

# Metabolomic Data Analysis with MetaboAnalyst 5.0

Name: guest6329621867019318963

November 24, 2022

## 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.<sup>1, 2</sup>

## 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)*);

---

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

<sup>2</sup>Nam 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

	Query	Match	HMDB	PubChem	KEGG	SMILES
1	HMDB0000714	Hippuric acid	HMDB0000714	464	C01586	C1=CC=C(C(=C1)C(=O)N
2	HMDB0000355	3-Hydroxymethylglutaric acid	HMDB0000355	1662	C03761	CC(CC(=O)O)(CC(=O)O)O
3	HMDB0011686	p-Cresol glucuronide	HMDB0011686	154035		CC1=CC=C(C(=C1)O)[C@H]
4	HMDB0000072	cis-Aconitic acid	HMDB0000072	643757	C00417	C(/C(=C/C(=O)O)/C(=O)O)
5	HMDB0006116	3-Hydroxyhippuric acid	HMDB0006116	450268		C1=CC(=CC(=C1)O)C(=O)O
6	HMDB0002643	3-(3-Hydroxyphenyl)-3-hydroxypropanoic acid	HMDB0002643	102959		C1=CC(=CC(=C1)O)C(=O)O
7	HMDB0341278	NA	NA	NA	NA	NA
8	HMDB0003099	1-Methyluric acid	HMDB0003099	69726	C16359	CN1C(=O)C2=C(NC(=O)O)C(=O)O
9	HMDB0000193	Isocitric acid	HMDB0000193	1198	C00311	C(C(C(C(=O)O)O)C(=O)O)O
10	HMDB0006344	Alpha-N-Phenylacetyl-L-glutamine	HMDB0006344	92258	C04148	C1=CC=C(C(=C1)CC(=O)O)N
11	HMDB0000152	Gentisic acid	HMDB0000152	3469	C00628	C1=CC(=C(C(=C1)O)C(=O)O)O
12	HMDB0013324	2-Octenoylcarnitine	HMDB0013324	53481667		CCCCC/C=C/C(=O)O[C@H]
13	HMDB0000440	3-Hydroxyphenylacetic acid	HMDB0000440	12122	C05593	C1=CC(=CC(=C1)O)CC(=O)O
14	HMDB0000875	Trigonelline	HMDB0000875	5570	C01004	C[N+](C)(C)C(=O)O
15	HMDB0002721	1-Methylinosine	HMDB0002721	65095		CN1C=NC2=C(C1=O)N=CN2
16	HMDB0006275	Dopamine 3-O-sulfate	HMDB0006275	122136	C13690	C1=CC(=C(C(=C1)CCN)OS(=O)(=O)O
17	HMDB0000893	Suberic acid	HMDB0000893	10457	C08278	C(CCCC(=O)O)CCC(=O)O
18	HMDB0000912	Succinyladenosine	HMDB0000912	20849086		C1=NC2=C(C(=N1)N[C@H](C2)C(=O)O)O
19	HMDB0000736	Isobutyryl-L-carnitine	HMDB0000736	168379		CC(C)C(=O)OC(CC(=O)O)N
20	HMDB0240751	NA	NA	NA	NA	NA
21	HMDB0142137	NA	NA	NA	NA	NA
22	HMDB0000730	Isobutyrylglycine	HMDB0000730	10855600		CC(C)C(=O)NCC(=O)O
23	HMDB0000812	N-Acetyl-L-aspartic acid	HMDB0000812	65065	C01042	CC(=O)N[C@@H](CC(=O)O)C(=O)O
24	HMDB0003464	4-Guanidinobutanoic acid	HMDB0003464	500	C01035	C(CC(=O)O)CN=C(N)N
25	HMDB0029992	Tetrahydropentoxylene	HMDB0029992	53481442		C1C(NC(C2=C1C3=CC=C(C2)O)O)O
26	HMDB0000925	Trimethylamine N-oxide	HMDB0000925	1145	C01104	C[N+](C)(C)[O-]
27	HMDB0011103	1,7-Dimethyluric acid	HMDB0011103	91611	C16356	CN1C2=C(C(NC1=O)NC(=O)O)C(=O)O
