

Metabolomic Data Analysis with MetaboAnalyst 6.0

Name: guest12923759116524456756

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](COP(=O)(O)COP

56	C02567	N1- Acetylspermine	HMDB0001186	916	C02567	CC(=O)NCCCNCCCCNCCCN
57	C02571	L-Acetylcarnitine	HMDB0000201	7045767	C02571	CC(=O)O[C@@H](CC(O)=O)C[N+](C)(C)C
58	C02862	Butyrylcarnitine	HMDB0002013	213144	C02862	CCCC(=O)O[C@@H](CC([O-])=O)C[N+](C)(C)C
59	C02990	Palmitoylcarnitine	HMDB0000222	11953816	C02990	CCCCCCCCCCCCCCCC(=O)O[C@@H](CC(O)=O)C[N+](C)(C)C
60	C03017	Propionylcarnitine	HMDB0000824	188824	C03017	CCC(=O)O[C@@H](CC(O)=O)C[N+](C)(C)C
61	C03546	D-myo-Inositol 4-phosphate	HMDB0001313	440043	C03546	O[C@@H]1[C@@H](O)[C@@H](O)[C@@H](OP(=O)(O)O)C@H
62	C03794	Adenylsuccinic acid	HMDB0000536	447145	C03794	O[C@@H]1[C@@H](COP(O)(O)=O)O[C@@H](O)C@H
63	C03889	NA	NA	NA	NA	NA
64	C04100	NA	NA	NA	NA	NA
65	C04230	CE(22:5(7Z,10Z,13Z,16Z,19Z))	HMDB0010375	24779458	C04230	CC\C=C/C\C=C/C\C=C/C\C=C/C\C=C/C\C=C/C
66	C05282	gamma-Glutamylglutamic acid	HMDB0011737	92865	C05282	N[C@@H](CCC(=O)N[C@@H](CCC(O)=O)O)C(=O)O
67	C05382	D-Sedoheptulose 7-phosphate	HMDB0001068	92042786	C05382	OC[C@]1(O)O[C@@H](COP(O)(O)=O)[C@@H](O)C@H
68	C05526	S-Glutathionyl-L-cysteine	METPA0607		C05526	
69	C05551	Penicillin G	HMDB0015186	5904	C05551	[H][C@]12SC(C)(C)[C@@H](N1C(=O)[C@@H](O)C(=O)C1=CC(=O)C2COC(=O)C2=CN=C1
70	C06525	Gentianine	HMDB0303030	354616	C06525	
71	C07005	Flunisolide	HMDB0014326	82153	C07005	[H][C@@]12C[C@@]3([H])[C@]4([H])C[C@H](O)C(=O)C1=CC(=O)C2=CC(=O)C3=CC(=O)C4=O
72	C07471	Methacholine	HMDB0015654	1993	C07471	CC(C[N+](C)(C)C)OC(C)=O
73	C07968	Diethylcarbamazine	HMDB0014849	3052	C07968	CCN(CC)C(=O)N1CCN(C)CC1
74	C10438	Cinnamic acid	HMDB0000567	5372954	C10438	OC(=O)\C=C/C1=CC=CC=C1
75	C11430	9-Hydroxyphenanthrene	HMDB0059801	10229	C11430	OC1=CC2=CC=CC=C2C2=CC=CC=C12
76	C13916	Cer(d18:1/14:0)	HMDB0011773	5282310	C13916	CCCCCCCCCCCCC\C=C/[C@@H](O)[C@@H](O)C@H
77	C16207	Anthraquinone	HMDB0248468	6780	C16207	O=C1C2=CC=CC=C2C(=O)C2=CC=CC=C1
78	C17925	Methylnoradrenaline	HMDB0002832	3917	C17925	CC(N)C(O)C1=CC(O)=C(O)C=C1
79	C19463	1,5-Naphthalenediamine	HMDB0244231	16720	C19463	NC1=CC=CC2=C1C=CC=C2N
80	C20387	Biotin sulfone	HMDB0004818	21252323	C20387	[H][C@]12CS(=O)(=O)[C@@H](CCCCC(=O)O)C(=O)C1=CC(=O)C2=CC(=O)C3=CC(=O)C4=CC(=O)C3=CC4
81	C21484	LysoPE(18:0/0:0)	HMDB0011130	9547068	C21484	[H][C@@](O)(COC(=O)CCCCCCCCCCCCC)C@H
82	C22599	NA	NA	NA	NA	NA

