

MetaComp User's Guide

Peng Zhai

November 14, 2017

Contents

1	Introduction	2
2	Contact information	2
3	Citing MetaComp	2
4	Prerequisites and installation	2
	4.1 Prerequisite	2
	4.2 Source code	2
	4.3 Installation	2
5	Input data	3
	5.1 Abundance profile matrix data	3
	5.2 Obtain profile from BLAST	4
	5.3 Obtain profile from HMMER	5
	5.4 Obtain profile from Kraken	6
	5.5 Obtain profile from MG-RAST	7
	5.6 Obtain profile from PhymmBL	8
	5.7 Obtain profile from MZmine	9
6	Multivariate statistics	10
	6.1 Cluster analysis	10
	6.2 Principal component analysis	11
7	Hypothesis testing	12
	7.1 Two samples test	12
	7.2 Multiple samples test	13
	7.3 Two groups of samples test	14
8	Environmental factors analysis	15

1 Introduction

MetaComp is a graphical software for analyzing meta-omic (i.e. metagenomics, metatranscriptomics, metaproteomics and metabolomics) profiles with related environmental information, such as phylogenetic profiles indicating the number of marker genes assigned to different taxonomic units or functional profiles indicating the number of sequences assigned to different subsystems or pathways. The aim of this document provide an easy but comprehensive introduction to MetaComp and show how it can be used to analyze meta-omic data. MetaComp is applicable to any meta-omics data by accepting abundance profile matrices (APM) saved as txt or BIOM format files [1]. Moreover, MetaComp can autmatically converts the output of several widely used platform into MetaComp-compatible input file.

2 Contact information

MetaComp is in active development. We encourage you to send any suggestions, comments and bug reports to hqzhu@pku.edu.cn. If reporting a bug, please provide as much information as possible and the related data which causes the bug. This will allow us to quickly resolve the issue.

3 Citing MetaComp

4 Prerequisites and installation

4.1 Prerequisite

- Windows 7 or higher version.
- Install Microsoft Office Excel 2010 or higher version.
- Install the required R packages using the following commands in the R console:

```
install.packages("pheatmap")
```

4.2 Source code

<https://github.com/pzhaipku/MetaComp>

4.3 Installation

- Windows: Download file “” setup from our website: <http://cqb.pku.edu.cn/ZhuLab/MetaComp/download.html>.

- Linux: Download file “” setup from our website: <http://cqb.pku.edu.cn/ZhuLab/MetaComp/download.html>. File Annotation.RData is the annotation information for Linux. Please put it in the same folder with file MetaComp.R. Finally, please input the following commands:

```
source("./MetaComp.R")
```

5 Input data

5.1 Abundance profile matrix data

MetaComp reads input file in text format, and the values in the file should be separated by tab. The first row of the file shows the name of samples, while the first column represents the selected statistical feature. The cell of the table indicates the hit number of one sample to the given feature. Users must select **Abundance profile matrix (.txt or .biom)** radio button from **Profile** dialog box in **Load Data** option within **File** menu before choose the input profile. Moreover, the .biom format input must be convert to biom.table format before loading.(Figure 1-3)

This format of input is the only format that can be load in Linux version. The command line is as follow:

```
input_data = readFeature(file pathway, featureType = "pfam" or "cog", format = "txt" or "biom")
```

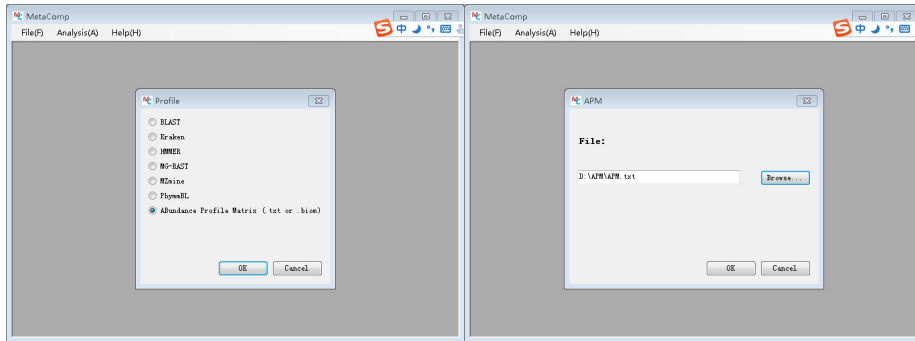


Figure 1:

