

Metabolomic Data Analysis with MetaboAnalyst 6.0

Name: guest2982651237338019278

January 8, 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 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

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 Quantitative Enrichment Analysis (QEA) which requires a concentration table. This is the most common data format generated from quantitative metabolomics studies. The phenotype label can be categorical (binary or multi-class) or continuous.

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 C

	Query	Match	HMDB	PubChem	KEGG	SMILES
1	C00003	NAD	HMDB0000902	5892	C00003	NC(=O)C1=C[N+](=CC=C1)[C@@H]1O[C@H]
2	C00008	ADP	HMDB0001341	6022	C00008	NC1=NC=NC2=C1N=CN2[C@@H]1O[C@H]
3	C00015	Uridine 5'-diphosphate	HMDB0000295	6031	C00015	O[C@H]1[C@@H](O)[C@@H](O)[C@@H]1CO
4	C00016	FAD	HMDB0001248	643975	C00016	CC1=CC2=C(C=C1C)N(C[C@H](O)[C@H](O))
5	C00019	S-Adenosylmethionine	HMDB0001185	34756	C00019	C[S+](CC[C@H](N)C(O)=O)[C@H]1O[C@H]
6	C00020	Adenosine monophosphate	HMDB0000045	6083	C00020	NC1=C2N=CN([C@H]3O[C@H](COP(O)(O)
7	C00029	Uridine diphosphate glucose	HMDB0000286	8629	C00029	OC[C@H]1O[C@H](OP(O)(=O)OP(O)(=O)O
8	C00043	Uridine diphosphate-N-acetylglucosamine	HMDB0000290	445675	C00043	CC(=O)N[C@@H]1[C@@H](O)[C@@H](O)[C@H]
9	C00051	Glutathione	HMDB0000125	124886	C00051	N[C@@H](CCC(=O)N)[C@@H](CS)C(=O)N
10	C00052	Uridine diphosphate galactose	HMDB0000302	18068	C00052	OC[C@H]1O[C@H](OP(O)(=O)OP(O)(=O)O
11	C00055	Cytidine monophosphate	HMDB0000095	6131	C00055	NC1=NC(=O)N(C=C1)[C@@H]1O[C@H](COP
12	C00062	L-Arginine	HMDB00000517	6322	C00062	N[C@@H](CCCN(C(N)=N)C(O)=O)
13	C00082	L-Tyrosine	HMDB0000158	6057	C00082	N[C@@H](CC1=CC=C(O)C=C1)C(O)=O
14	C00103	Glucose 1-phosphate	HMDB0001586	65533	C00103	OC[C@H]1O[C@H](OP(O)(O)=O)[C@H](O)
15	C00105	Uridine 5'-monophosphate	HMDB0000288	6030	C00105	O[C@H]1[C@@H](O)[C@@H](O)[C@@H]1CO
16	C00120	Biotin	HMDB0000030	171548	C00120	[H][C@]12CS[C@@H](CCCCC(O)=O)[C@@H]
17	C00127	Oxidized glutathione	HMDB0003337	65359	C00127	N[C@@H](CCC(=O)N)[C@@H](CSSC[C@H](
18	C00144	Guanosine monophosphate	HMDB0001397	6804	C00144	NC1=NC2=C(N=CN2)[C@H]2O[C@H](CO
19	C00147	Adenine	HMDB0000034	190	C00147	NC1=C2NC=NC2=NC=N1
20	C00157	Phosphatidylcholine			C00157	
21	C00158	Citric acid	HMDB0000094	311	C00158	OC(=O)CC(O)(CC(O)=O)C(O)=O
22	C00167	Uridine diphosphate glucuronic acid	HMDB0000935	17473	C00167	O[C@@H]1[C@@H](COP(O)(=O)OP(O)(=O)O
23	C00170	5'-Methylthioadenosine	HMDB0001173	439176	C00170	CSC[C@H]1O[C@H](C[C@H](O)[C@@H]1O)N
