

Metabolomic Data Analysis with MetaboAnalyst 5.0

Name: guest12749090168457640060

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.^{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 Ma

	Query	Match	HMDB	PubChem	KEGG	SMILES
1	HMDB0000714	Hippuric acid	HMDB0000714	464	C01586	<chem>C1=CC=C(C=C1)C(=O)NCC(=O)O</chem>
2	HMDB0000784	Azelaic acid	HMDB0000784	2266	C08261	<chem>C(CCCC(=O)O)CCCC(=O)O</chem>
3	HMDB0000893	Suberic acid	HMDB0000893	10457	C08278	<chem>C(CCCC(=O)O)CCC(=O)O</chem>
4	HMDB0061384	NA	NA	NA	NA	NA
5	HMDB0000072	cis-Aconitic acid	HMDB0000072	643757	C00417	<chem>C(/C(=C/C(=O)O)/C(=O)O)C(=O)O</chem>
6	HMDB0000132	Guanine	HMDB0000132	764	C00242	<chem>C1=NC2=C(N1)C(=O)N=C(N2)N</chem>
7	HMDB0006116	3-Hydroxyhippuric acid	HMDB0006116	450268		<chem>C1=CC(=CC(=C1)O)C(=O)NC(=O)O</chem>
8	HMDB0006275	Dopamine 3-O-sulfate	HMDB0006275	122136	C13690	<chem>C1=CC(=C(C=C1CCN)OS(=O)(=O)O)C</chem>
9	HMDB0002643	3-(3-Hydroxyphenyl)-3-hydroxypropanoic acid	HMDB0002643	102959		<chem>C1=CC(=CC(=C1)O)C(C(=O)O)C(=O)O</chem>
10	HMDB0000193	Isocitric acid	HMDB0000193	1198	C00311	<chem>C(C(C(C(=O)O)O)C(=O)O)C(=O)O</chem>
11	HMDB0000446	N-Alpha-acetyllysine	HMDB0000446	192590	C12989	<chem>CC(=O)NC(CCCCN)C(=O)O</chem>
12	HMDB0011103	1,7-Dimethyluric acid	HMDB0011103	91611	C16356	<chem>CN1C2=C(C(NC1=O)NC(=O)N(C2)=O)C</chem>
13	HMDB0000881	Xanthurenic acid	HMDB0000881	5699	C02470	<chem>C1=CC2=C(C(=C1)O)NC(=CC2=O)C(=O)O</chem>
14	HMDB0000730	Isobutyrylglycine	HMDB0000730	10855600		<chem>CC(C)C(=O)NCC(=O)O</chem>
15	HMDB0003099	1-Methyluric acid	HMDB0003099	69726	C16359	<chem>CN1C(=O)C2=C(NC(=O)N2)N(C(=O)O)C</chem>
16	HMDB0000875	Trigonelline	HMDB0000875	5570	C01004	<chem>C[N+](=O)C=C(C(=C1)C(=O)O)C1=O</chem>
17	HMDB0000152	Gentisic acid	HMDB0000152	3469	C00628	<chem>C1=CC(=C(C=C1O)C(=O)O)C(=O)O</chem>
18	HMDB0000157	Hypoxanthine	HMDB0000157	790	C00262	<chem>C1=NC2=C(C(N1C(=O)N=CN2)=O)C</chem>
19	HMDB0002730	Nicotinamide N-oxide	HMDB0002730	72661		<chem>C1=CC(=C(N1[+](=O)[O-])C(=O)O)C</chem>
20	HMDB0000678	Isovalerylglycine	HMDB0000678	546304		<chem>CC(C)CC(=O)NCC(=O)O</chem>
21	HMDB0000206	N6-Acetyl-L-lysine	HMDB0000206	92832	C02727	<chem>CC(=O)NCCCC[C@@H](C(=O)O)C</chem>
22	HMDB0010327	Dehydroisoandrosterone 3-glucuronide	HMDB0010327	53480448	C03033	<chem>CC12CCC3C(C1CCC2=O)CC=C3OC(=O)C</chem>
23	HMDB0003072	Quinic acid	HMDB0003072	6508	C00296	<chem>OC1C[C@@](O)(C[C@@H](O)[C@H](O)C</chem>
24	HMDB0000094	Citric acid	HMDB0000094	311	C00158	<chem>C(C(=O)O)C(CC(=O)O)C(=O)O</chem>
25	HMDB0014611	Quinine	HMDB0014611	3034034	C06526	<chem>COC1=CC2=C(C(=CN=C2C=C1)C</chem>
26	HMDB0000440	3-Hydroxyphenylacetic acid	HMDB0000440	12122	C05593	<chem>C1=CC(=CC(=C1)O)CC(=O)O</chem>
27	HMDB0000391	7-Ketodeoxycholic acid	HMDB0000391	188292	C04643	<chem>C[C@H](CCC(=O)O)[C@H]1CC[C@H](O)C</chem>