Figure 2:

Feature	ID1	ID2	ID3
O00001	2	11	4
O00002	2	1	1
O00003	1	0	2
O00004	7	5	5
O00005	0	2	0
O00006	5	11	3
O00007	2	2	2
O00008	7	8	5
O00009	1	2	1
O00010	1	2	2
O00012	4	1	4
O00013	6	5	4
O00014	4	3	1
O00015	6	7	1
O00016	2	3	1
O00017	3	4	1
O00018	5	1	2
O00019	5	6	1
O00020	2	2	0
O00021	6	4	1

Figure 3:

5.2 Obtain profile from BLAST

MetaComp also accepts meta-omics profiles obtained from BLAST ([2], <https://blast.ncbi.nlm.nih.gov/Blast.cgi>) result. MetaComp works directly with BLAST result obtained by clicking on download in result web page, followed by selecting Hit Table(text) output type choice. Moreover, the BLAST result file can be obtained from table format (-outfmt 7). MetaComp can convert these BLAST results to standard Abundance profile matrices (APM) data through selecting **BLAST** radio button from **Profile** dialog box in **Load Data** option within **File** menu. After opening up the **BLAST** dialog box, you can select the BLAST result files you wish to input. (Figure 4-6)

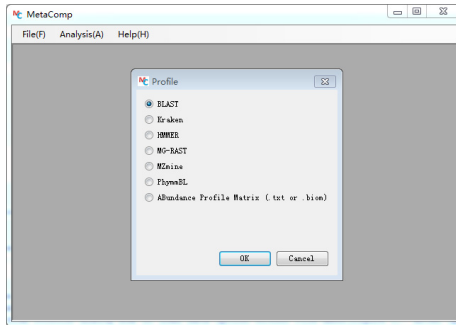


Figure 4:

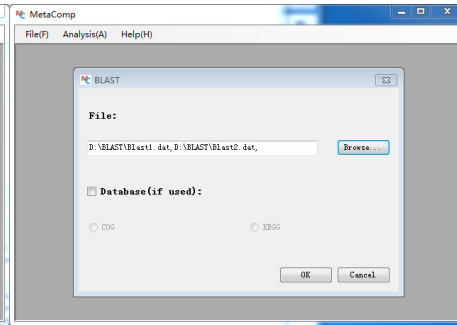


Figure 5:

Feature	File1	File2
act_C01_06180	1	1
act_C1aaa_2359	1	1
alg_B01_17360	1	1
bth_B7_2146	1	1
can_B07_08150	1	1
cke_C00_0537	1	1
leg_B07_11010	2	2
lva_B0001_2004	2	2
och_C1aaa_3950	1	1
rhm_B07_21400	1	1
rhc_B00W_27195	1	1
rxx_B01_08930	2	2
rhm_B01_21300	1	1
rhm_B01_21410	1	1
fpn_P78_14660	2	2
lcy_B07_18420	1	1
rwn_C01_25950	2	2
shg_B07_01170	1	1
shg_B07_01160	1	1
shh_A1_08130	3	3
osp_0409_1637	1	1

Figure 6:

5.3 Obtain profile from HMMER

The input profile can also be acquired from HMMER ([3], <http://hmmer.org/>). After downloading **hmmer-3.1b2.tar.gz** from <http://hmmer.org/> and unpacking it, you can get the desired results from **hmmsearch** command. MetaComp can convert these file into MetaComp-compatible profiles through selecting **HMMER** radio button convert these BLAST results to standard Abundance profile matrices (APM) data through selecting **BLAST** radio button from **Profile** dialog box in **Load Data** option within **File** menu. Click on the **OK** button after selecting the result file you wish to convert. (Figure 7-9)

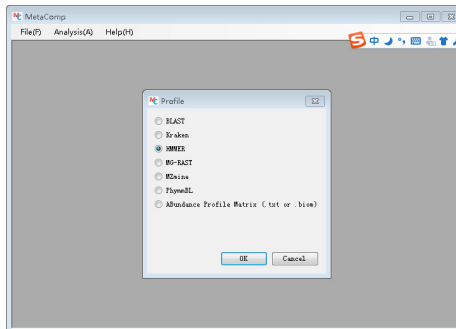


Figure 7:

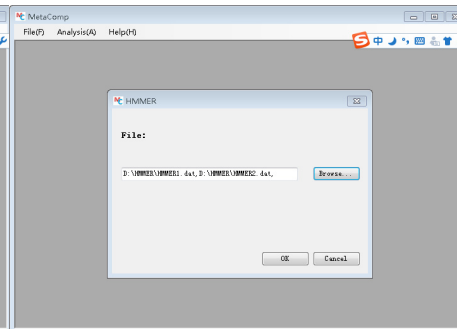


Figure 8:

Feature	F1341	F1342
PF0979.12	4613	0
PF0979.13	4667	0
PF0511.10	0	10663
PF0538.9	0	10659
PF0463.11	0	1209
PF0452.11	0	1582
PF0158.24	0	4884
PF1432.4	0	4613

Figure 9:

5.4 Obtain profile from Kraken

Kraken ([4], <http://ccb.jhu.edu/software/kraken/>) result files are achieved from kraken-translate command. The selection of Kraken result file can initiate after choosing **Kraken** radio button from **Profile** dialog box in **Load Data** option within **File** menu. (Figure 10-12)

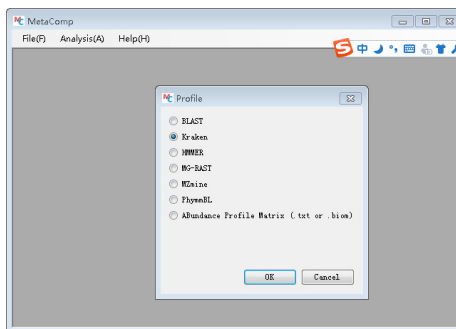


Figure 10:

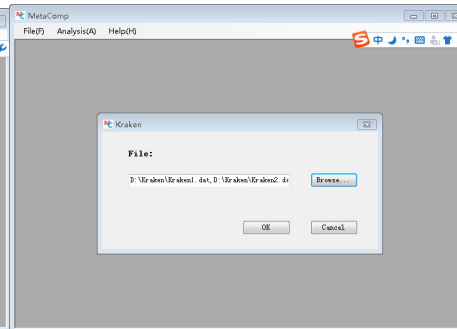


Figure 11:

Feature	File1	File2
root:cellular organism...	15	21
root:cellular organism...	277	620
root:cellular organism...	2813	5900
root:cellular organism...	1554	3563
root:cellular organism...	892	1978
root:cellular organism...	659	1598
root:cellular organism...	821	1830
root:cellular organism...	1307	4303
root:cellular organism...	1384	3245
root:cellular organism...	222	461
root:cellular organism...	1495	3387
root:cellular organism...	5825	13698
root:cellular organism...	340	754
root:cellular organism...	95	232
root:cellular organism...	1575	3460
root:cellular organism...	172	427
root:cellular organism...	1416	3200
root:cellular organism...	245	613
root:cellular organism...	82	172
root:cellular organism...	642	1563
root:cellular organism...	792	1775

Figure 12:

5.5 Obtain profile from MG-RAST

MetaComp provides support for analyzing MG-RAST taxonomic or functional profiles. Visit the MG-RAST website ([5], <http://metagenomics.anl.gov/>) and browse the list of public metagenomes. Profiles for multiple samples can be obtained and downloaded as tab-separated values (tsv) file using the table data visualization. To work with MG-RAST profiles, they must be converted into a MetaComp-compatible profile. From within MetaComp, select the **MG-RAST** radio button from **Profile** dialog box in **Load Data** option within **File** menu. This opens up the **MG-RAST** dialog box. Click on the **OK** button after selecting the MG-RAST profile you wish to convert. (Figure 13-15)

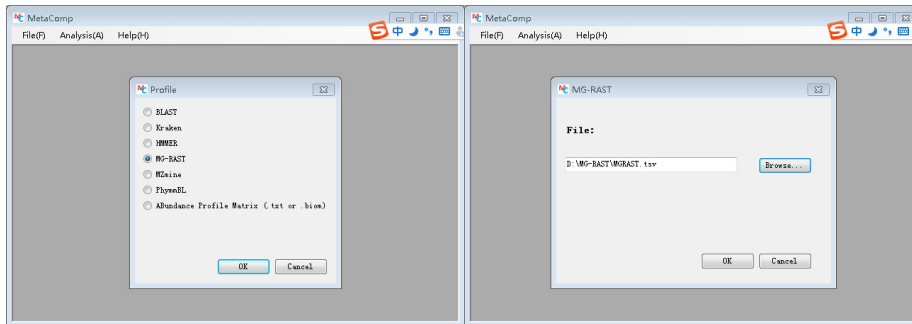


Figure 13:

Figure 14:

	042PHM_CAGATC_1009_B1_001	042PHM_CAGATC_1009_B2_001	042PHM_ACTTGA_ID1
Animals	0	2	4
Archaea	0	0	0
Bacteria	4895212	4697247	3729119
Eukaryota	63955	63195	544753
Viruses	9693	8292	7824
Other sequences	0	1	1

Figure 15:

5.6 Obtain profile from PhymmBL

PhymmBL([6], <http://www.cbcb.umd.edu/software/phymm/>) result files are achieved from `scoreReads.pl` command. The selection of PhymmBL result file can initiate after choosing **PhymmBL** radio button from **Profile** dialog box in **Load Data** option within **File** menu. (Figure 16-18)

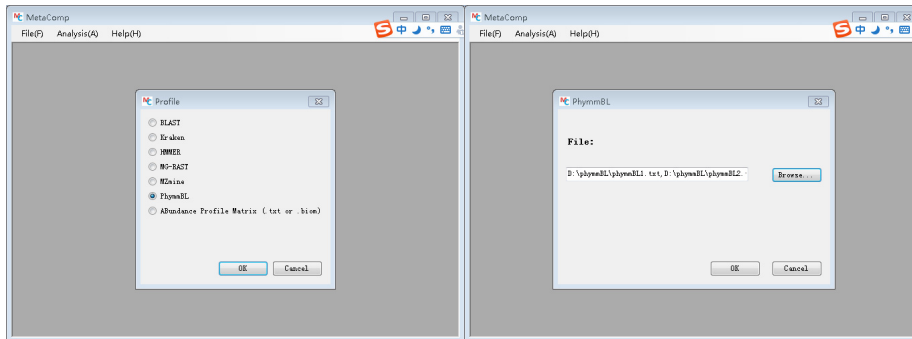


Figure 16:

Figure 17:

Feature	File1	File2
Terriglobus_grossus_BCM	51	185
Bacteroides_fragilis_B	9618	9227
Streptococcus_gammalis	251	328
Listeria_monocytogenes	167	148
Bacillus_cereus_B4264	278	325
Deinifactoribius_bacul	419	1019
Bacillus_thuringiensis	2113	2481
Alteromonas_aeolicus	110	91
Brachyspira_glandricol	53	49
Enterobacter_cloacae_s	151	149
Deinifactoribius_erie	859	1112
Cronobacter_subaeolis	581	678
Thermaplasma_acidophil	453	460
Histococcus_salsus	356	384
Deinifactoribius_alle	467	595
Candidatus_Mitrospiru	152	249
Tropomyces_sitotetraci	811	1039
Legionella_sp._J77918	361	596
Yersinia_pestis_0102038	916	1129
Escherichia_coli_B121L	136	250
Streptococcus_suis_021	319	305

Figure 18:

5.7 Obtain profile from MZmine

MZmine([7], <http://mzmine.github.io/>) result files are achieved as Figure 19. The selection of MZmine result file can initiate after choosing **MZmine** radio button from **Profile** dialog box in **Load Data** option within **File** menu. (Figure 20-22)

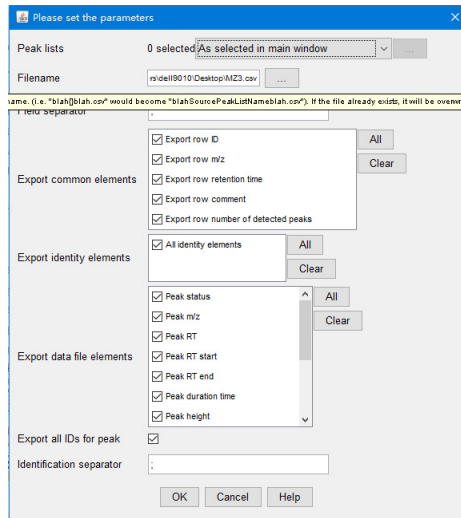


Figure 19:

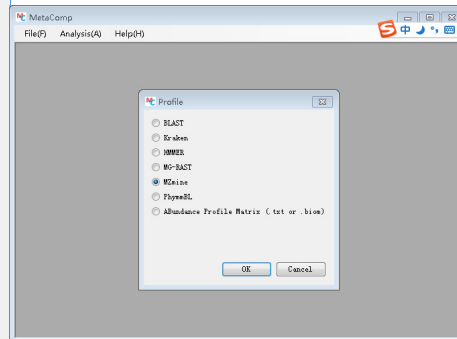


Figure 20:

All these input format example mentioned above can be download from:

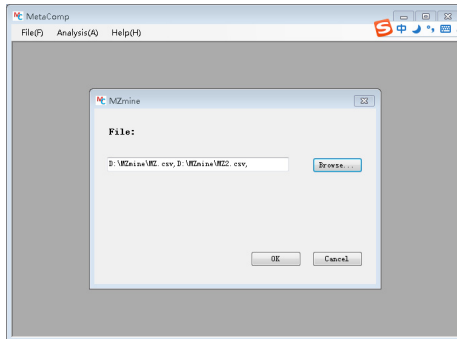


Figure 21:

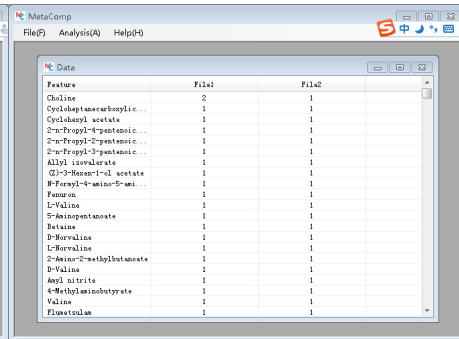


Figure 22:

6 Multivariate statistics

6.1 Cluster analysis

Cluster analysis can be performed in two models: K-means clustering and hierarchical clustering. K-means clustering model requires users to input the cluster number. Cluster analysis can be operated through the **Clustering analysis** dialog in the **Analysis** menu. (Figure 23-26)

Linux commands line:

Hcluster(input_data) (for hierarchical clustering)

KMeans(input_data, cluster number) (for k-means clustering)

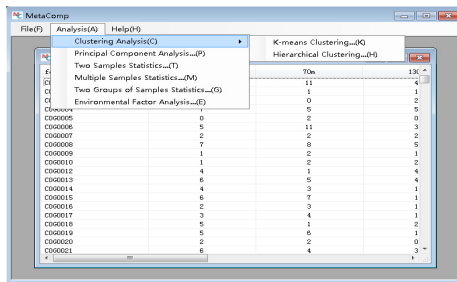


Figure 23:

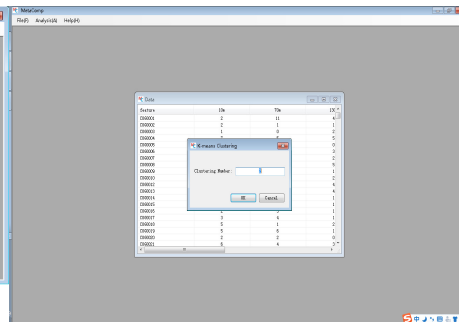


Figure 24:

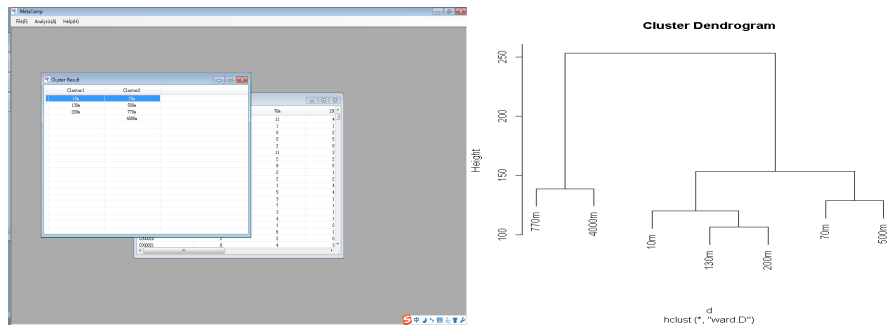


Figure 25: Result of k-means cluster. Figure 26: Result of hierarchical clustering.

