

Metabolomic Data Analysis with MetaboAnalyst 6.0

Name: guest2559499329594951229

January 7, 2025

1 Background

MSEA or Metabolite Set Enrichment Analysis is a way to identify biologically meaningful patterns that are significantly enriched in quantitative metabolomic data. In conventional approaches, metabolites are evaluated individually for their significance under conditions of study. Those compounds that have passed certain significance level are then combined to see if any meaningful patterns can be discerned. In contrast, MSEA directly investigates if a set of functionally related metabolites without the need to preselect compounds based on some arbitrary cut-off threshold. It has the potential to identify subtle but consistent changes among a group of related compounds, which may go undetected with the conventional approaches.

Essentially, MSEA is a metabolomic version of the popular GSEA (Gene Set Enrichment Analysis) software with its own collection of metabolite set libraries as well as an implementation of user-friendly web-interfaces. GSEA is widely used in genomics data analysis and has proven to be a powerful alternative to conventional approaches. For more information, please refer to the original paper by Subramanian A, and a nice review paper by Nam D, Kim SY.^{1, 2}

2 MSEA Overview

Metabolite set enrichment analysis consists of four steps - data input, data processing, data analysis, and results download. Different analysis procedures are performed based on different input types. In addition, users can also browse and search the metabolite set libraries as well as upload their self-defined metabolite sets for enrichment analysis. Users can also perform metabolite name mapping between a variety of compound names, synonyms, and major database identifiers.

3 Data Input

There are three enrichment analysis algorithms offered by MSEA. Accordingly, three different types of data inputs are required by these three approaches:

- A list of important compound names - entered as a one column data (*Over Representation Analysis (ORA)*);
- A single measured biofluid (urine, blood, CSF) sample- entered as tab separated two-column data with the first column for compound name, and the second for concentration values (*Single Sample Profiling (SSP)*);

¹Subramanian *Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles.*, Proc Natl Acad Sci USA. 2005 102(43): 15545-50

²Nam D, Kim SY. *Gene-set approach for expression pattern analysis*, Briefings in Bioinformatics. 2008 9(3): 189-197.

- A compound concentration table - entered as a comma separated (.csv) file with the each sample per row and each metabolite concentration per column. The first column is sample names and the second column for sample phenotype labels (*Quantitative Enrichment Analysis (QEA)*)

You selected Over Representation Analysis (ORA) which requires a list of compound names as input.

4 Data Process

The first step is to standardize the compound labels. It is an essential step since the compound labels will be subsequently compared with compounds contained in the metabolite set library. MSEA has a built-in tool to convert between compound common names, synonyms, identifiers used in HMDB ID, PubChem, ChEBI, BiGG, METLIN, KEGG, or Reactome. **Table 1** shows the conversion results. Note: *1* indicates exact match, *2* indicates approximate match, and *0* indicates no match. A text file contain the result can be found the downloaded file *name_map.csv*

