

Metabolomic Data Analysis with MetaboAnalyst 5.0

Name: guest15282839817347858290

November 29, 2023

1 Data Processing and Normalization

1.1 Reading and Processing the Raw Data

MetaboAnalyst accepts a variety of data types generated in metabolomic studies, including compound concentration data, binned NMR/MS spectra data, NMR/MS peak list data, as well as MS spectra (NetCDF, mzXML, mzDATA). Users need to specify the data types when uploading their data in order for MetaboAnalyst to select the correct algorithm to process them. Table 1 summarizes the result of the data processing steps.

1.1.1 Reading Peak Intensity Table

The peak intensity table should be uploaded in comma separated values (.csv) format. Samples can be in rows or columns, with class labels immediately following the sample IDs.

Samples are in rows and features in columns The uploaded file is in comma separated values (.csv) format. The uploaded data file contains 54 (samples) by 637 (peaks(mz/rt)) data matrix.

1.1.2 Data Integrity Check

Before data analysis, a data integrity check is performed to make sure that all the necessary information has been collected. The class labels must be present and contain only two classes. If samples are paired, the class label must be from $-n/2$ to -1 for one group, and 1 to $n/2$ for the other group (n is the sample number and must be an even number). Class labels with same absolute value are assumed to be pairs. Compound concentration or peak intensity values should all be non-negative numbers. By default, all missing values, zeros and negative values will be replaced by the half of the minimum positive value found within the data (see next section)

1.1.3 Missing value imputations

Too many zeroes or missing values will cause difficulties for downstream analysis. MetaboAnalyst offers several different methods for this purpose. The default method replaces all the missing and zero values with a small values (the half of the minimum positive values in the original data) assuming to be the detection limit. The assumption of this approach is that most missing values are caused by low abundance metabolites (i.e. below the detection limit). In addition, since zero values may cause problem for data normalization (i.e. log), they are also replaced with this small value. User can also specify other methods, such as replace by mean/median, or use K-Nearest Neighbours (KNN), Probabilistic PCA (PPCA), Bayesian PCA (BPCA) method, Singular Value Decomposition (SVD) method to impute the missing values ¹. Please choose the one that is the most appropriate for your data.

¹Stacklies W, Redestig H, Scholz M, Walther D, Selbig J. *pcaMethods: a bioconductor package, providing PCA methods for incomplete data.*, Bioinformatics 2007 23(9):1164-1167

Zero or missing values were replaced by $1/5$ of the min positive value for each variable.

1.1.4 Data Filtering

The purpose of the data filtering is to identify and remove variables that are unlikely to be of use when modeling the data. No phenotype information are used in the filtering process, so the result can be used with any downstream analysis. This step can usually improves the results. Data filter is strongly recommended for datasets with large number of variables (> 250) datasets contain much noise (i.e.chemometrics data). Filtering can usually improve your results².

*For data with number of variables < 250 , this step will reduce 5% of variables; For variable number between 250 and 500, 10% of variables will be removed; For variable number btween 500 and 1000, 25% of variables will be removed; And 40% of variabed will be removed for data with over 1000 variables. The None option is only for less than 5000 features. Over that, if you choose None, the IQR filter will still be applied. In addition, the maximum allowed number of variables is **10000***

No data filtering is performed.

²Hackstadt AJ, Hess AM. *Filtering for increased power for microarray data analysis*, BMC Bioinformatics. 2009; 10: 11.

Table 1: Summary of data processing results

	Features (positive)	Missing/Zero	Features (processed)
C9.1	637	0	637
C9.2	637	0	637
C9.3	637	0	637
C9.4	637	0	637
D9.1	637	0	637
D9.2	637	0	637
D9.3	637	0	637
D9.4	637	0	637
F9.1	637	0	637
F9.2	637	0	637
F9.3	637	0	637
F9.4	637	0	637
X9.1	637	0	637
X9.2	637	0	637
X9.3	637	0	637
X9.4	637	0	637
C12.1	637	0	637
C12.3	637	0	637
C12.4	637	0	637
D12.1	637	0	637
D12.2	637	0	637
D12.3	637	0	637
D12.4	637	0	637
F12.1	637	0	637
F12.2	637	0	637
F12.3	637	0	637
F12.4	637	0	637
X12.1	637	0	637
X12.2	637	0	637
X12.3	637	0	637
BLANK_2_Dup1	635	2	637
BLANK_2_Dup2	635	2	637
BLANK_2_Dup3	635	2	637
BLANK_2	635	2	637
QC.1.Dup	637	0	637
QC.1	637	0	637
QC.2	637	0	637
QC.3	637	0	637
C15.1	637	0	637
C15.2	637	0	637
C15.3	637	0	637
C15.4	637	0	637
D15.1	636	1	637
D15.2	637	0	637
D15.3	636	1	637
D15.4	636	1	637
F15.1	636	1	637
F15.2	636	1	637
F15.3	637	0	637
F15.4	637	0	637
X15.1	637	0	637
X15.2	637	0	637
X15.3	637	0	637
X15.4	637	0	637

1.2 Data Normalization

The data is stored as a table with one sample per row and one variable (bin/peak/metabolite) per column. The normalization procedures implemented below are grouped into four categories. Sample specific normalization allows users to manually adjust concentrations based on biological inputs (i.e. volume, mass); row-wise normalization allows general-purpose adjustment for differences among samples; data transformation and scaling are two different approaches to make features more comparable. You can use one or combine both to achieve better results.