6.2 Principal component analysis

Principal component analysis (PCA) is applied in two model: whole data analysis model and clustering analysis model. Whole data analysis model is the model we common used. Also, clustering analysis model can apply PCA within the clustering information. The example of clustering information can be download from PCA can be applied through the **Principal component analysis** dialog in **Analysis** menu. (Figure 27-29)

Linux commands line:

PCA(input_data, ShowsampleName="text" or "NA")

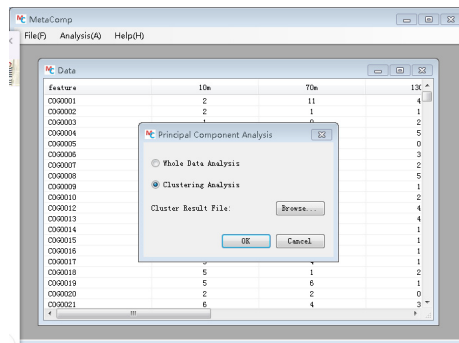


Figure 27:

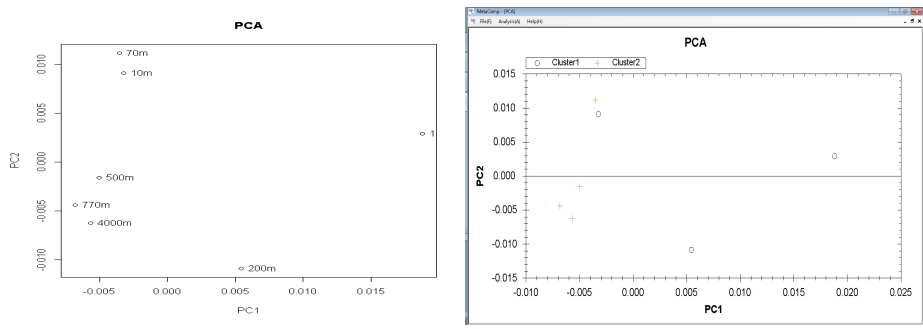


Figure 28: Result of whole data analysis
 Figure 29: Result of clustering analysis model.

7 Hypothesis testing

7.1 Two samples test

To analyze a pair of samples, click on the **Two samples Statistic** dialog in **Analysis** menu. In this dialog, you can choose a favorable statistical test, p-value and data type. Moreover, you can choose the database you require if the feature in your profile is Pfam or COG database.(Figure 30-32)

```
Linux commands line:
result=twoSamplesComp(input_data)
plotTopVar(result)
```

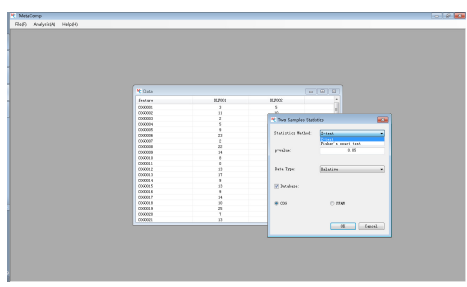


Figure 30:

Rank	Feature ID	MF002	p-value	bonferroni	Annotation
1	C00002	8	2.019637	1.0	0.99988 [D] aspartate/asparagine/serine/threonine tyrosylase, arginase family
3	C00080	4	0.12208	1.0	0.99988 [D] endonuclease restriction factor
4	C00009	23	15.017123	1.0	0.99988 [D] ribitol/ribose dehydrogenase
5	C00010	6	50.048256	1.0	0.99988 [D] isochrysochrome (tetraether patch paperfatty)
6	C00019	10	10.010220	1.0	0.99988 [D] trans-methyl chaperonin/immunoglobulin chaperone
7	C00011	20	14.019541	1.0	0.99988 [D] phosphoglycerate dehydrogenase and related dehydrogenase
8	C00079	7	3.011722	1.0	0.99988 [D] BMP family
9	C00060	13	27.011734	1.0	0.99988 [D] isochrysochrome
10	C00022	1	5.016058	1.0	0.99988 [D] 2-methylsuccinate (4S)-epimerase
11	C00004	14	27.01779	1.0	0.99988 [D] cytoplasmic phosphoribosyltransferase and related transferase
12	C00021	12	9.017212	1.0	0.99988 [D] transferase
13	C00071	15	9.017212	1.0	0.99988 [D] ribonucleic chaperone (small heat shock protein)
14	C00014	10	22.018094	1.0	0.99988 [D] adenylsuccinate synthase
15	C00016	8	17.019959	1.0	0.99988 [D] aspartate/serine/threonine dehydrogenase
16	C00006	8	5.022873	1.0	0.99988 [D] riboflavin synthase beta-chain
17	C00012	6	5.022873	1.0	0.99988 [D] penicillin synthase
18	C00041	7	15.02210	1.0	0.99988 [D] predicted sugar kinase
19	C00086	12	24.02210	1.0	0.99988 [D] DNA-directed RNA polymerase, beta' subunit/19S D0 subunit
20	C00018	10	25.022604	1.0	0.99988 [D] RNA transcription factor/epidemiological field (only MF cycle/indirect)
21	C00004	5	11.027254	1.0	0.99988 [D] amonin peroxidase
22	C00043	4	2.027926	1.0	0.99988 [D] poly(3-hydroxybutyrate) deacetylase and related deacetylase
23	C00044	7	14.02925	1.0	0.99988 [D] DNA and RNA restriction endonuclease
24	C00011	9	17.029801	1.0	0.99988 [D] penicillin synthase
25	C00024	8	15.030052	1.0	0.99988 [D] histidinol-5-phosphate synthase
26	C00027	11	9.030362	1.0	0.99988 [D] propanoate dehydrogenase
27	C00008	10	8.030826	1.0	0.99988 [D] ribonucleic protein D0
28	C00004	12	23.032548	1.0	0.99988 [D] penicillin acetyltransferase
29	C00027	16	26.034688	1.0	0.99988 [D] predicted ATPase of the P0-loop superfamily implicated in cell cycle

Figure 31: Analysis result(excel).

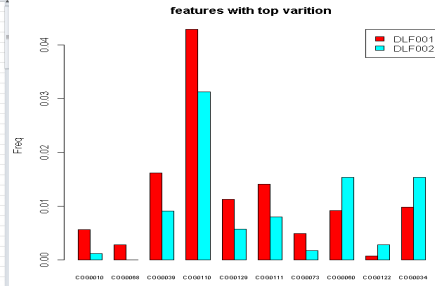


Figure 32: Analysis result(figure).

7.2 Multiple samples test

To analyze multiple samples, click on the **Multiple samples Statistic** dialog in **Analysis** menu. In this dialog, you can choose a favorable statistical test, p-value and data type. Just like Two samples test, you can choose the database you require if the feature in your profile is Pfam or COG database. Also you can select the most favorable visualizations you demanded.(Figure 33-35)

Linux commands line:

```
result=twoSamplesComp(input_data)
```

```
plotTopVar(result)
```

```
plotClust(result)
```

```
plotMDS(result, ShowsampleName = "legend", "text", "both" or "NA")
```

```
plotHeatMap(input_data, show_rownames = T or F, cluster_rows = T or F)
```

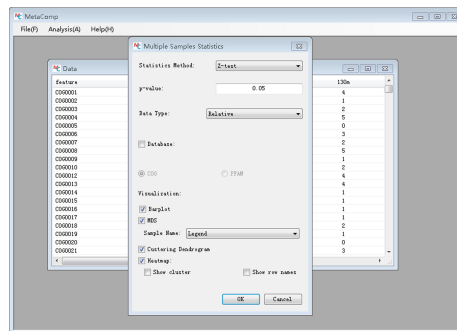


Figure 33:

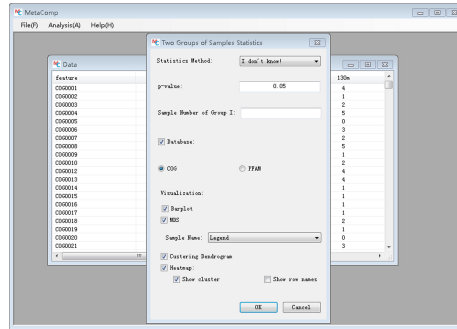


Figure 36:

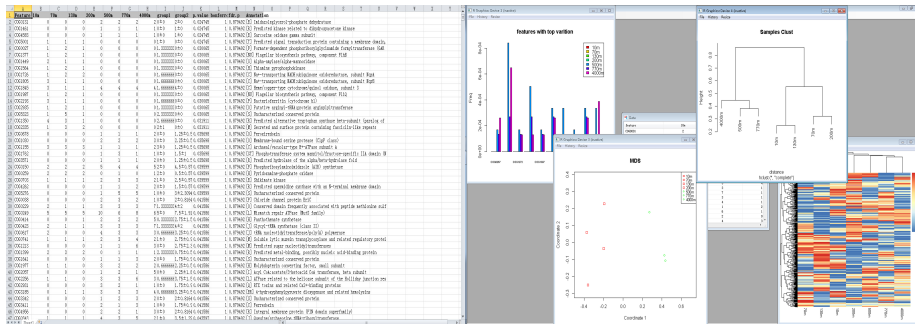


Figure 37: Analysis result(excel).