24	C00199	D-Ribulose 5-phosphate	HMDB0000618	439184	C00199	OCC(=O)[C@H](O)[C@@H](O)COP(O)(O)=O
25	C00242	Guanine	HMDB0000132	764	C00242	NC1=NC(=O)C2=C(N1)N=CN2
26	C00262	Hypoxanthine	HMDB0000157	790	C00262	OC1=NC=NC2=C1NC=N2
27	C00294	Inosine	HMDB0000195	6021	C00294	OC[C@H]1O[C@H](C[C@H](O)[C@@H]1O)N1
28	C00299	Uridine	HMDB0000296	6029	C00299	OC[C@H]1O[C@H](C[C@H](O)[C@@H]1O)N1
29	C00319	Sphingosine	HMDB0000252	5280335	C00319	CCCCCCCCCCCCC\C=C/[C@@H](O)[C@@H]
30	C00325	GDP-L-fucose	HMDB0001095	439211	C00325	C[C@@H]1OC(OP(O)(=O)OP(O)(=O)OC[C@H]
31	C00350	PE(14:0/20:1(11Z))	HMDB0000834	52924120	C00350	[H][C@@](COC(=O)CCCCCCCCCCCCC)(C
32	C00360	Deoxyadenosine monophosphate	HMDB0000905	12599	C00360	NC1=NC=NC2=C1N=CN2[C@H]1C[C@H](COP
33	C00362	2'-Deoxyguanosine 5'-monophosphate	HMDB0001044	65059	C00362	NC1=NC2=C(N=CN2)[C@H]2C[C@H](O)[C@H]
34	C00364	5-Thymidylic acid	HMDB0001227	9700	C00364	CC1=CN(C[C@H]2C[C@H](O)[C@@H](COP(O)(
35	C00378	Thiamine	HMDB0000235	1130	C00378	CC1=C(CC(O)S[C=[N+]1CC1=CN=C(C)N=
36	C00380	Cytosine	HMDB0000630	597	C00380	NC1=CC=NC(=O)N1
37	C00385	Xanthine	HMDB0000292	1188	C00385	O=C1NC2=C(NC=N2)C(=O)N1
38	C00387	Guanosine	HMDB0000133	6802	C00387	NC1=NC2=C(N=CN2)[C@@H]2O[C@H](CO
39	C00463	Indole	HMDB0000738	798	C00463	N1C=CC2=C1C=CC=C2
40	C00487	L-Carnitine	HMDB0000062	10917	C00487	C[N+](C)(C)[C@H](O)CC(O)=O
41	C00491	L-Cystine	HMDB0000192	67678	C00491	N[C@@H](CSSC[C@H](N)C(O)=O)C(O)=O
42	C00570	CDP-ethanolamine	HMDB0001564	123727	C00570	NCCOP(O)(=O)OP(O)(=O)OC[C@H]1O[C@H]
43	C00588	Phosphorylcholine	HMDB0001565	1014	C00588	C[N+](C)(C)CCOP(O)(O)=O
44	C00612	N1-Acetylspemidine	HMDB0001276	496	C00612	CC(=O)NCCCCNCCCN
45	C00670	Glycerophosphocholine	HMDB0000086	657272	C00670	C[N+](C)(C)CCOP([O-])(=O)OC[C@H](O)
46	C00836	Sphinganine	HMDB0000269	91486	C00836	CCCCCCCCCCCCCCC[C@@H](O)[C@@H](O)
47	C00864	Pantothenic acid	HMDB0000210	6613	C00864	CC(C)(CO)[C@@H](O)C(=O)NCCC(O)=O
48	C00946	Adenosine 2'-phosphate	HMDB0011617	53481006	C00946	NC1=NC=NC2=C1N=CN2C1O[C@H](CO)1
49	C01586	Hippuric acid	HMDB00000714	464	C01586	OC(=O)CNC(=O)C1=CC=CC=C1
50	C01780	Aldosterone	HMDB0000037	5839	C01780	[H][C@@]12CC[C@H](C(=O)CO)[C@]1(C)[C@H]
51	C02301	O-Acylcarnitine		5355	C02301	
52	C02305	Phosphocreatine	HMDB0001511	587	C02305	CN(CC(O)=O)C(=N)NP(O)(O)=O
53	C02494	1-Methyladenosine	HMDB0003331	27476	C02494	CN1C=NC2=C(N=CN2)[C@@H]2O[C@H](CO
54	C02567	N1-Acetylspermine	HMDB0001186	916	C02567	CC(=O)NCCCCNCCCN
55	C02571	L-Acetylcarнитine	HMDB0000021	7045767	C02571	CC(=O)O[C@H](CC(O)N)C[C@H](O)C[N+](C)(C)C

