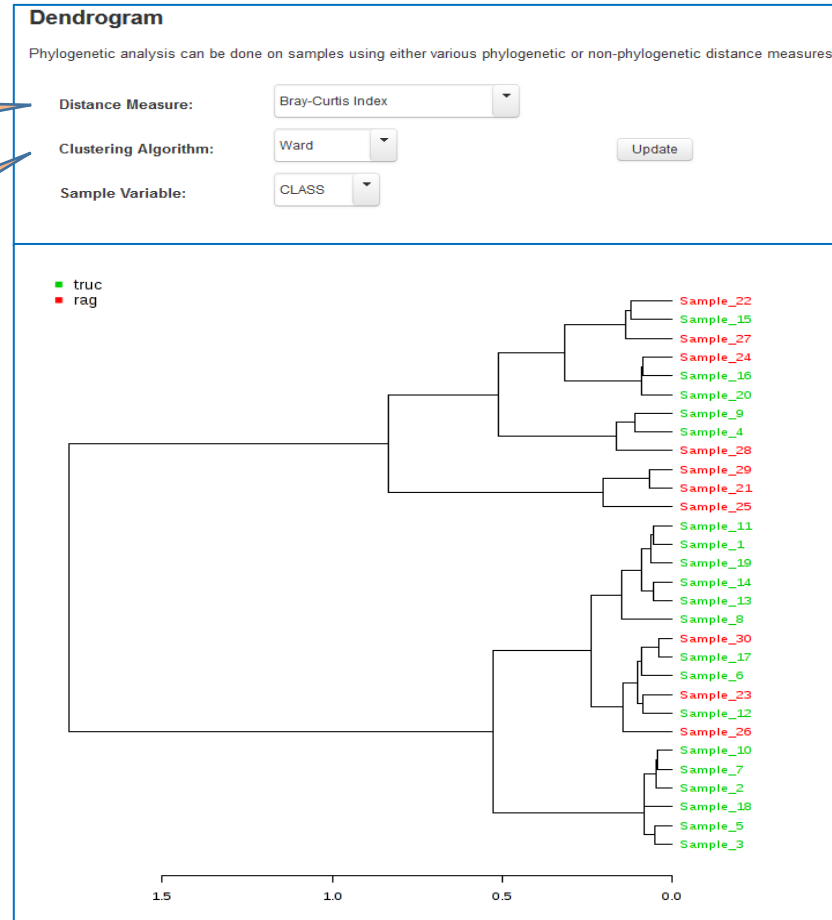
# B. Clustering Analysis



## 2. Heatmap
- Visualize the relative patterns of high-abundance features against a background of features that are mostly low-abundance or absent.
- Various distance and clustering methods supported.(both sample and feature-wise)
- Provides a summary of normalized user's data.

# B. Clustering Analysis



**Chose from different distance measure.**

**Chose from different clustering algorithm.**

## 3. Dendrogram

- Performs phylogenetic analysis on samples using ordination based distance measures.  (support for 5 most widely used)

# B. Clustering analysis



3 most common method supported for performing correlation analysis

## 4. Correlation analysis
• Helps in identifying biologically or biochemically meaningful relationship between features. (genes)

# B. Clustering analysis



**5. Pattern Search**

- Helps in identifying or search for a pattern based on correlation analysis on defined pattern.
- Pattern can be defined based on either feature (gene) of interest or based on predefined or custom profile of experimental factors.