28	HMDB0000625	Gluconic acid	HMDB0000625	10690	C00257	<chem>C([C@H]([C@H]([C@H]([C@H](O)CO)O)O)O)O</chem>
29	HMDB0003464	4-Guanidinobutanoic acid	HMDB0003464	500	C01035	<chem>C(CC(=O)O)CN=C(N)N</chem>
30	HMDB0000917	Urocholic acid	HMDB0000917	122340	C17644	<chem>C[C@H](CCC(=O)O)[C@H]1CC[C@H](O)C</chem>
31	HMDB0000669	Ortho-Hydroxyphenylacetic acid	HMDB0000669	11970	C05852	<chem>C1=CC=C(C(=C1)CC(=O)O)C(=O)O</chem>
32	HMDB0000191	L-Aspartic acid	HMDB0000191	5960	C00049	<chem>C([C@@H](C(=O)O)N)C(=O)O</chem>
33	HMDB0000355	3-Hydroxymethylglutaric acid	HMDB0000355	1662	C03761	<chem>CC(C(=O)O)C(=O)O</chem>
34	HMDB0001406	Niacinamide	HMDB0001406	936	C00153	<chem>C1=CC(=CN=C1)C(=O)N</chem>
35	HMDB0000418	18-Hydroxycortisol	HMDB0000418	44263343		<chem>C[C@]12CCC(=O)C=C1CC[C@]2(C)O</chem>
36	HMDB0029992	Tetrahydropentoxylene	HMDB0029992	53481442		<chem>C1C(NC(C2=C1C3=CC=C(C=C3)O)O)O</chem>
37	HMDB0001987	2-Hydroxy-2-methylbutyric acid	HMDB0001987	95433		<chem>CCC(C)(C(=O)O)O</chem>
38	HMDB0000306	Tyramine	HMDB0000306	5610	C00483	<chem>C1=CC(=CC=C1CCN)O</chem>
39	HMDB0000842	Quinaldic acid	HMDB0000842	7124	C06325	<chem>C1=CC=C2C(=C1)C=C(CC(=N2)O)C</chem>
40	HMDB0000259	Serotonin	HMDB0000259	5202	C00780	<chem>C1=CC2=C(C(=C1O)C(=CN2)O)C</chem>
41	HMDB0013678	4-Hydroxyhippuric acid	HMDB0013678	151012		<chem>C1=CC(=CC=C1C(=O)NCC(=O)O)C</chem>
42	HMDB0001713	m-Coumaric acid	HMDB0001713	637541	C12621	<chem>C1=CC(=CC(=C1)O)/C=C/C(C1=O)C</chem>
43	HMDB0000133	Guanosine	HMDB0000133	6802	C00387	<chem>C1=NC2=C(N1[C@H]3[C@@H](O)C</chem>
44	HMDB0060001	NA	NA	NA	NA	NA
45	HMDB0013713	N-acetyltryptophan	HMDB0013713	700653		<chem>[H][C@@](CC1=CN2C=CC=CC=C2N1)C</chem>
46	HMDB0000262	Thymine	HMDB0000262	1135	C00178	<chem>CC1=CN(C(=O)NC1=O)C</chem>
47	HMDB0001297	Norcotinine	HMDB0001297	413		<chem>C1CC(=O)NC1C2=CN=CC=C2</chem>
48	HMDB0001991	7-Methylxanthine	HMDB0001991	68374	C16353	<chem>CN1C=NC2=C1C(=O)NC(=O)N2C</chem>
49	HMDB0002024	Imidazoleacetic acid	HMDB0002024	96215	C02835	<chem>C1=C(NC=N1)CC(=O)O</chem>
50	HMDB0000912	Succinyladenosine	HMDB0000912	20849086		<chem>C1=NC2=C(C(=N1)N[C@@H](O)C</chem>
51	HMDB0000138	Glycocholic acid	HMDB0000138	23617285	C01921	<chem>C[C@H](CCC(=O)O)NCC(=O)O</chem>
52	HMDB0028933	Leucyl-Leucine	HMDB0028933	76807	C11332	<chem>CC(C)CC(N)C(=O)NC(CC(C)C)C(=O)O</chem>
53	HMDB0001982	3,7-Dimethyluric acid	HMDB0001982	83126	C16360	<chem>CN1C2=C(NC1=O)N(C(=O)O)C</chem>
54	HMDB0000162	L-Proline	HMDB0000162	145742	C00148	<chem>C1C[C@H](NC1)C(=O)O</chem>
55	HMDB0000512	N-Acetyl-L-phenylalanine	HMDB0000512	74839	C03519	<chem>CC(=O)N[C@@H](CC1=CC=C(C=C1)C</chem>
56	HMDB0004827	Proline betaine	HMDB0004827	7016563	C10172	<chem>C[N+](C)(CCC[C@H]1C(=O)[O-])C</chem>
57	HMDB0000097	Choline	HMDB0000097	305	C00114	<chem>C[N+](C)(C)CCO</chem>
58	HMDB0001860	Paraxanthine	HMDB0001860	4687	C13747	<chem>CN1C=NC2=C1C(=O)N(C(=O)O)C</chem>