Table 1: Result from Compound Name Mapping

	Query	Match	HMDB	PubChem	KEGG	SMILES
1	C02301	O-Acylcarnitine		5355	C02301	
2	C00062	L-Arginine	HMDB0000517	6322	C00062	<chem>N[C@@H](CCCN(C)=N)C(O)=O</chem>
3	C05526	S-Glutathionyl-L-cysteine	METPA0607		C05526	
4	C00307	Citicoline	HMDB0001413	13805	C00307	<chem>C[N+](C)(C)CCOP(O)(=O)OP(O)(=O)OC[C@H]1O</chem>
5	C00127	Oxidized glutathione	HMDB0003337	65359	C00127	<chem>N[C@@H](CCC(=O)N[C@@H](CSSC[C@H](NC(=O)C</chem>
6	C00020	Adenosine monophosphate	HMDB0000045	6083	C00020	<chem>NC1=C2N=CN([C@@H]3O[C@H](COP(O)(O)=O)[C</chem>
7	C00144	Guanosine monophosphate	HMDB0001397	6804	C00144	<chem>NC1=NC2=C(N=CN2[C@@H]2O[C@H](COP(O)(O)=</chem>
8	C05635	5-Hydroxyindoleacetic acid	HMDB0000763	1826	C05635	<chem>OC(=O)CC1=CNC2=C1C=C(O)C=C2</chem>
9	C20387	Biotin sulfone	HMDB0004818	21252323	C20387	<chem>[H][C@]12CS(=O)(=O)[C@@H](CCCCC(O)=O)[C@@</chem>
10	C01586	Hippuric acid	HMDB0000714	464	C01586	<chem>OC(=O)CNC(=O)C1=CC=CC=C1</chem>
11	C05551	Penicillin G	HMDB0015186	5904	C05551	<chem>[H][C@]12SC(C)(C)[C@@H](N1C(=O)[C@H]2NC(=O</chem>
12	C00319	Sphingosine	HMDB0000252	5280335	C00319	<chem>CCCCCCCCCCCCC\C=C\[C@@H](O)[C@@H](N)CO</chem>
13	C04230	CE(22:5(7Z,10Z,13Z,16Z,19Z))	HMDB0010375	24779458	C04230	<chem>CC\C=C/C\C=C/C\C=C/C\C=C/C\C=C/C\CCCC</chem>
14	C04100	NA	NA	NA	NA	NA
15	C02990	Palmitoylcarnitine	HMDB0000222	11953816	C02990	<chem>CCCCCCCCCCCCCCCC(=O)O[C@H](CC(O)=O)C[</chem>
16	C21484	LysoPE(18:0/0:0)	HMDB0011130	9547068	C21484	<chem>[H][C@@](O)(COC(=O)CCCCCCCCCCCCCCCCC)C</chem>
17	C05382	D-Sedoheptulose 7-phosphate	HMDB0001068	92042786	C05382	<chem>OC[C@]1(O)O[C@H](COP(O)(O)=O)[C@@H](O)[C@</chem>
18	C00052	Uridine diphosphategalactose	HMDB0000302	18068	C00052	<chem>OC[C@H]1O[C@H](OP(O)(=O)OP(O)(=O)OC[C@H]2O</chem>
19	C00105	Uridine 5'-monophosphate	HMDB0000288	6030	C00105	<chem>O[C@H]1[C@@H](O)[C@@H](O[C@@H]1COP(O)(O)=O</chem>
20	C00158	Citric acid	HMDB0000094	311	C00158	<chem>OC(=O)CC(O)(CC(O)=O)C(O)=O</chem>
21	C00360	Deoxyadenosine monophosphate	HMDB0000905	12599	C00360	<chem>NC1=NC=NC2=C1N=CN2[C@H]1C[C@H](O)[C@@H</chem>
22	C00157	Phosphatidylcholine			C00157	
23	C06525	Gentianine	HMDB0303030	354616	C06525	<chem>C=CC1=C2CCOC(=O)C2=CN=C1</chem>
24	C00550	2beta-Hydroxytestosterone	HMDB0012654	53481791	C00550	<chem>C[C@]12CCC3C(CCC4=CC(=O)[C@@H](O)C[C@]34</chem>
25	C00380	Cytosine	HMDB0000630	597	C00380	<chem>NC1=CC=NC(=O)N1</chem>
26	C02862	Butyrylcarnitine	HMDB0002013	213144	C02862	<chem>CCCC(=O)O[C@H](CC([O-])=O)C[N+](C)(C)C</chem>
27	C00864	Pantothenic acid	HMDB0000210	6613	C00864	<chem>CC(C)(CO)[C@@H](O)C(=O)NCCC(O)=O</chem>
28	C00029	Uridine diphosphate glucose	HMDB0000286	8629	C00029	<chem>OC[C@H]1O[C@H](OP(O)(=O)OP(O)(=O)OC[C@H]2O</chem>
29	C00147	Adenine	HMDB0000034	190	C00147	<chem>NC1=C2NC=NC2=NC=N1</chem>
30	C00284	Edetic Acid	HMDB0015109	6049	C00284	<chem>OC(=O)CN(CCN(CC(O)=O)CC(O)=O)CC(O)=O</chem>
31	C00570	CDP-ethanolamine	HMDB0001564	123727	C00570	<chem>NCCOP(O)(=O)OP(O)(=O)OC[C@H]1O[C@H]([C@H</chem>
32	C14550	4-Nonylphenol	HMDB0038982	1752	C14550	<chem>CCCCCCCCC1=CC=C(O)C=C1</chem>
33	C00946	Adenosine 2'-phosphate	HMDB0011617	53481006	C00946	<chem>NC1=NC=NC2=C1N=CN2C1O[C@H](CO)[C@@H](C</chem>

The second step is to check concentration values. For SSP analysis, the concentration must be measured in *umol* for blood and CSF samples. The urinary concentrations must be first converted to *umol/mmol_creatinine* in order to compare with reported concentrations in literature. No missing or negative values are allowed in SSP analysis. The concentration data for QEA analysis is more flexible. Users can upload either the original concentration data or normalized data. Missing or negative values are allowed (coded as *NA*) for QEA.