# C. (a) Differential abundance analysis



**Univariate Statistical Comparisons**

Chose from parametric or non-parametric statistical tests

Experimental factor — Phenotype

Statistical method — Mann-Whitney/Kruskal-Wallis

Adjusted p-value cutoff — 0.05

Submit | Network Mapping ➤

Click here to visualize the differential genes in metabolic networks

Click on "Details" to see group-wise data distribution for each individual feature

| Name ⇕ | Pvalues ⇕ | FDR ⇕ | Statistics ⇕ | View |
|---|---|---|---|---|
| K00002 | 1.0825E-5 | 0.0012511 | 100.0 | 🖼 Details |
| K00012 | 1.0825E-5 | 0.0012511 | 100.0 | 🖼 Details |
| K00024 | 1.0825E-5 | 0.0012511 | 100.0 | 🖼 Details |
| K00018 | 1.0825E-5 | 0.0012511 | 0.0 | 🖼 Details |
| K00016 | 1.0825E-5 | 0.0012511 | 0.0 | 🖼 Details |
| K00021 | 1.0825E-5 | 0.0012511 | 0.0 | 🖼 Details |
| K00015 | 1.0825E-5 | 0.0012511 | 100.0 | 🖼 Details |
| K00052 | 1.4939E-4 | 0.0066798 | 0.0 | 🖼 Details |

Differential abundant genes (KO) are highlighted in orange color

## 1. Univariate Statistical Comparisons

- t-test/ANOVA (parametric) or Mann-Whitney/KW test (non-parametric) can be done.
- Depending upon no. of sample groups, statistical test is chosen from parametric or non parametric test options.
- P-values adjusted using **FDR** method.

# C. (a) Differential Abundance Analysis



**2. metagenomeSeq**
- R package which aims to detect differential abundant features in microbiome experiments with an explicit design.
- Accounts for **under-sampling** and **sparsity** in such data.
- Performs zero-inflated Gaussian fit (**fitZIG**) or fit-Feature (**fitFeature**) on data after normalizing the data through **cumulative sum scaling** (CSS) method (novel approach)
- fitFeature model is recommended over fitZIG for two groups comparison.
- Very sensitive and specific in nature.(fails with very low sample size)

# C. (a) Differential Abundance Analysis



**Chose from different Experimental factors**

**Click to perform Functional Enrichment Analysis on differentially abundant features**

**Click on "Details" to see group-wise data distribution for each individual feature**

Differential abundance analysis methods

| Experimental factor | Phenotype |
| Algorithm | EdgeR |
| Adjusted p-value cutoff | 0.05 |

Submit    Network Mapping

| Name ⇕ | log2FC ⇕ | logCPM ⇕ | Pvalues ⇕ | FDR ⇕ | View |
|---|---|---|---|---|---|
| K00029 | 12.699 | 12.077 | 4.5188E-62 | 2.9733E-59 | Details |
| K00051 | 13.296 | 11.601 | 7.3507E-62 | 2.9733E-59 | Details |
| K00030 | 13.101 | 11.38 | 5.5589E-53 | 1.499E-50 | Details |
| K00048 | -11.468 | 9.4128 | 1.0103E-49 | 2.0434E-47 | Details |
| K00045 | -10.343 | 7.8578 | 2.4431E-46 | 3.9529E-44 | Details |
| K00044 | -13.115 | 12.076 | 8.5303E-45 | 1.1502E-42 | Details |
| K00025 | -12.393 | 11.596 | 1.2923E-44 | 1.4935E-42 | Details |
| K00024 | -12.216 | 12.214 | 3.9404E-37 | 3.9847E-35 | Details |

## 3. EdgeR

* Developed for RNAseq data analysis.
* Powerful statistical method (outperforms others methods with appropriate data filtration and normalization techniques);
* By default, **RLE** (Relative Log Expression) normalization is performed on the data.

**Note**: If no significant gene will be identified using p-value cut-off, then top 500 genes based on their p-values will be used for network analysis.