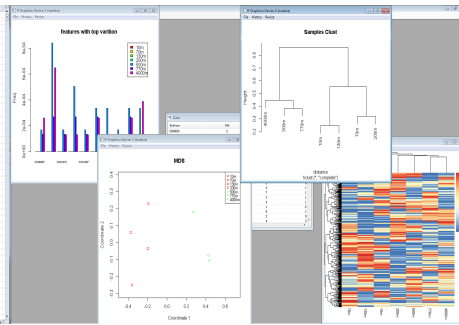


Figure 38: Analysis result(figure).

8 Environmental factors analysis

To operate environmental factors analysis, click on the **Environmental factors analysis** dialog in **Analysis** menu. In this dialog, you need to load the environmental factors information, input the p-value, choose whether you require to include the cross term of environmental factors and load Pfam or COG database while analysing. The example of environmental factors information can be download from.(Figure 39-40)

Linux commands line:

EnvironmentFactor(input_data,environment factor file pathway,Feature number)

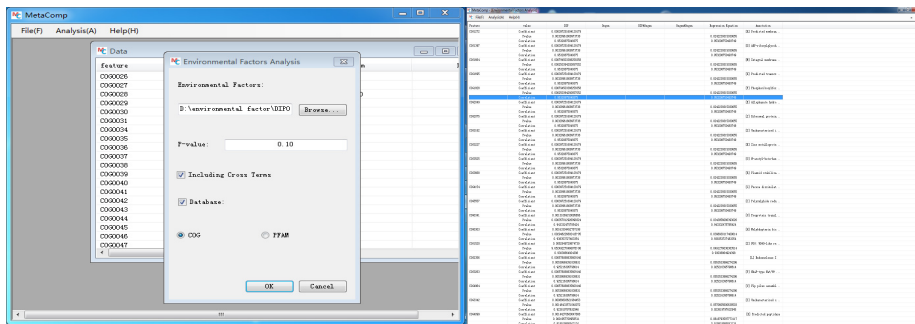


Figure 39: Analysis result(excel).

Figure 40: Analysis result(figure).

Bibliography

- [1] Daniel McDonald, Jose C Clemente, Justin Kuczynski, Jai Ram Rideout, Jesse Stombaugh, Doug Wendel, Andreas Wilke, Susan Huse, John Hufnagle, Folker Meyer, et al. The biological observation matrix (biom) format or: how i learned to stop worrying and love the ome-ome. *GigaScience*, 1(1):1, 2012.
- [2] Stephen F Altschul, Warren Gish, Webb Miller, Eugene W Myers, and David J Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, 1990.
- [3] Jaina Mistry, Robert D Finn, Sean R Eddy, Alex Bateman, and Marco Punta. Challenges in homology search: Hmmer3 and convergent evolution of coiled-coil regions. *Nucleic acids research*, 41(12):e121–e121, 2013.
- [4] Derrick E Wood and Steven L Salzberg. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome biology*, 15(3):1, 2014.
- [5] Elizabeth M Glass, Jared Wilkening, Andreas Wilke, Dionysios Antonopoulos, and Folker Meyer. Using the metagenomics rast server (mg-rast) for analyzing shotgun metagenomes. *Cold Spring Harbor Protocols*, 2010(1):pdb-prot5368, 2010.
- [6] Arthur Brady and Steven L Salzberg. Phymm and phymmbl: metagenomic phylogenetic classification with interpolated markov models. *Nature methods*, 6(9):673–676, 2009.
- [7] Tomáš Pluskal, Sandra Castillo, Alejandro Villar-Briones, and Matej Orešič. Mzmine 2: Modular framework for processing, visualizing, and analyzing mass spectrometry-based molecular profile data. *BMC Bioinformatics*, 11(1):395, 2010.