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

Quantitative enrichment analysis (QEA) will be performed when the user uploads a concentration table. The enrichment analysis is performed using package **globaltest**³. It uses a generalized linear model to estimate a *Q-statistic* for each metabolite set, which describes the correlation between compound concentration profiles, X, and clinical outcomes, Y. The *Q statistic* for a metabolite set is the average of the Q statistics for each metabolite in the set. **Figure 2** below summarizes the result.

³Jelle J. Goeman, Sara A. van de Geer, Floor de Kort and Hans C. van Houwelingen. *A global test for groups of genes: testing association with a clinical outcome*, Bioinformatics Vol. 20 no. 1 2004, pages 93-99

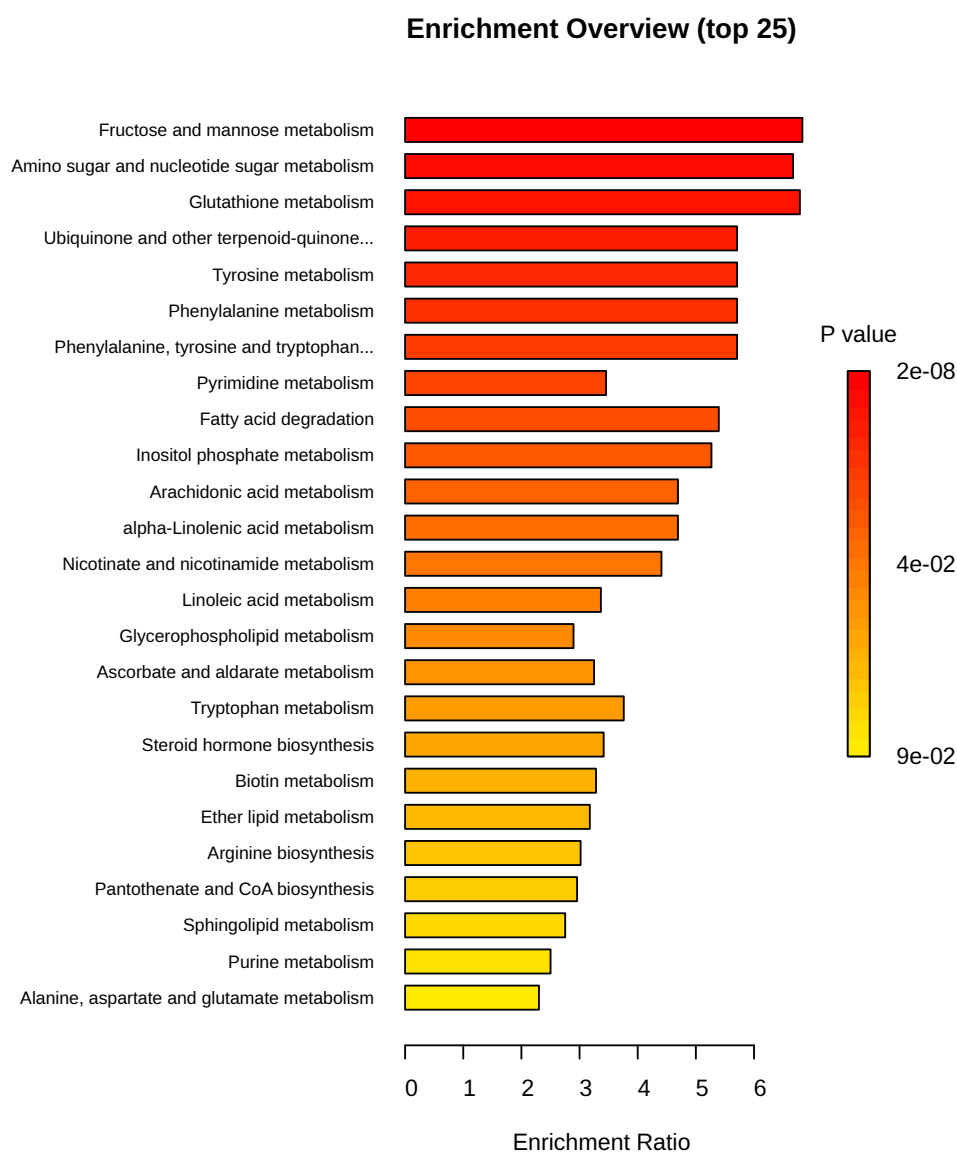


Figure 1: Summary plot for Quantitative Enrichment Analysis (QEA).

Table 2: Result from Quantitative Enrichment Analysis

	Total Cmpd	Hits	Statistic Q	Expected Q	Raw p	Holm p	FDR
Fructose and mannose metabolism	20	2	97.61	14.29	1.98E-08	7.13E-07	7.13E-07
Amino sugar and nucleotide sugar metabolism	42	5	95.32	14.29	4.40E-08	1.54E-06	7.92E-07
Glutathione metabolism	28	2	96.99	14.29	5.35E-06	1.82E-04	6.42E-05
Ubiquinone and other terpenoid-quinone biosynthesis	18	1	81.56	14.29	2.11E-03	6.97E-02	1.09E-02
Tyrosine metabolism	42	1	81.56	14.29	2.11E-03	6.97E-02	1.09E-02
Phenylalanine metabolism	8	1	81.56	14.29	2.11E-03	6.97E-02	1.09E-02
Phenylalanine, tyrosine and tryptophan biosynthesis	4	1	81.56	14.29	2.11E-03	6.97E-02	1.09E-02
Pyrimidine metabolism	39	5	49.38	14.29	3.69E-03	1.07E-01	1.66E-02
Fatty acid degradation	39	1	77.07	14.29	4.15E-03	1.16E-01	1.66E-02
Inositol phosphate metabolism	30	1	75.26	14.29	5.25E-03	1.42E-01	1.89E-02
Arachidonic acid metabolism	44	1	67.03	14.29	1.29E-02	3.36E-01	3.88E-02