59	HMDB0002894	5-Methylcytosine	HMDB0002894	65040	C02376	<chem>CC1=C(NC(=O)N=C1)N</chem>
60	HMDB0000661	Glutaric acid	HMDB0000661	743	C00489	<chem>C(CC(=O)O)CC(=O)O</chem>
61	HMDB0000866	N-Acetyl-L-tyrosine	HMDB0000866	68310	C01657	<chem>CC(=O)N[C@@H](CC1=CC=CC=C1)C(=O)O</chem>
62	HMDB0000020	p-Hydroxyphenylacetic acid	HMDB0000020	127	C00642	<chem>C1=CC(=CC=C1CC(=O)O)O</chem>
63	HMDB0003331	1-Methyladenosine	HMDB0003331	27476	C02494	<chem>CN1C=NC2=C(C1=N)N=CN2</chem>
64	HMDB0005923	N4-Acetylcytidine	HMDB0005923	107461		<chem>CC(=O)NC1=NC(=O)N(C=C1)N</chem>
65	HMDB0003157	Guanidinosuccinic acid	HMDB0003157	439918	C03139	<chem>C([C@@H](C(=O)O)N=C(N)N)C(=O)O</chem>
66	HMDB0062179	NA	NA	NA	NA	NA
67	HMDB0000172	L-Isoleucine	HMDB0000172	6306	C00407	<chem>CC[C@H](C)[C@@H](C(=O)O)N</chem>
68	HMDB0000050	Adenosine	HMDB0000050	60961	C00212	<chem>C1=NC2=C(C(=N1)N)N=CN2</chem>
69	HMDB0000754	3-Hydroxyisovaleric acid	HMDB0000754	69362	C20827	<chem>CC(C)(CC(=O)O)O</chem>
70	HMDB0000230	N-Acetylneuraminic acid	HMDB0000230	445063	C19910	<chem>CC(=O)N[C@@H]1[C@H](C[C@@H](O)C(=O)N)C[C@@H](O)C1=O</chem>
71	HMDB0000715	Kynurenic acid	HMDB0000715	3845	C01717	<chem>C1=CC=C2C(=C1)C(=O)C=C2</chem>
72	HMDB0000118	Homovanillic acid	HMDB0000118	1738	C05582	<chem>COc1c(C(=CC(=C1)CC(=O)O)O)cc(C)c1</chem>
73	HMDB0004824	N2,N2-Dimethylguanosine	HMDB0004824	92919		<chem>CN(C)C1=NC(=O)C2=C(N1)N=CN2</chem>
74	HMDB0241300	NA	NA	NA	NA	NA
75	HMDB0002035	4-Hydroxycinnamic acid	HMDB0002035	637542	C00811	<chem>C1=CC(=CC=C1/C=C/C(=O)O)O</chem>
76	HMDB0000729	Alpha-Hydroxyisobutyric acid	HMDB0000729	11671		<chem>CC(C)(C(=O)O)O</chem>
77	HMDB0062640	3-hydroxy-2-isobutyrate	HMDB0062640	87	C01188	<chem>CC(CO)C(O)=O</chem>
78	HMDB0001434	3-Methoxytyrosine	HMDB0001434	1670		<chem>COc1c(C(=CC(=C1)CC(C(=O)O)O)O)cc(C)c1</chem>
79	HMDB0000201	L-Acetylcarnitine	HMDB0000201	7045767	C02571	<chem>CC(=O)OC(CC(=O)[O-])C[N+](=O)[O-]</chem>
80	HMDB0000824	Propionylcarnitine	HMDB0000824	107738	C03017	<chem>CCC(=O)OC(CC(=O)[O-])C[N+](=O)[O-]</chem>
81	HMDB0004620	N-a-Acetyl-L-arginine	HMDB0004620	67427		<chem>CC(=O)N[C@@H](CCCNC(=O)N)C(=O)O</chem>
82	HMDB0000630	Cytosine	HMDB0000630	597	C00380	<chem>C1=C(NC(=O)N=C1)N</chem>
83	HMDB0000092	Dimethylglycine	HMDB0000092	673	C01026	<chem>CN(C)CC(=O)O</chem>
84	HMDB0000732	Hydroxykynurenine	HMDB0000732	89	C02794	<chem>C1=CC(=C(C(=C1)O)N)C(=O)O</chem>
85	HMDB0002432	Sumiki's acid	HMDB0002432	80642	C20448	<chem>C1=C(OC(=C1)C(=O)O)CO</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