# C. (a) Differential Abundance Analysis



**Chose from different Experimental factors**

**Click to perform Functional Enrichment Analysis on differentially abundant features**

**Click on "Details" to see group-wise data distribution for each individual feature**

Differential abundance analysis methods

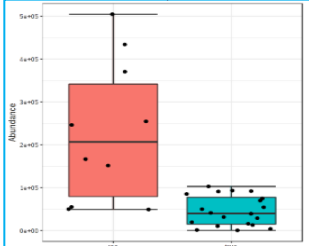| Experimental factor | Phenotype |
| Algorithm | DESeq2 |
| Adjusted p-value cutoff | 0.05 |

Submit    Network Mapping

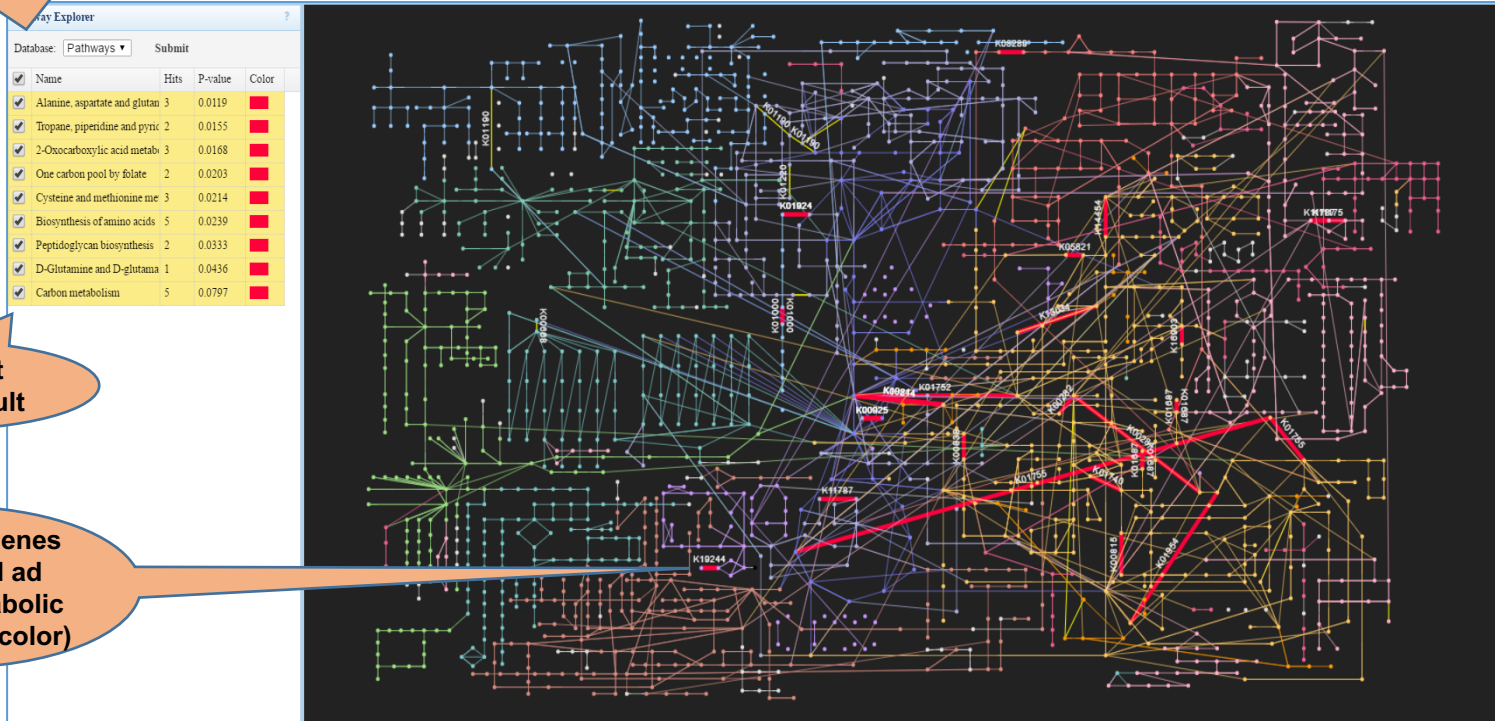| Name ⇕ | log2FC ⇕ | lfcSE ⇕ | Pvalues ⇕ | FDR ⇕ | View |
|---|---|---|---|---|---|
| K00045 | -9.9405 | 0.44313 | 1.8948E-111 | 1.5329E-108 | Details |
| K00029 | 10.93 | 0.51924 | 2.2519E-98 | 9.1089E-96 | Details |
| K00030 | 10.886 | 0.55785 | 8.2343E-85 | 2.2205E-82 | Details |
| K00048 | -10.14 | 0.52195 | 4.617E-84 | 9.3379E-82 | Details |
| K00051 | 10.788 | 0.57258 | 3.4896E-79 | 5.6462E-77 | Details |
| K00044 | -10.481 | 0.57848 | 2.3151E-73 | 3.1216E-71 | Details |
| K00024 | -9.8971 | 0.57073 | 2.3003E-67 | 2.6585E-65 | Details |
| K00025 | -9.9696 | 0.57633 | 4.8305E-67 | 4.8848E-65 | Details |