28	HMDB0001411	Cotinine N-oxide	HMDB0001411	9815514		CN1[C@@H](CCC1=O)C2=CC=CC=C2
29	HMDB0003072	Quinic acid	HMDB0003072	6508	C00296	OC1C[C@@H](O)(C[C@@H](O)C(=O)O)C(=O)O
30	HMDB0000512	N-Acetyl-L-phenylalanine	HMDB0000512	74839	C03519	CC(=O)N[C@@H](CC1=CC=C(C=C1)O)C(=O)O
31	HMDB0003331	1-Methyladenosine	HMDB0003331	27476	C02494	CN1C=NC2=C(C1=N)N=CN2
32	HMDB0013678	4-Hydroxyhippuric acid	HMDB0013678	151012		C1=CC(=CC(=C1)O)CC(=O)O
33	HMDB0061112	3-Carboxy-4-methyl-5-propyl-2-furanpropionic acid	HMDB0061112	123979		CCCC1=C(C(C(C1)O)O)O
34	HMDB0013676	2,6-Dihydroxybenzoic acid	HMDB0013676	9338	C21298	C1=CC(=C(C(=C1)O)C(=O)O)O
35	HMDB0001406	Niacinamide	HMDB0001406	936	C00153	C1=CC(=CN=C1)C(=O)N
36	HMDB0004824	N2,N2-Dimethylguanosine	HMDB0004824	92919		CN(C)C1=NC(=O)C2=C(C(=N1)N)C(=O)O
37	HMDB60001	NA	NA	NA	NA	NA
38	HMDB0002802	Cortisone	HMDB0002802	225609	C00762	C[C@]12CCC(=O)C=C1CCC2=O
39	HMDB0001563	1-Methylguanosine	HMDB0001563	96373	C04545	CN1C(=O)C2=C(C(N=C1N)N)C(=O)O
40	HMDB0001297	Norcotinine	HMDB0001297	413		C1CC(=O)NC1C2=CC=CC=C2
41	HMDB0001713	m-Coumaric acid	HMDB0001713	637541	C12621	C1=CC(=CC(=C1)O)/C=C/C
42	HMDB0013713	N-acetyltryptophan	HMDB0013713	700653		[H][C@@H](CC1=CC=CC=C1)C(=O)O
43	HMDB0000092	Dimethylglycine	HMDB0000092	673	C01026	CN(C)CC(=O)O
44	HMDB0000138	Glycocholic acid	HMDB0000138	23617285	C01921	C[C@@H](CCC(=O)O)NCC(=O)O
45	HMDB0000159	L-Phenylalanine	HMDB0000159	6140	C00079	C1=CC=C(C(=C1)C[C@@H](O)C(=O)O)O
46	HMDB0000262	Thymine	HMDB0000262	1135	C00178	CC1=CNC(=O)NC1=O
47	HMDB0000162	L-Proline	HMDB0000162	145742	C00148	C1C[C@H](NC1)C(=O)O
48	HMDB0028933	Leucyl-Leucine	HMDB0028933	76807	C11332	CC(C)CC(N)C(=O)NC(CCC(=O)O)C(=O)O
49	HMDB00191	L-Aspartic acid	HMDB0000191	5960	C00049	C([C@@H](C(=O)O)N)C(=O)O
50	HMDB0002035	4-Hydroxycinnamic acid	HMDB0002035	637542	C00811	C1=CC(=CC(=C1)/C=C/C(=O)O)O
51	HMDB0000133	Guanosine	HMDB0000133	6802	C00387	C1=NC2=C(C(N1[C@H]3[C@H](O)C(=O)O)O)O
52	HMDB0000020	p-Hydroxyphenylacetic acid	HMDB0000020	127	C00642	C1=CC(=CC(=C1)O)CC(=O)O
53	HMDB0002024	Imidazoleacetic acid	HMDB0002024	96215	C02835	C1=C(NC(=N1)CC(=O)O)O
54	HMDB0000661	Glutaric acid	HMDB0000661	743	C00489	C(CC(=O)O)CC(=O)O
55	HMDB0002432	Sumiki's acid	HMDB0002432	80642	C20448	C1=C(OC(=C1)C(=O)O)O
56	HMDB0001860	Paraxanthine	HMDB0001860	4687	C13747	CN1C=NC2=C1C(=O)N(C)C(=O)O
57	HMDB0000732	Hydroxykynurenine	HMDB0000732	89	C02794	C1=CC(=C(C(=C1)O)N)C(=O)O
58	HMDB0000669	Ortho-Hydroxyphenylacetic acid	HMDB0000669	11970	C05852	C1=CC=C(C(=C1)CC(=O)O)O