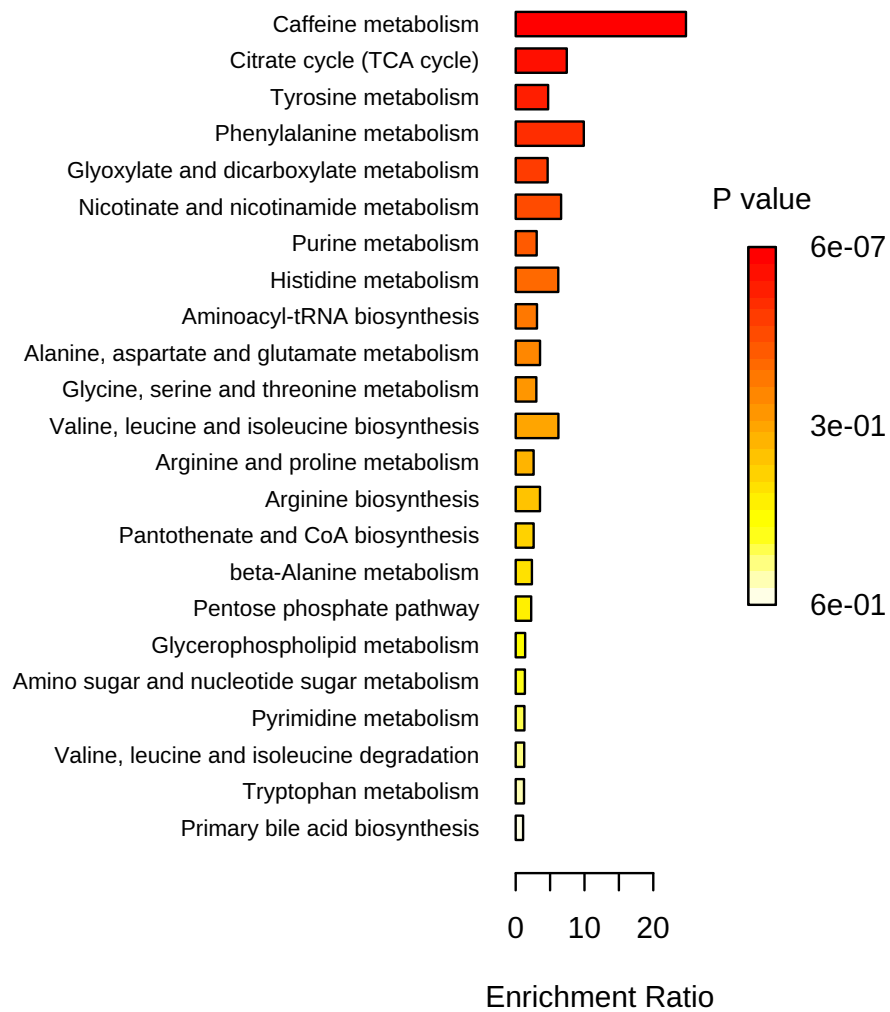


Figure 1: Summary Plot for Over Representation Analysis (ORA)

Table 2: Result from Over Representation Analysis

	total	expected	hits	Raw p	Holm p	FDR
Caffeine metabolism	10	0.20	5	5.63E-07	4.73E-05	4.73E-05
Citrate cycle (TCA cycle)	20	0.40	3	6.73E-03	5.59E-01	2.49E-01
Tyrosine metabolism	42	0.85	4	8.90E-03	7.29E-01	2.49E-01
Phenylalanine metabolism	10	0.20	2	1.60E-02	1.00E+00	3.37E-01
Glyoxylate and dicarboxylate metabolism	32	0.65	3	2.49E-02	1.00E+00	4.17E-01
Nicotinate and nicotinamide metabolism	15	0.30	2	3.52E-02	1.00E+00	4.17E-01
Purine metabolism	65	1.31	4	3.89E-02	1.00E+00	4.17E-01
Histidine metabolism	16	0.32	2	3.97E-02	1.00E+00	4.17E-01
Aminoacyl-tRNA biosynthesis	48	0.97	3	6.98E-02	1.00E+00	6.51E-01
Alanine, aspartate and glutamate metabolism	28	0.56	2	1.08E-01	1.00E+00	9.05E-01
Glycine, serine and threonine metabolism	33	0.67	2	1.42E-01	1.00E+00	1.00E+00
Valine, leucine and isoleucine biosynthesis	8	0.16	1	1.51E-01	1.00E+00	1.00E+00
Arginine and proline metabolism	38	0.77	2	1.77E-01	1.00E+00	1.00E+00