The normalization consists of the following options:

1. Row-wise procedures:
 - Sample specific normalization (i.e. normalize by dry weight, volume)
 - Normalization by the sum
 - Normalization by the sample median
 - Normalization by a reference sample (probabilistic quotient normalization)³
 - Normalization by a pooled or average sample from a particular group
 - Normalization by a reference feature (i.e. creatinine, internal control)
 - Quantile normalization
2. Data transformation :
 - Log transformation (base 10)
 - Square root transformation
 - Cube root transformation
3. Data scaling:
 - Mean centering (mean-centered only)
 - Auto scaling (mean-centered and divided by standard deviation of each variable)
 - Pareto scaling (mean-centered and divided by the square root of standard deviation of each variable)
 - Range scaling (mean-centered and divided by the value range of each variable)

Figure 1 shows the effects before and after normalization.

³Dieterle F, Ross A, Schlotterbeck G, Senn H. *Probabilistic quotient normalization as robust method to account for dilution of complex biological mixtures. Application in 1H NMR metabonomics*, 2006, Anal Chem 78 (13);4281 - 4290

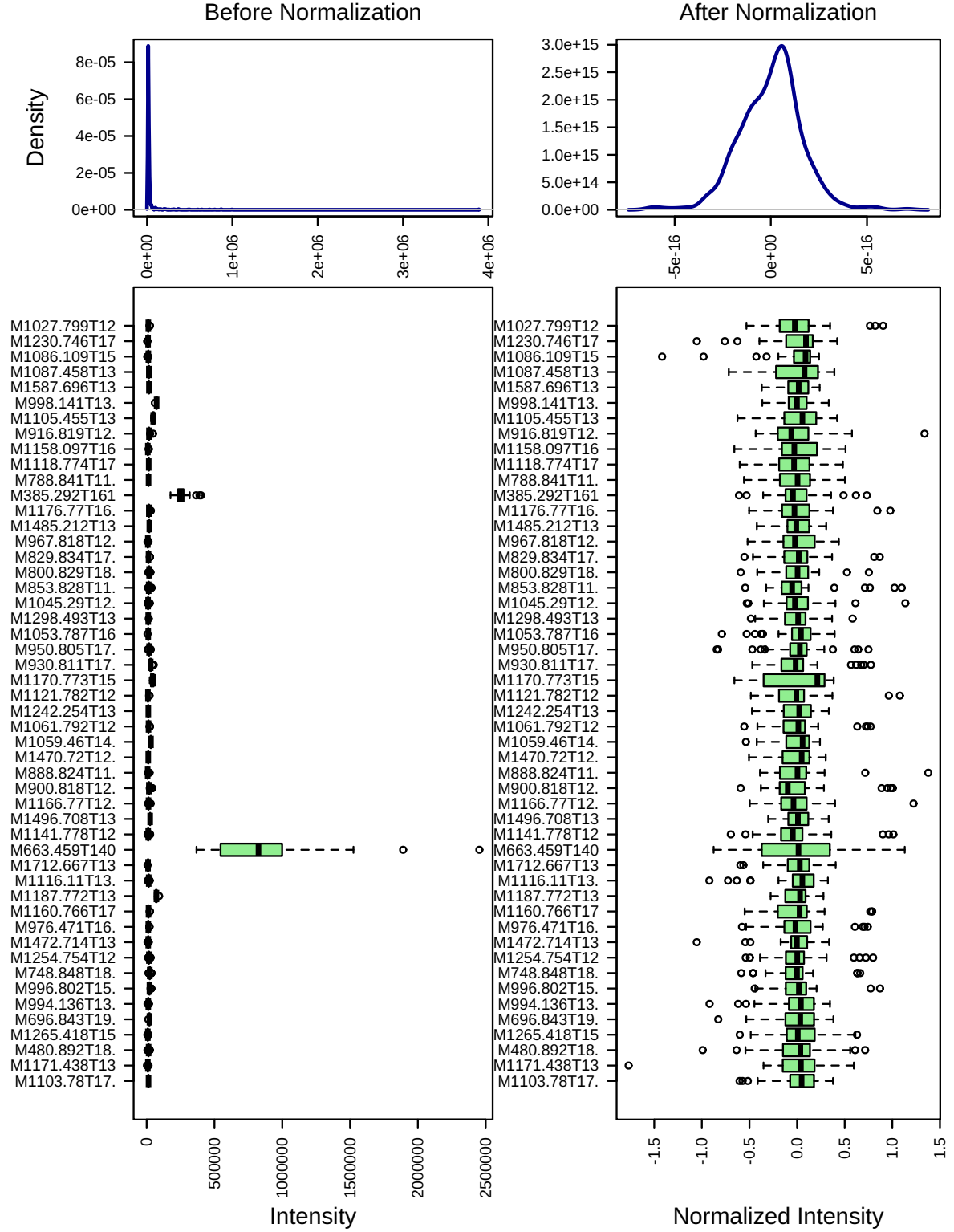


Figure 1: Box plots and kernel density plots before and after normalization. The boxplots show at most 50 features due to space limit. The density plots are based on all samples. Selected methods : Row-wise normalization: Normalization by a reference feature; Data transformation: Log10 Normalization; Data scaling: Pareto Scaling.