59	HMDB0000230	N-Acetylneuraminic acid	HMDB0000230	445063	C19910	<chem>CC(=O)N[C@@H]1[C@H](O)</chem>
60	HMDB0000491	3-Methyl-2-oxovaleric acid	HMDB0000491	47	C00671	<chem>CCC(C)C(=O)C(=O)O</chem>
61	HMDB0001987	2-Hydroxy-2-methylbutyric acid	HMDB0001987	95433		<chem>CCC(C)(C(=O)O)O</chem>
62	HMDB0000630	Cytosine	HMDB0000630	597	C00380	<chem>C1=C(NC(=O)N=C1)N</chem>
63	HMDB0000407	2-Hydroxy-3-methylbutyric acid	HMDB0000407	99823		<chem>CC(C)C(C(=O)O)O</chem>
64	HMDB0000842	Quinaldic acid	HMDB0000842	7124	C06325	<chem>C1=CC=C2C(=C1)C=CC(=O)N2</chem>
65	HMDB0001046	Cotinine	HMDB0001046	408		<chem>CN1C(CCC1=O)C2=CN=C(C2)C</chem>
66	HMDB0000118	Homovanillic acid	HMDB0000118	1738	C05582	<chem>COC1=C(C=CC(=C1)CC(=O)O)O</chem>
67	HMDB0004827	Proline betaine	HMDB0004827	7016563	C10172	<chem>C[N+](C)(CCC[C@H]1C(=O)N1)C</chem>
68	HMDB0000172	L-Isoleucine	HMDB0000172	6306	C00407	<chem>CC[C@H](C)[C@H](C(=O)O)C</chem>
69	HMDB0005923	N4-Acetylcytidine	HMDB0005923	107461		<chem>CC(=O)NC1=NC(=O)N(C(=O)N1)C</chem>
70	HMDB0002730	Nicotinamide N-oxide	HMDB0002730	72661		<chem>C1=CC(=C[N+](=C1)[O-])C</chem>
71	HMDB0000073	Dopamine	HMDB0000073	681	C03758	<chem>C1=CC(=C(C=C1CCN)O)C</chem>
72	HMDB0061384	NA	NA	NA	NA	NA
73	HMDB0000158	L-Tyrosine	HMDB0000158	6057	C00082	<chem>C1=CC(=CC=C1C[C@@H](O)C</chem>
74	HMDB0244966	NA	NA	NA	NA	NA
75	HMDB0000201	L-Acetylcarnitine	HMDB0000201	7045767	C02571	<chem>CC(=O)OC(CC(=O)O)[O-]C</chem>
76	HMDB0000177	L-Histidine	HMDB0000177	6274	C00135	<chem>C1=C(NC=N1)C[C@@H](O)C</chem>
77	HMDB0012296	Trimethylaminoacetone	HMDB0012296	151806		<chem>CC(=O)C[N+](C)(C)C</chem>
78	HMDB0000391	7-Ketodeoxycholic acid	HMDB0000391	188292	C04643	<chem>C[C@H](CCC(=O)O)[C@H](O)C</chem>
79	HMDB0000050	Adenosine	HMDB0000050	60961	C00212	<chem>C1=NC2=C(C(=N1)N)N=C(N2)C</chem>
80	HMDB0061684	N-Acetyl isoleucine	HMDB0061684	7036275		<chem>CC[C@H](C)[C@H](NC(=O)C)C</chem>
81	HMDB0001434	3-Methoxytyrosine	HMDB0001434	1670		<chem>COC1=C(C=CC(=C1)CC(=O)O)O</chem>
82	HMDB0001476	3-Hydroxyanthranilic acid	HMDB0001476	86	C00632	<chem>C1=CC(=C(C(=C1)O)N)C</chem>
83	HMDB0010319	Inodxyl glucuronide	HMDB0010319	2733785	C03033	<chem>C1=CC=C2C(=C1)C(=CN2)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.