5 Selection of Metabolite Set Library

Before proceeding to enrichment analysis, a metabolite set library has to be chosen. There are seven built-in libraries offered by MSEA:

- Metabolic pathway associated metabolite sets (*currently contains 99 entries*);
- Disease associated metabolite sets (reported in blood) (*currently contains 344 entries*);
- Disease associated metabolite sets (reported in urine) (*currently contains 384 entries*);
- Disease associated metabolite sets (reported in CSF) (*currently contains 166 entries*);
- Metabolite sets associated with SNPs (*currently contains 4598 entries*);
- Predicted metabolite sets based on computational enzyme knockout model (*currently contains 912 entries*);
- Metabolite sets based on locations (*currently contains 73 entries*);
- Drug pathway associated metabolite sets (*currently contains 461 entries*);

In addition, MSEA also allows user-defined metabolite sets to be uploaded to perform enrichment analysis on arbitrary groups of compounds which researchers want to test. The metabolite set library is simply a two-column comma separated text file with the first column for metabolite set names and the second column for its compound names (**must use HMDB compound name**) separated by "; ". Please note, the built-in libraries are mainly from human studies. The functional grouping of metabolites may not be valid. Therefore, for data from subjects other than human being, users are suggested to upload their self-defined metabolite set libraries for enrichment analysis.

6 Enrichment Analysis

Over Representation Analysis (ORA) is performed when a list of compound names is provided. The list of compound list can be obtained through conventional feature selection methods, or from a clustering algorithm, or from the compounds with abnormal concentrations detected in SSP, to investigate if some biologically meaningful patterns can be identified.

ORA was implemented using the *hypergeometric test* to evaluate whether a particular metabolite set is represented more than expected by chance within the given compound list. One-tailed p values are provided after adjusting for multiple testing. **Figure 2** below summarizes the result.