## 4. DESeq2
- Developed for RNAseq data analysis.
- Uses negative binomial generalized linear models to estimate **dispersion** and **logarithmic fold changes**.

**Note**: If no significant gene will be identified using p-value cut-off, then top 500 genes based on their p-values will be used for network analysis.

# C. (b) Network and Functional Enrichment Analysis



- Significant genes from differential analysis are mapped to KO IDs;
- Functional enrichment analysis is performed;( KEGG modules or pathways)
- The enriched pathways or modules can be interactively visualized within the metabolic networks.

# D. Biomarker analysis



Click here to visualize the differential genes in metabolic networks

Linear Discriminant Analysis (LDA) Effect Size (LEfSe) ?

Chose from different Experimental factors

Experimental factor          Phenotype ▼
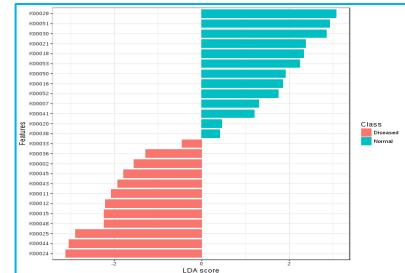
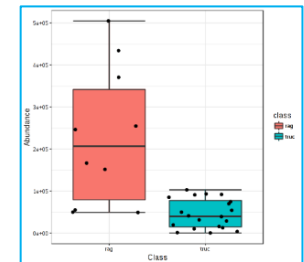Adjusted p-value cutoff       0.05                    Submit        Network Mapping ➔

Log LDA score                 1.0

Effect size (LDA score) of differential features

| Name ⇕ | Pvalues ⇕ | FDR ⇕ | Diseased ⇕ | Normal ⇕ | LDAscore ⇕ | View |
|---|---|---|---|---|---|---|
| K00052 | 1.2795E-4 | 0.0057752 | 0.10758 | 110.515 | 1.75 | 🖼 Details |
| K00051 | 1.2795E-4 | 0.0057752 | 0.128649 | 1680.31 | 2.92 | 🖼 Details |
| K00038 | 1.2978E-4 | 0.0057752 | 0.0178571 | 3.11296 | 0.406 | 🖼 Details |
| K00041 | 1.3988E-4 | 0.0057752 | 0.0775576 | 30.0144 | 1.2 | 🖼 Details |
| K00036 | 1.3988E-4 | 0.0057752 | 36.273 | 0.0443712 | -1.28 | 🖼 Details |
| K00050 | 1.4828E-4 | 0.0057752 | 0.0961246 | 158.97 | 1.91 | 🖼 Details |
| K00043 | 1.4828E-4 | 0.0057752 | 163.708 | 0.0618665 | -1.92 | 🖼 Details |

Click on "Details" to see group-wise data distribution for each individual feature



## 1. LEfSe

- compare the metagenomics (16S or shotgun) abundance profiles between samples in different state.
- performs a set of statistical tests for detecting differentially abundant features (**KW sum-rank test:** statistical significance) and biomarker discovery.(**Linear Discriminant analysis:** Effect Size)
- Network and functional enrichment analysis can also be performed on DE genes.