Arginine biosynthesis	14	0.28	1	2.49E-01	1.00E+00	1.00E+00
Pantothenate and CoA biosynthesis	19	0.38	1	3.23E-01	1.00E+00	1.00E+00
beta-Alanine metabolism	21	0.42	1	3.50E-01	1.00E+00	1.00E+00
Pentose phosphate pathway	22	0.44	1	3.63E-01	1.00E+00	1.00E+00
Glycerophospholipid metabolism	36	0.73	1	5.24E-01	1.00E+00	1.00E+00
Amino sugar and nucleotide sugar metabolism	37	0.75	1	5.34E-01	1.00E+00	1.00E+00
Pyrimidine metabolism	39	0.79	1	5.53E-01	1.00E+00	1.00E+00
Valine, leucine and isoleucine degradation	40	0.81	1	5.62E-01	1.00E+00	1.00E+00
Tryptophan metabolism	41	0.83	1	5.71E-01	1.00E+00	1.00E+00
Primary bile acid biosynthesis	46	0.93	1	6.14E-01	1.00E+00	1.00E+00

7 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"conc\", \"msetora\", FALSE)"
[2] "compd.vec<-c(\"HMDB0000714\", \"HMDB0000784\", \"HMDB0000893\", \"HMDB0061384\", \"HMDB0000072\", \"")
[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, \"guest12749090168457640060\")\\n"
```

The report was generated on Thu Nov 24 07:19:33 2022 with R version 4.2.2 (2022-10-31), OS system: Linux, version: -Ubuntu SMP Thu Oct 13 08:03:55 UTC 2022 .