## Metabolite Sets Enrichment Overview

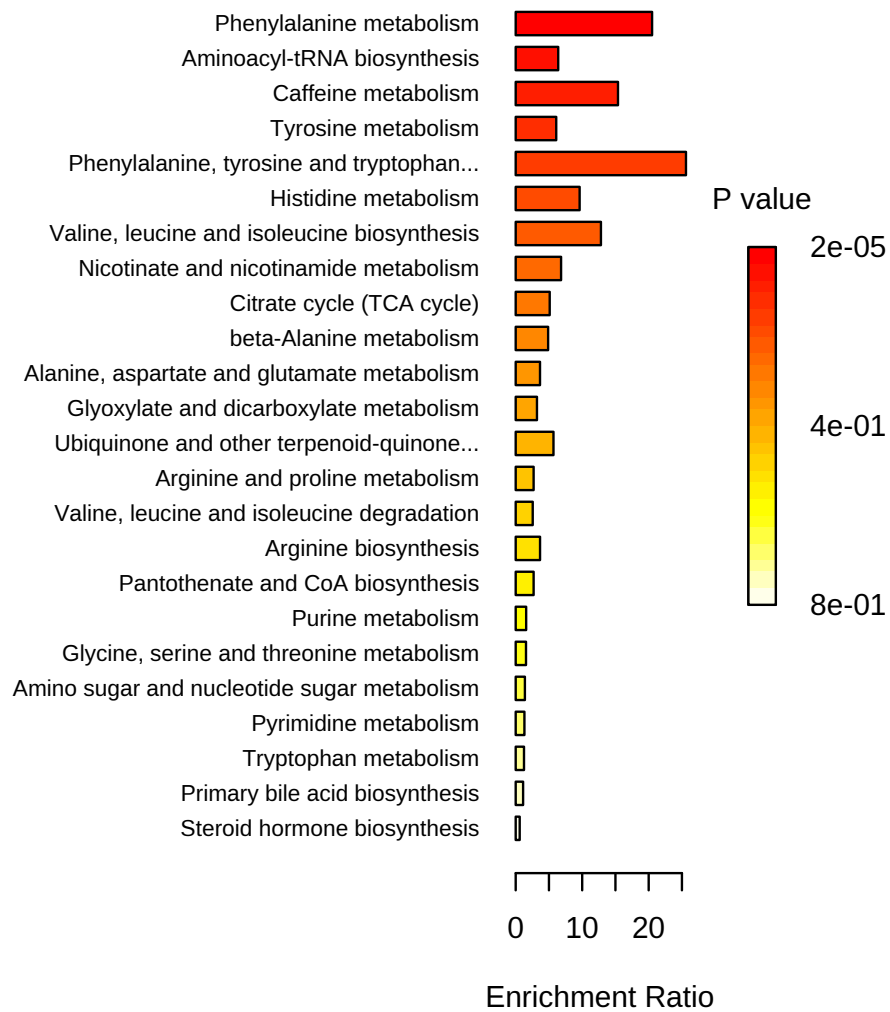


Figure 1: Summary Plot for Over Representation Analysis (ORA)