2 Statistical and Machine Learning Data Analysis

MetaboAnalyst offers a variety of methods commonly used in metabolomic data analyses. They include:

1. Univariate analysis methods:
 - Fold Change Analysis
 - T-tests
 - Volcano Plot
 - One-way ANOVA and post-hoc analysis
 - Correlation analysis
2. Multivariate analysis methods:
 - Principal Component Analysis (PCA)
 - Partial Least Squares - Discriminant Analysis (PLS-DA)
3. Robust Feature Selection Methods in microarray studies
 - Significance Analysis of Microarray (SAM)
 - Empirical Bayesian Analysis of Microarray (EBAM)
4. Clustering Analysis
 - Hierarchical Clustering
 - Dendrogram
 - Heatmap
 - Partitional Clustering
 - K-means Clustering
 - Self-Organizing Map (SOM)
5. Supervised Classification and Feature Selection methods
 - Random Forest
 - Support Vector Machine (SVM)

Please note: some advanced methods are available only for two-group sample analysis.

2.1 One-way ANOVA

Univariate analysis methods are the most common methods used for exploratory data analysis. For multi-group analysis, MetaboAnalyst provides one-way Analysis of Variance (ANOVA). As ANOVA only tells whether the overall comparison is significant or not, it is usually followed by post-hoc analyses in order to identify which two levels are different. MetaboAnalyst provides two most commonly used methods for this purpose - Fisher's least significant difference method (Fisher's LSD) and Tukey's Honestly Significant Difference (Tukey's HSD). The univariate analyses provide a preliminary overview about features that are potentially significant in discriminating the conditions under study.

Figure 2 shows the important features identified by ANOVA analysis. Table 2 shows the details of these features. The **post-hoc Sig. Comparison** column shows the comparisons between different levels that are significant given the p value threshold.

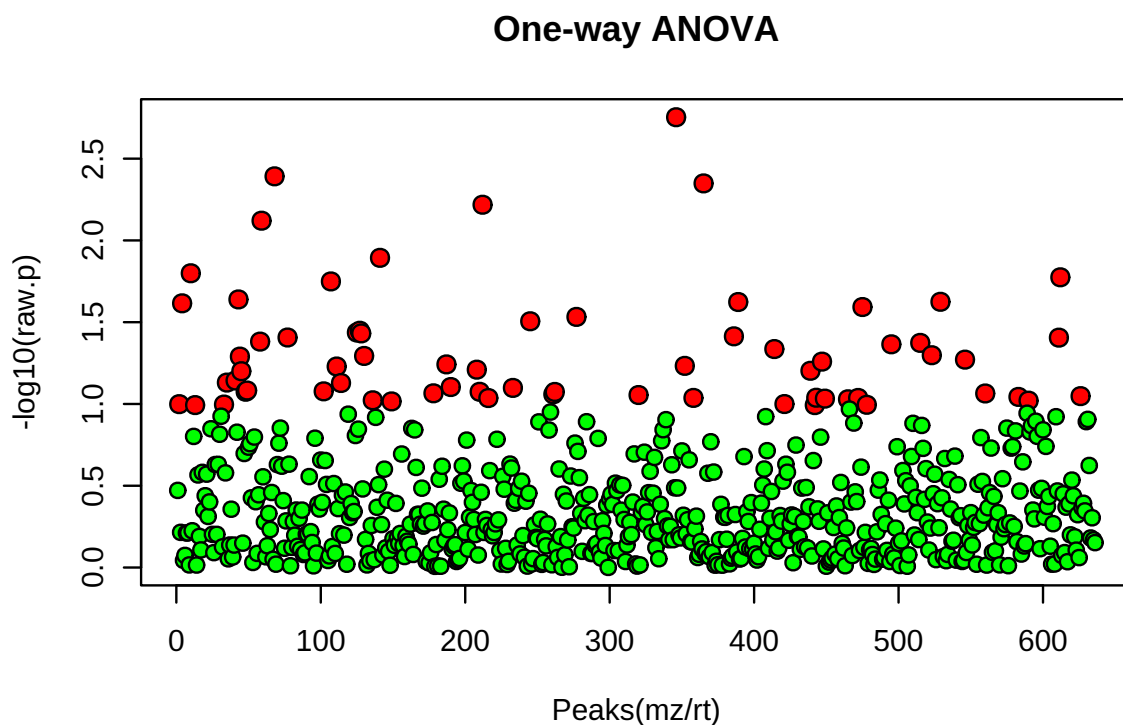


Figure 2: Important features selected by ANOVA plot with p value threshold 0.95057.