# D. Biomarker analysis



**User can choose from no. of trees to be used for classification**
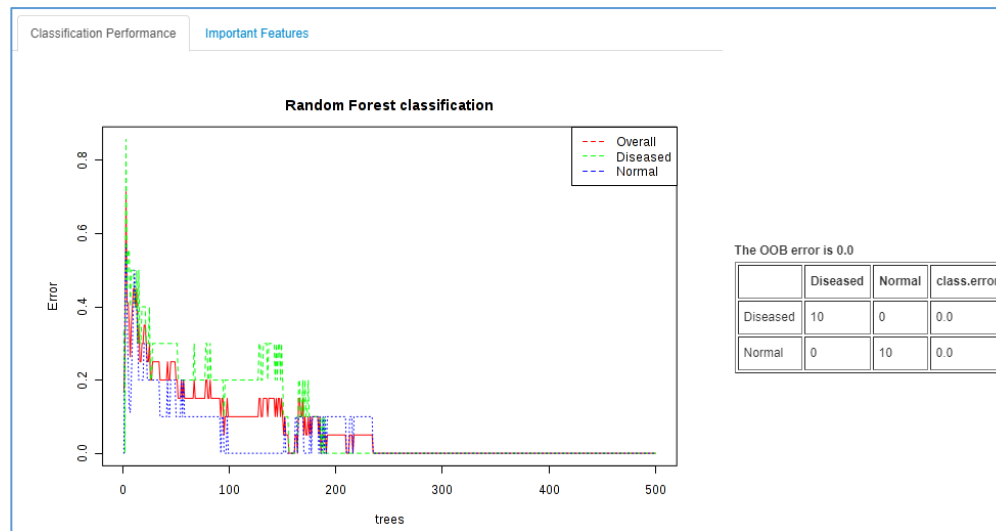
**No. of predictors for each node**

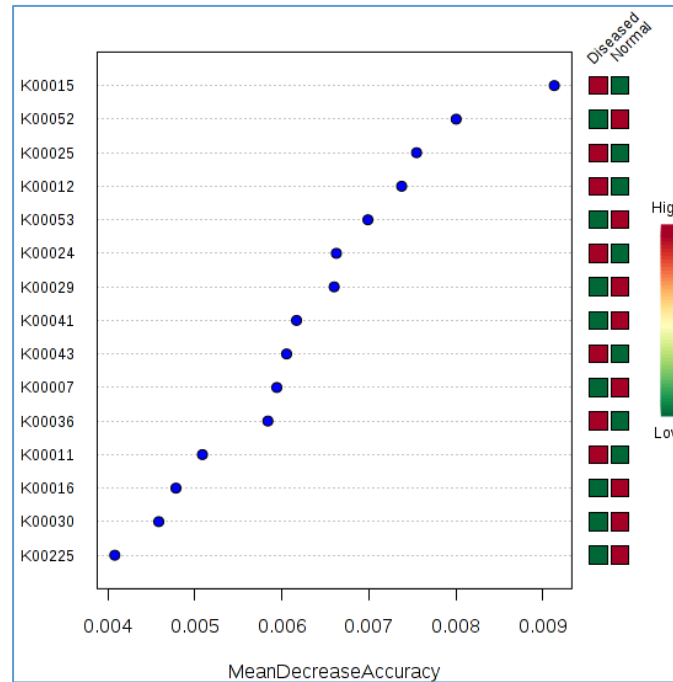## 2. Random forests

- Ensemble learning method used for classification, regression and other tasks.
- It operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) of the individual trees.
- Random forests correct for decision trees habit of overfitting to their training set.

# D. Biomarker analysis



**Most important features for classification of data into provided class groups**

## 2. Random Forest
- It provides estimates of what variables are important in the classification of data
- It computes proximities between pairs of cases that can be used in clustering, locating outliers, or give interesting views of the data

==END==