alpha-Linolenic acid metabolism	13	1	67.03	14.29	1.29E-02	3.36E-01	3.88E-02
Nicotinate and nicotinamide metabolism	15	1	62.98	14.29	1.87E-02	4.49E-01	5.18E-02
Linoleic acid metabolism	5	2	48.10	14.29	2.30E-02	5.28E-01	5.90E-02
Glycerophospholipid metabolism	36	4	41.38	14.29	2.84E-02	6.25E-01	6.82E-02
Ascorbate and aldarate metabolism	9	2	46.42	14.29	3.05E-02	6.41E-01	6.86E-02
Tryptophan metabolism	41	1	53.72	14.29	3.86E-02	7.72E-01	8.17E-02
Steroid hormone biosynthesis	87	1	48.78	14.29	5.40E-02	1.00E+00	1.08E-01
Biotin metabolism	10	1	46.90	14.29	6.09E-02	1.00E+00	1.15E-01
Ether lipid metabolism	20	1	45.38	14.29	6.70E-02	1.00E+00	1.21E-01
Arginine biosynthesis	14	1	43.12	14.29	7.69E-02	1.00E+00	1.28E-01
Pantothenate and CoA biosynthesis	20	1	42.25	14.29	8.10E-02	1.00E+00	1.28E-01
Sphingolipid metabolism	32	2	39.33	14.29	8.52E-02	1.00E+00	1.28E-01
Purine metabolism	70	12	35.73	14.29	8.53E-02	1.00E+00	1.28E-01
Alanine, aspartate and glutamate metabolism	28	2	32.88	14.29	8.93E-02	1.00E+00	1.29E-01
Arginine and proline metabolism	36	3	26.48	14.29	1.44E-01	1.00E+00	1.99E-01
Lysine degradation	30	1	26.74	14.29	1.89E-01	1.00E+00	2.53E-01
Pentose and glucuronate interconversions	19	3	20.82	14.29	2.38E-01	1.00E+00	3.06E-01
Galactose metabolism	27	2	20.20	14.29	2.72E-01	1.00E+00	3.37E-01
Starch and sucrose metabolism	18	1	18.23	14.29	2.91E-01	1.00E+00	3.50E-01
Cysteine and methionine metabolism	33	3	11.79	14.29	4.66E-01	1.00E+00	5.42E-01
Riboflavin metabolism	4	1	5.51	14.29	5.76E-01	1.00E+00	6.48E-01
Citrate cycle (TCA cycle)	20	1	1.03	14.29	8.11E-01	1.00E+00	8.59E-01
Glyoxylate and dicarboxylate metabolism	31	1	1.03	14.29	8.11E-01	1.00E+00	8.59E-01
Thiamine metabolism	7	1	0.63	14.29	8.52E-01	1.00E+00	8.65E-01
Pentose phosphate pathway	23	1	0.52	14.29	8.65E-01	1.00E+00	8.65E-01