Table 2: Top 50 features identified by One-way ANOVA and post-hoc analysis

	Peaks(mz/rt)	f.value	p.value	-log10(p)	FDR	Tukey's HSD
1	M1085.78T15.949	7.3733	0.0017650	2.7533	0.94990	Second-First; Third-First
2	M681.468T1391.964	6.2772	0.0040561	2.3919	0.94990	Second-First
3	M1107.445T13.643	6.1489	0.0044806	2.3487	0.94990	Second-First; Third-Second
4	M937.806T17.14	5.7654	0.0060502	2.2182	0.95056	Third-First; Third-Second
5	M663.461T1374.938	5.4820	0.0075741	2.1207	0.95056	Second-First
6	M839.337T903.707	4.8335	0.0127790	1.8935	0.95056	Third-First; Third-Second
7	M309.131T17.119	4.5689	0.0158780	1.7992	0.95056	Second-First
8	M1576.695T13.11	4.5001	0.0168060	1.7745	0.95056	Third-Second
9	M774.828T14.53	4.4298	0.0178130	1.7493	0.95056	Third-First
10	M633.474T1424.866	4.1260	0.0229500	1.6392	0.95056	Third-Second
11	M1355.738T13.11	4.0863	0.0237280	1.6247	0.95056	Second-First
12	M1130.78T14.113	4.0826	0.0238020	1.6234	0.95056	Third-First
13	M247.083T835.562	4.0608	0.0242420	1.6154	0.95056	Third-First
14	M1246.09T13.111	3.9988	0.0255410	1.5928	0.95056	Second-First
15	M1021.326T836.385	3.8346	0.0293480	1.5324	0.95056	Third-First
16	M977.8T13.648	3.7622	0.0312110	1.5057	0.95056	Third-Second
17	M817.26T833.71	3.6061	0.0356590	1.4478	0.95056	Third-First
18	M817.177T13.11	3.5768	0.0365670	1.4369	0.95056	Third-First
19	M817.264T834.906	3.5621	0.0370300	1.4315	0.95056	Third-First
20	M1126.276T16.125	3.5143	0.0385800	1.4136	0.95056	Third-First
21	M711.854T15.533	3.4950	0.0392270	1.4064	0.95056	Third-First
22	M1576.691T13.111	3.4927	0.0393040	1.4056	0.95056	Third-Second
23	M663.459T1408.312	3.4279	0.0415610	1.3813	0.95056	Second-First
24	M1337.74T13.11	3.4070	0.0423180	1.3735	0.95056	Second-First
25	M1282.251T13.11	3.3839	0.0431680	1.3648	0.95056	Third-First
26	M1163.774T15.006	3.3067	0.0461530	1.3358	0.95056	Third-Second
27	M1350.741T12.108	3.2068	0.0503360	1.2981	0.95056	Third-First
28	M824.345T903.272	3.1950	0.0508540	1.2937	0.95056	
29	M633.482T1643.428	3.1845	0.0513210	1.2897	0.95056	
30	M1397.732T14.114	3.1354	0.0535650	1.2711	0.95056	
31	M1200.766T12.171	3.1034	0.0550860	1.2590	0.95056	
32	M902.806T18.126	3.0586	0.0572870	1.2419	0.95056	
33	M1092.287T12.108	3.0350	0.0584810	1.2330	0.95056	
34	M794.18T13.113	3.0250	0.0589950	1.2292	0.95056	
35	M930.811T17.116	2.9741	0.0616920	1.2098	0.95056	
36	M1188.441T13.111	2.9562	0.0626720	1.2029	0.95056	
37	M634.485T1424.862	2.9494	0.0630450	1.2004	0.95056	
38	M631.468T1459.998	2.8008	0.0718770	1.1434	0.95056	
39	M591.475T1442.976	2.7680	0.0739940	1.1308	0.95056	
40	M799.252T819.392	2.7612	0.0744380	1.1282	0.95056	
41	M908.151T15.062	2.6973	0.0787840	1.1036	0.95056	
42	M967.818T12.107	2.6833	0.0797680	1.0982	0.95056	
43	M648.465T1441.075	2.6424	0.0827310	1.0823	0.95056	
44	M771.505T1357.807	2.6291	0.0837110	1.0772	0.95056	
45	M647.463T1427.898	2.6225	0.0842040	1.0747	0.95056	
46	M931.815T16.119	2.6218	0.0842580	1.0744	0.95056	
47	M996.798T17.124	2.6199	0.0844050	1.0736	0.95056	
48	M894.49T15.115	2.5981	0.0860560	1.0652	0.95056	
49	M1436.722T13.111	2.5955	0.0862620	1.0642	0.95056	
50	M996.139T12.108	2.5829	0.0872320	1.0593	0.95056	

2.2 Principal Component Analysis (PCA)

PCA is an unsupervised method aiming to find the directions that best explain the variance in a data set (X) without referring to class labels (Y). The data are summarized into much fewer variables called *scores* which are weighted average of the original variables. The weighting profiles are called *loadings*. The PCA analysis is performed using the `prcomp` package. The calculation is based on singular value decomposition.

The Rscript `chemometrics.R` is required. Figure 3 is pairwise score plots providing an overview of the various separation patterns among the most significant PCs; Figure 4 is the scree plot showing the variances explained by the selected PCs; Figure 5 shows the 2-D scores plot between selected PCs; Figure 6 shows the biplot between the selected PCs. Interactive 3-D scores plots are not included here and can be directly downloaded from website.