Table 2: Result from Over Representation Analysis

	total	expected	hits	Raw p	Holm p	FDR
Phenylalanine metabolism	10	0.20	4	2.29E-05	1.93E-03	1.93E-03
Aminoacyl-tRNA biosynthesis	48	0.94	6	2.28E-04	1.89E-02	9.57E-03
Caffeine metabolism	10	0.20	3	7.36E-04	6.04E-02	2.06E-02
Tyrosine metabolism	42	0.82	5	1.03E-03	8.34E-02	2.16E-02
Phenylalanine, tyrosine and tryptophan biosynthesis	4	0.08	2	2.16E-03	1.73E-01	3.63E-02
Histidine metabolism	16	0.31	3	3.17E-03	2.51E-01	4.44E-02
Valine, leucine and isoleucine biosynthesis	8	0.16	2	9.60E-03	7.49E-01	1.15E-01
Nicotinate and nicotinamide metabolism	15	0.29	2	3.31E-02	1.00E+00	3.47E-01
Citrate cycle (TCA cycle)	20	0.39	2	5.64E-02	1.00E+00	5.17E-01
beta-Alanine metabolism	21	0.41	2	6.15E-02	1.00E+00	5.17E-01
Alanine, aspartate and glutamate metabolism	28	0.55	2	1.02E-01	1.00E+00	7.78E-01
Glyoxylate and dicarboxylate metabolism	32	0.62	2	1.28E-01	1.00E+00	8.93E-01
Ubiquinone and other terpenoid-quinone biosynthesis	9	0.18	1	1.63E-01	1.00E+00	1.00E+00
Arginine and proline metabolism	38	0.74	2	1.68E-01	1.00E+00	1.00E+00
Valine, leucine and isoleucine degradation	40	0.78	2	1.83E-01	1.00E+00	1.00E+00
Arginine biosynthesis	14	0.27	1	2.42E-01	1.00E+00	1.00E+00
Pantothenate and CoA biosynthesis	19	0.37	1	3.14E-01	1.00E+00	1.00E+00
Purine metabolism	65	1.27	2	3.65E-01	1.00E+00	1.00E+00
Glycine, serine and threonine metabolism	33	0.65	1	4.82E-01	1.00E+00	1.00E+00
Amino sugar and nucleotide sugar metabolism	37	0.72	1	5.22E-01	1.00E+00	1.00E+00
Pyrimidine metabolism	39	0.76	1	5.41E-01	1.00E+00	1.00E+00
Tryptophan metabolism	41	0.80	1	5.59E-01	1.00E+00	1.00E+00
Primary bile acid biosynthesis	46	0.90	1	6.02E-01	1.00E+00	1.00E+00
Steroid hormone biosynthesis	85	1.66	1	8.22E-01	1.00E+00	1.00E+00

## 7 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"conc\", \"msetora\", FALSE)"
[2] "compd.vec<-c(\"HMDB0000714\", \"HMDB0000355\", \"HMDB0011686\", \"HMDB0000072\", \"HMDB0006116\", \"
[3] "mSet<-Setup.MapData(mSet, compd.vec);"
[4] "mSet<-CrossReferencing(mSet, \"hmdb\");"
[5] "mSet<-CreateMappingResultTable(mSet)"
[6] "mSet<-SetMetabolomeFilter(mSet, F);"
[7] "mSet<-SetCurrentMsetLib(mSet, \"kegg_pathway\", 2);"
[8] "mSet<-CalculateHyperScore(mSet)"
[9] "mSet<-PlotORA(mSet, \"ora_0\", \"net\", \"png\", 72, width=NA)"
[10] "mSet<-PlotEnrichDotPlot(mSet, \"ora\", \"ora_dot_0\", \"png\", 72, width=NA)"
[11] "mSet<-CalculateHyperScore(mSet)"
[12] "mSet<-PlotORA(mSet, \"ora_1\", \"net\", \"png\", 72, width=NA)"
[13] "mSet<-PlotEnrichDotPlot(mSet, \"ora\", \"ora_dot_1\", \"png\", 72, width=NA)"
[14] "mSet<-SaveTransformedData(mSet)"
[15] "mSet<-PreparePDFReport(mSet, \"guest6329621867019318963\")\n"
[16] "mSet<-PreparePDFReport(mSet, \"guest6329621867019318963\")\n"
```

---

The report was generated on Thu Nov 24 07:17:49 2022 with R version 4.2.2 (2022-10-31), OS system:  
Linux, version: -Ubuntu SMP Thu Oct 13 08:03:55 UTC 2022 .