7 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"conc\", \"msetqea\", FALSE)"
[2] "mSet<-Read.TextData(mSet, \"Replacing_with_your_file_path\", \"rowu\", \"disc\");"
[3] "mSet<-SanityCheckData(mSet)"
[4] "mSet<-ReplaceMin(mSet);"
[5] "mSet<-CrossReferencing(mSet, \"kegg\");"
[6] "mSet<-CreateMappingResultTable(mSet)"
[7] "mSet<-PreparePrenormData(mSet)"
[8] "mSet<-SanityCheckData(mSet)"
[9] "mSet<-FilterVariable(mSet, \"F\", 25, \"iqr\", 0, \"mean\", 0)"
[10] "mSet<-PreparePrenormData(mSet)"
[11] "mSet<-Normalization(mSet, \"NULL\", \"NULL\", \"NULL\", ratio=FALSE, ratioNum=20)"
[12] "mSet<-PlotNormSummary(mSet, \"norm_0_\", \"png\", 72, width=NA)"
[13] "mSet<-PlotSampleNormSummary(mSet, \"snorm_0_\", \"png\", 72, width=NA)"
[14] "mSet<-SetMetabolomeFilter(mSet, F);"
[15] "mSet<-SetCurrentMsetLib(mSet, \"kegg_pathway\", 2);"
[16] "mSet<-CalculateGlobalTestScore(mSet)"
[17] "mSet<-PlotQEA.Overview(mSet, \"qea_0_\", \"net\", \"png\", 72, width=NA)"
[18] "mSet<-PlotEnrichDotPlot(mSet, \"qea\", \"qea_dot_0_\", \"png\", 72, width=NA)"
[19] "mSet<-SaveTransformedData(mSet)"
[20] "mSet<-PreparePDFReport(mSet, \"guest2982651237338019278\")\n"
```

The report was generated on Wed Jan 8 09:24:32 2025 with R version 4.3.2 (2023-10-31), OS system: Linux.