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

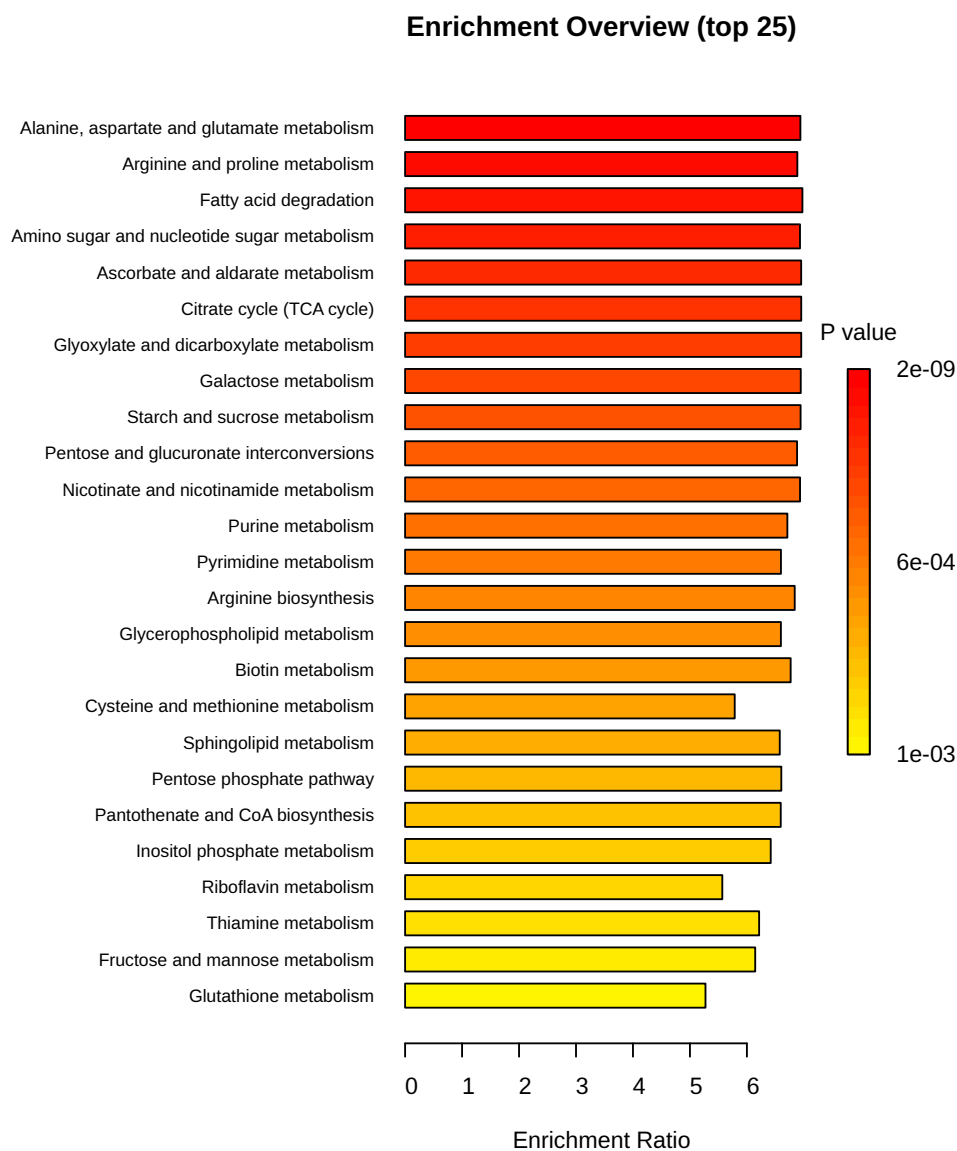


Figure 1: Summary plot for Quantitative Enrichment Analysis (QEA).