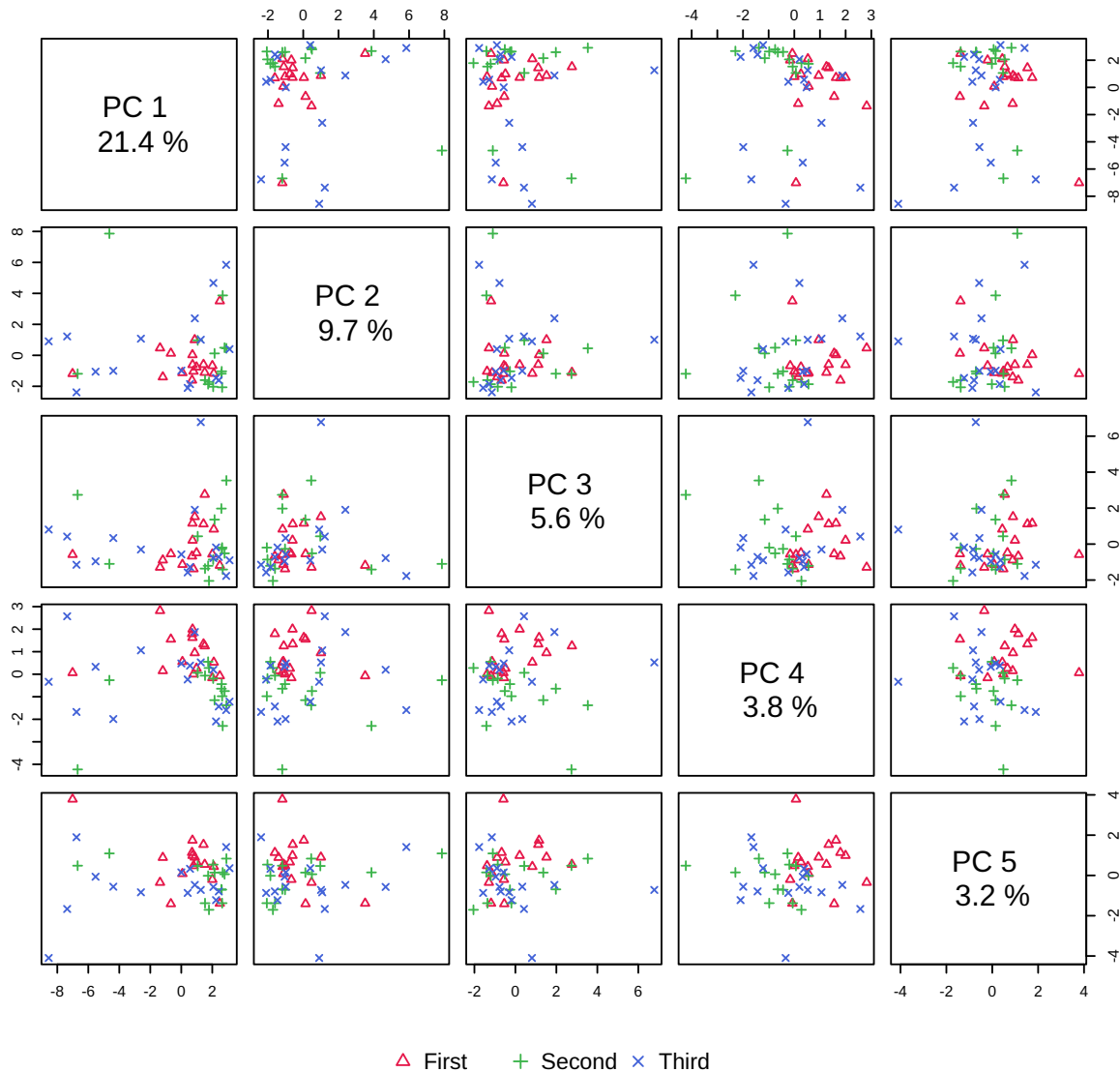


Figure 3: Pairwise score plots between the selected PCs. The explained variance of each PC is shown in the corresponding diagonal cell.