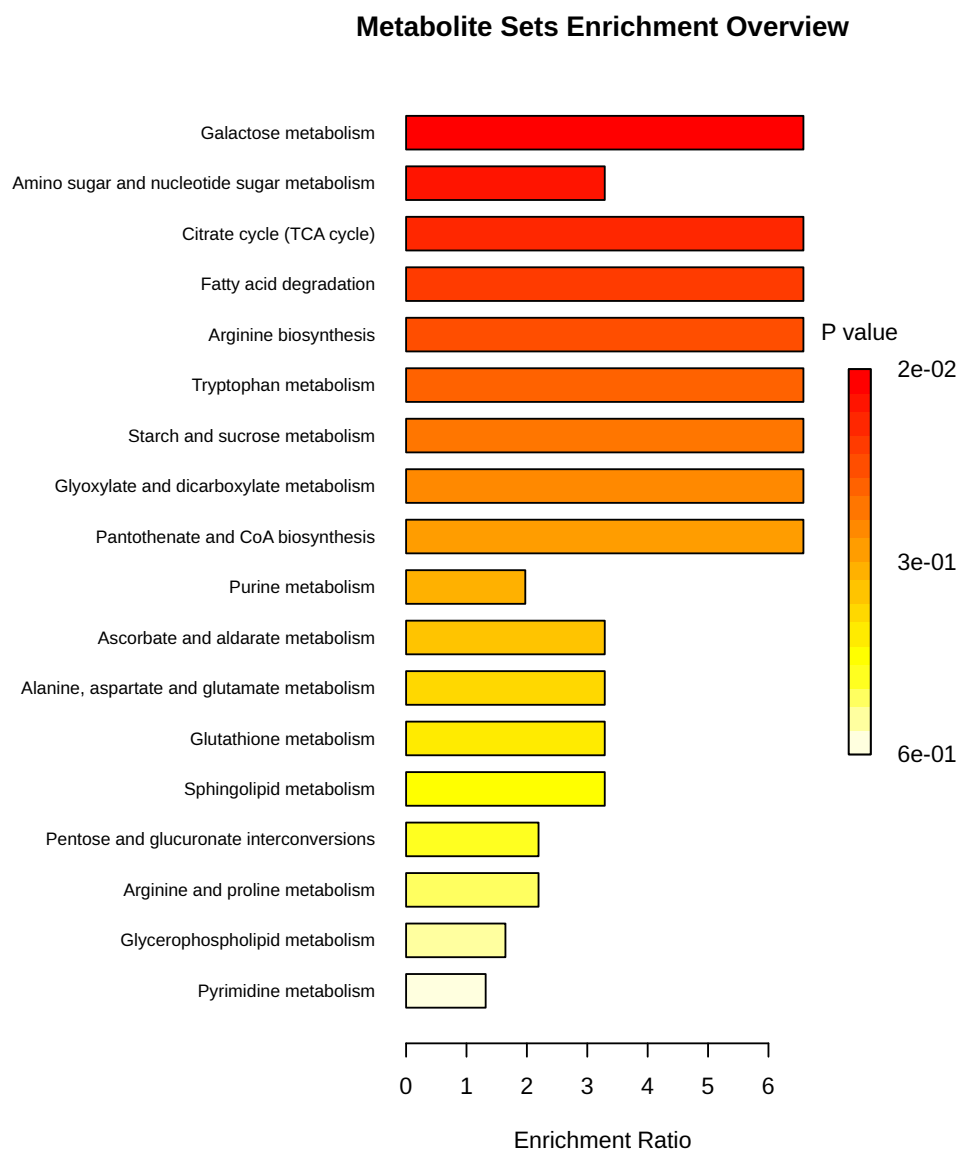


Figure 1: Summary Plot for Over Representation Analysis (ORA)

Table 2: Result from Over Representation Analysis

	total	expected	hits	Raw p	Holm p	FDR
Galactose metabolism	2	0.30	2	2.23E-02	1.00E+00	1.00E+00
Amino sugar and nucleotide sugar metabolism	4	0.61	2	1.10E-01	1.00E+00	1.00E+00
Citrate cycle (TCA cycle)	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Fatty acid degradation	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Arginine biosynthesis	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Tryptophan metabolism	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Starch and sucrose metabolism	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Glyoxylate and dicarboxylate metabolism	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Pantothenate and CoA biosynthesis	1	0.15	1	1.52E-01	1.00E+00	1.00E+00
Purine metabolism	10	1.52	3	1.79E-01	1.00E+00	1.00E+00
Ascorbate and aldarate metabolism	2	0.30	1	2.82E-01	1.00E+00	1.00E+00
Alanine, aspartate and glutamate metabolism	2	0.30	1	2.82E-01	1.00E+00	1.00E+00
Glutathione metabolism	2	0.30	1	2.82E-01	1.00E+00	1.00E+00
Sphingolipid metabolism	2	0.30	1	2.82E-01	1.00E+00	1.00E+00
Pentose and glucuronate interconversions	3	0.46	1	3.92E-01	1.00E+00	1.00E+00
Arginine and proline metabolism	3	0.46	1	3.92E-01	1.00E+00	1.00E+00
Glycerophospholipid metabolism	4	0.61	1	4.86E-01	1.00E+00	1.00E+00
Pyrimidine metabolism	5	0.76	1	5.66E-01	1.00E+00	1.00E+00

7 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"conc\", \"msetora\", FALSE)"
[2] "compd.vec<-c(\"C02301\", \"C00062\", \"C05526\", \"C00307\", \"C00127\", \"C00020\", \"C00144\", \"C05
[3] "mSet<-Setup.MapData(mSet, compd.vec);"
[4] "mSet<-CrossReferencing(mSet, \"kegg\");"
[5] "mSet<-CreateMappingResultTable(mSet)"
[6] "mSet<-Setup.HMDBReferenceMetabolome(mSet, \"BP.txt\");"
[7] "mSet<-SetMetabolomeFilter(mSet, T);"
[8] "mSet<-SetCurrentMsetLib(mSet, \"kegg_pathway\", 2);"
[9] "mSet<-CalculateHyperScore(mSet)"
[10] "mSet<-PlotORA(mSet, \"ora_0\", \"net\", \"png\", 72, width=NA)"
[11] "mSet<-PlotEnrichDotPlot(mSet, \"ora\", \"ora_dot_0\", \"png\", 72, width=NA)"
[12] "mSet<-CalculateHyperScore(mSet)"
[13] "mSet<-PlotORA(mSet, \"ora_1\", \"net\", \"png\", 72, width=NA)"
[14] "mSet<-PlotEnrichDotPlot(mSet, \"ora\", \"ora_dot_1\", \"png\", 72, width=NA)"
[15] "mSet<-SaveTransformedData(mSet)"
[16] "mSet<-PreparePDFReport(mSet, \"guest2559499329594951229\")\n"
```

The report was generated on Tue Jan 7 09:30:07 2025 with R version 4.3.2 (2023-10-31), OS system: Linux.