Table 2: Result from Quantitative Enrichment Analysis

	Total Cmpd	Hits	Statistic Q	Expected Q	Raw p	Holm p	FDR
Alanine, aspartate and glutamate metabolism	28	2	99.13	14.29	1.90E-09	6.47E-08	6.47E-08
Arginine and proline metabolism	36	3	98.37	14.29	4.80E-09	1.58E-07	8.16E-08
Fatty acid degradation	39	1	99.64	14.29	1.47E-08	4.70E-07	1.67E-07
Amino sugar and nucleotide sugar metabolism	42	4	99.03	14.29	3.72E-08	1.15E-06	3.16E-07
Ascorbate and aldarate metabolism	9	2	99.33	14.29	6.78E-08	2.03E-06	4.36E-07
Citrate cycle (TCA cycle)	20	1	99.34	14.29	8.98E-08	2.60E-06	4.36E-07
Glyoxylate and dicarboxylate metabolism	31	1	99.34	14.29	8.98E-08	2.60E-06	4.36E-07
Galactose metabolism	27	2	99.22	14.29	1.15E-07	3.11E-06	4.90E-07
Starch and sucrose metabolism	18	1	99.18	14.29	1.72E-07	4.47E-06	6.49E-07
Pentose and glucuronate interconversions	19	3	98.30	14.29	2.41E-07	6.03E-06	8.20E-07
Nicotinate and nicotinamide metabolism	15	1	99.04	14.29	2.76E-07	6.63E-06	8.54E-07
Purine metabolism	70	11	95.86	14.29	6.27E-07	1.44E-05	1.78E-06
Pyrimidine metabolism	39	5	94.25	14.29	3.25E-06	7.15E-05	8.50E-06
Arginine biosynthesis	14	1	97.71	14.29	3.80E-06	7.98E-05	9.22E-06
Glycerophospholipid metabolism	36	5	94.24	14.29	8.56E-06	1.71E-04	1.94E-05
Biotin metabolism	10	1	96.69	14.29	1.15E-05	2.18E-04	2.44E-05
Cysteine and methionine metabolism	33	3	82.66	14.29	1.94E-05	3.50E-04	3.88E-05
Sphingolipid metabolism	32	2	93.96	14.29	5.06E-05	8.61E-04	9.57E-05
Pentose phosphate pathway	23	1	94.34	14.29	5.77E-05	9.24E-04	1.03E-04
Pantothenate and CoA biosynthesis	20	1	94.20	14.29	6.23E-05	9.34E-04	1.06E-04
Inositol phosphate metabolism	30	1	91.69	14.29	1.85E-04	2.59E-03	3.00E-04
Riboflavin metabolism	4	2	79.55	14.29	2.69E-04	3.50E-03	4.16E-04
Thiamine metabolism	7	1	88.78	14.29	4.61E-04	5.53E-03	6.81E-04
Fructose and mannose metabolism	20	1	87.80	14.29	5.96E-04	6.55E-03	8.44E-04
Glutathione metabolism	28	2	75.33	14.29	1.13E-03	1.13E-02	1.53E-03
Ubiquinone and other terpenoid-quinone biosynthesis	18	1	78.32	14.29	3.48E-03	3.14E-02	4.00E-03
Tyrosine metabolism	42	1	78.32	14.29	3.48E-03	3.14E-02	4.00E-03
Phenylalanine metabolism	8	1	78.32	14.29	3.48E-03	3.14E-02	4.00E-03
Phenylalanine, tyrosine and tryptophan biosynthesis	4	1	78.32	14.29	3.48E-03	3.14E-02	4.00E-03
Arachidonic acid metabolism	44	1	77.77	14.29	3.76E-03	3.14E-02	4.00E-03
Linoleic acid metabolism	5	1	77.77	14.29	3.76E-03	3.14E-02	4.00E-03
alpha-Linolenic acid metabolism	13	1	77.77	14.29	3.76E-03	3.14E-02	4.00E-03
Lysine degradation	30	1	59.30	14.29	2.54E-02	5.08E-02	2.62E-02
Ether lipid metabolism	20	1	14.24	14.29	3.57E-01	3.57E-01	3.57E-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, \"guest12923759116524456756\")\n"
```

The report was generated on Wed Jan 8 09:39:11 2025 with R version 4.3.2 (2023-10-31), OS system: Linux.