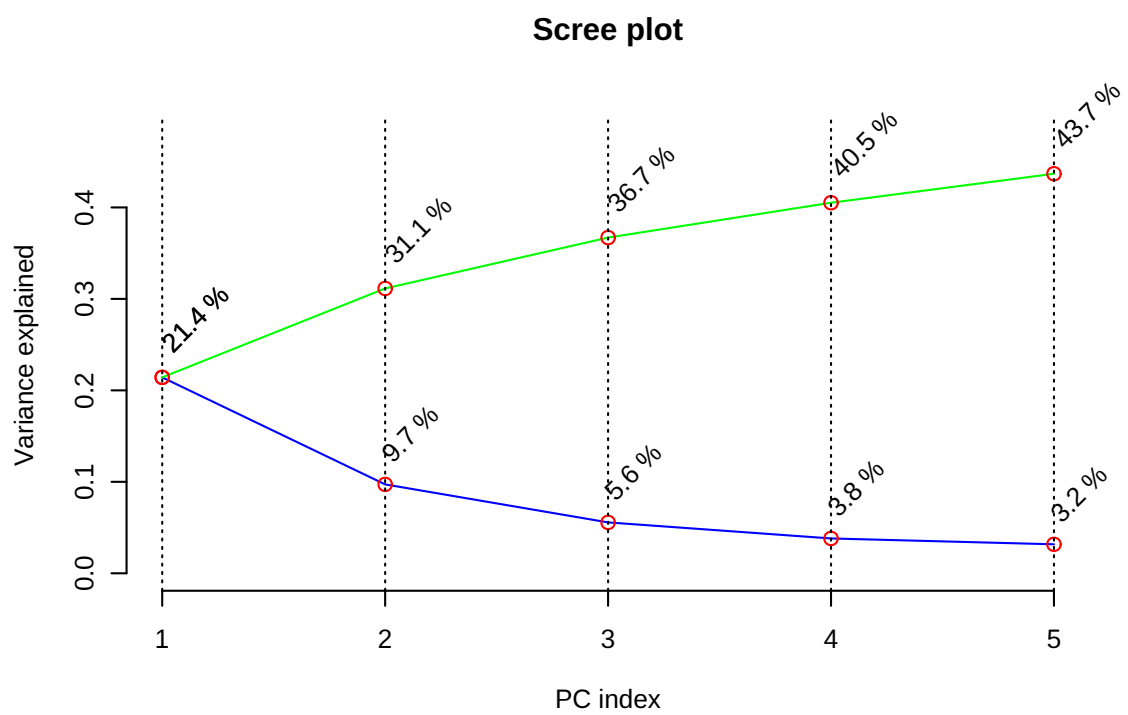
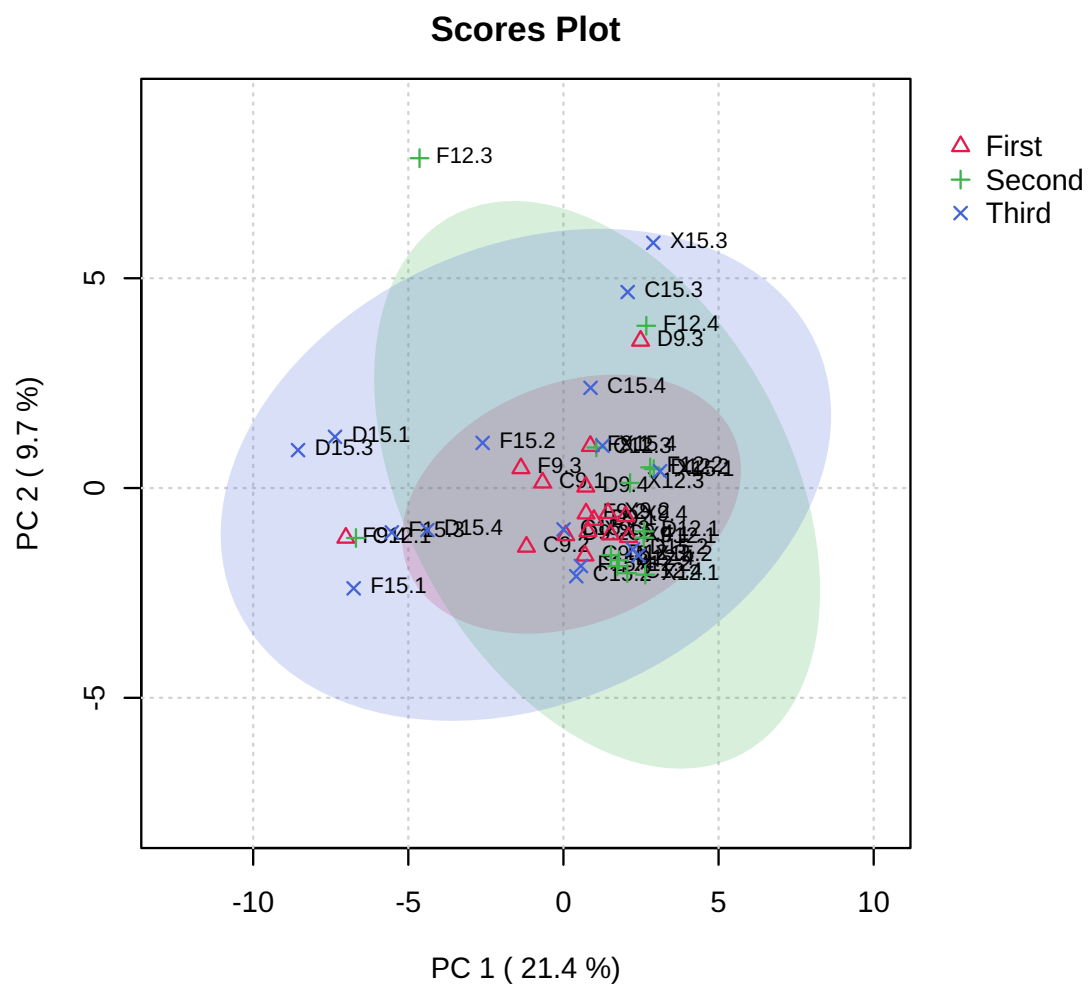


Figure 4: Scree plot shows the variance explained by PCs. The green line on top shows the accumulated variance explained; the blue line underneath shows the variance explained by individual PC.



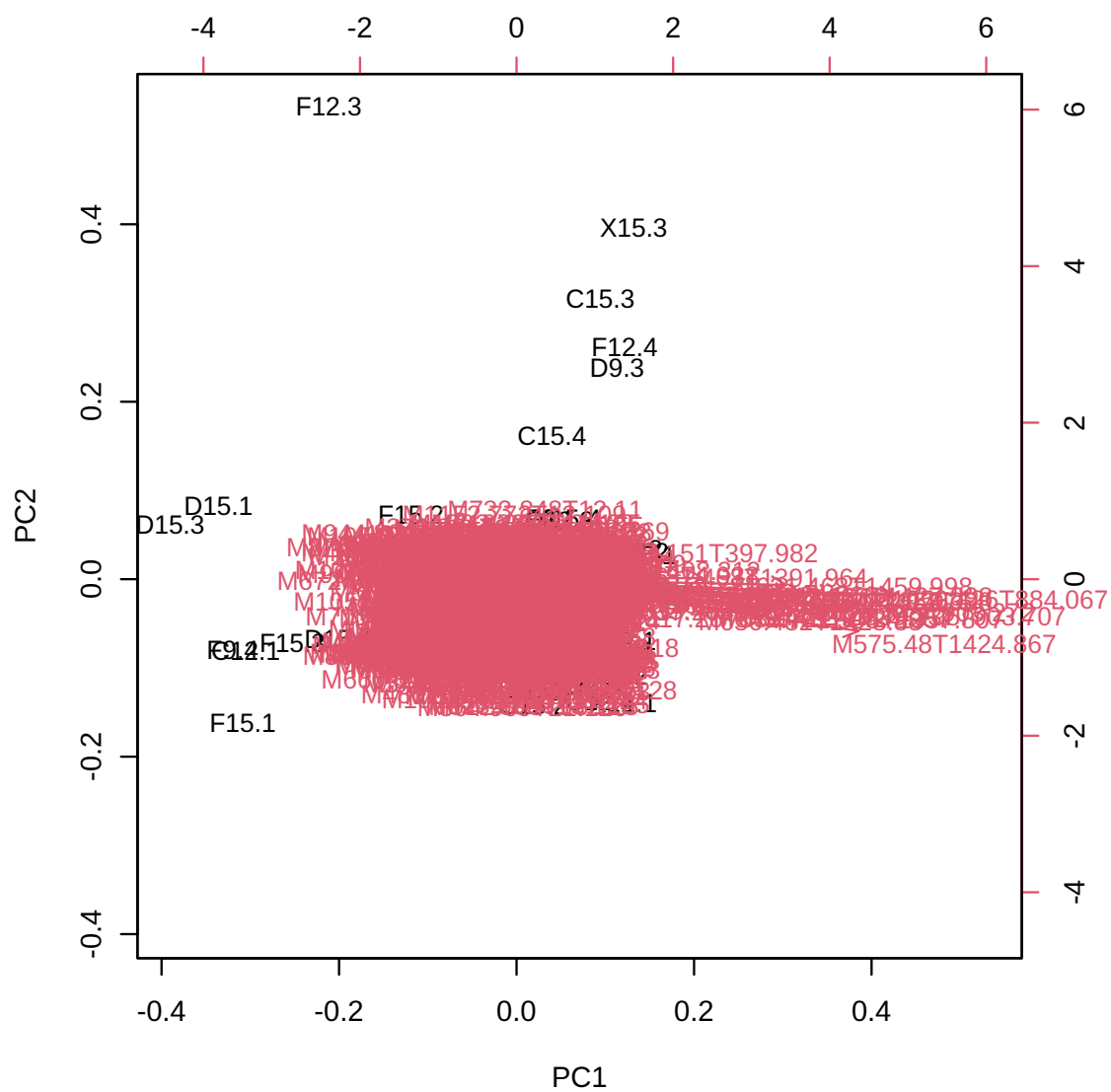


Figure 6: PCA biplot between the selected PCs. Note, you may want to test different centering and scaling normalization methods for the biplot to be displayed properly.

2.3 Hierarchical Clustering

In (agglomerative) hierarchical cluster analysis, each sample begins as a separate cluster and the algorithm proceeds to combine them until all samples belong to one cluster. Two parameters need to be considered when performing hierarchical clustering. The first one is similarity measure - Euclidean distance, Pearson's correlation, Spearman's rank correlation. The other parameter is clustering algorithms, including average linkage (clustering uses the centroids of the observations), complete linkage (clustering uses the farthest pair of observations between the two groups), single linkage (clustering uses the closest pair of observations) and Ward's linkage (clustering to minimize the sum of squares of any two clusters). Heatmap is often presented as a visual aid in addition to the dendrogram.

Hierarchical clustering is performed with the `hclust` function in package `stat`. Figure 7 shows the clustering result in the form of a dendrogram. Figure 8 shows the clustering result in the form of a heatmap.

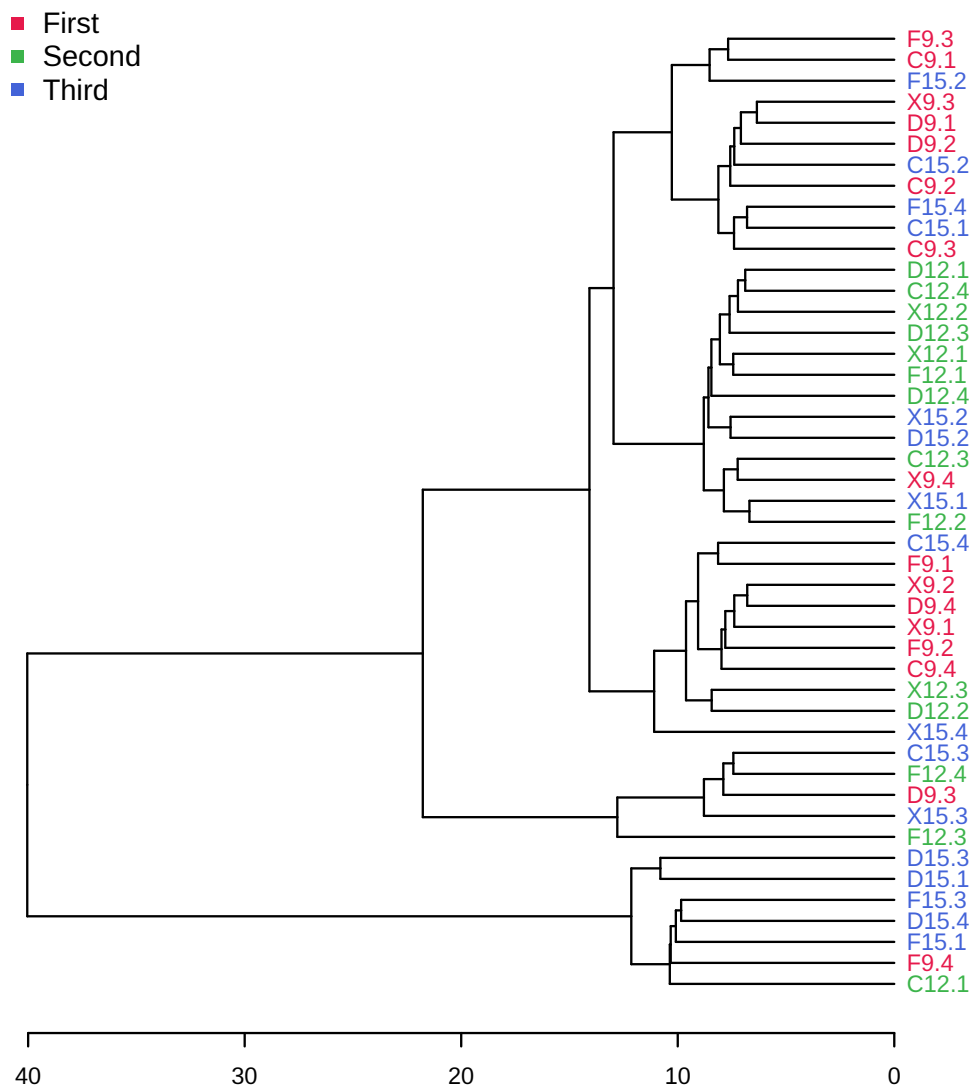


Figure 7: Clustering result shown as dendrogram (distance measure using `euclidean`, and clustering algorithm using `ward.D`).

3 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"pktable\", \"stat\", FALSE)"
[2] "mSet<-Read.TextData(mSet, \"Replacing_with_your_file_path\", \"rowu\", \"disc\");"
[3] "mSet<-SanityCheckData(mSet)"
[4] "mSet<-ReplaceMin(mSet);"
[5] "mSet<-SanityCheckData(mSet)"
[6] "mSet<-FilterVariable(mSet, \"F\", 25, \"iqr\", 0)"
[7] "mSet<-PreparePrenormData(mSet)"
[8] "mSet<-GetGroupNames(mSet, \"\")"
[9] "feature.nm.vec <- c(\"\")"
[10] "smp.nm.vec <- c(\"\")"
[11] "grp.nm.vec <- c(\"First\", \"Second\", \"Third\")"
[12] "mSet<-UpdateData(mSet, T)"
[13] "mSet<-PreparePrenormData(mSet)"
[14] "mSet<-Normalization(mSet, \"CompNorm\", \"LogNorm\", \"ParetoNorm\", \"Sodium_Formate\", ratio
[15] "mSet<-PlotNormSummary(mSet, \"norm_0\", \"png\", 72, width=NA)"
[16] "mSet<-PlotSampleNormSummary(mSet, \"snorm_0\", \"png\", 72, width=NA)"
[17] "mSet<-PlotHeatMap(mSet, \"heatmap_0\", \"png\", 72, width=NA, \"norm\", \"row\", \"euclidean\"
[18] "mSet<-PlotHeatMap(mSet, \"heatmap_1\", \"png\", 72, width=NA, \"norm\", \"row\", \"euclidean\"
[19] "mSet<-ANOVA.Anal(mSet, F, 0.05, FALSE)"
[20] "mSet<-PlotANOVA(mSet, \"aov_0\", \"png\", 72, width=NA)"
[21] "mSet<-ANOVA.Anal(mSet, F, 0.1043, FALSE)"
[22] "mSet<-PlotANOVA(mSet, \"aov_1\", \"png\", 72, width=NA)"
[23] "mSet<-ANOVA.Anal(mSet, F, 1.0, FALSE)"
[24] "mSet<-PlotANOVA(mSet, \"aov_2\", \"png\", 72, width=NA)"
[25] "mSet<-Calculate.ANOVA.posthoc(mSet, \"tukey\", 0.05)"
[26] "mSet<-Calculate.ANOVA.posthoc(mSet, \"tukey\", 0.05)"
[27] "mSet<-ANOVA.Anal(mSet, F, 0.95057, FALSE)"
[28] "mSet<-PlotANOVA(mSet, \"aov_3\", \"png\", 72, width=NA)"
[29] "mSet<-Calculate.ANOVA.posthoc(mSet, \"tukey\", 0.05)"
[30] "mSet<-PCA.Anal(mSet)"
[31] "mSet<-PlotPCAPairSummary(mSet, \"pca_pair_0\", \"png\", 72, width=NA, 5)"
[32] "mSet<-PlotPCAScree(mSet, \"pca_scee_0\", \"png\", 72, width=NA, 5)"
[33] "mSet<-PlotPCA2DScore(mSet, \"pca_score2d_0\", \"png\", 72, width=NA, 1,2,0.95,0,0)"
[34] "mSet<-PlotPCALoading(mSet, \"pca_loading_0\", \"png\", 72, width=NA, 1,2);"
[35] "mSet<-PlotPCABiplot(mSet, \"pca_biplot_0\", \"png\", 72, width=NA, 1,2)"
[36] "mSet<-PlotPCA3DLoading(mSet, \"pca_loading3d_0\", \"json\", 1,2,3)"
[37] "mSet<-PlotPCA2DScore(mSet, \"pca_score2d_1\", \"png\", 72, width=NA, 1,2,0.95,1,0)"
[38] "mSet<-PlotSubHeatMap(mSet, \"heatmap_2\", \"png\", 72, width=NA, \"norm\", \"row\", \"euclidean\"
[39] "mSet<-PlotHCTree(mSet, \"tree_0\", \"png\", 72, width=NA, \"euclidean\", \"ward.D\")"
[40] "mSet<-SaveTransformedData(mSet)"
[41] "mSet<-PreparePDFReport(mSet, \"guest15282839817347858290\")\n"
```

The report was generated on Wed Nov 29 12:12:19 2023 with R version 4.2.2 (2022-10-31), OS system:
Linux, version: -Ubuntu SMP Mon May 15 15:18:26 UTC 2023 .