

Metabolomic Data Analysis with MetaboAnalyst 5.0

Name: guest1848273810588172982

April 27, 2022

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 columns and features in rows. The uploaded file is in comma separated values (.csv) format. The uploaded data file contains 9 (samples) by 30 (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 filtering was applied

Table 1: Summary of data processing results

	Features (positive)	Missing/Zero	Features (processed)
CO_1	30	0	30
CO_2	30	0	30
CO_3	30	0	30
KY_1	30	0	30
KY_2	30	0	30
KY_3	30	0	30
OR_1	30	0	30
OR_2	30	0	30
OR_3	30	0	30

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

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

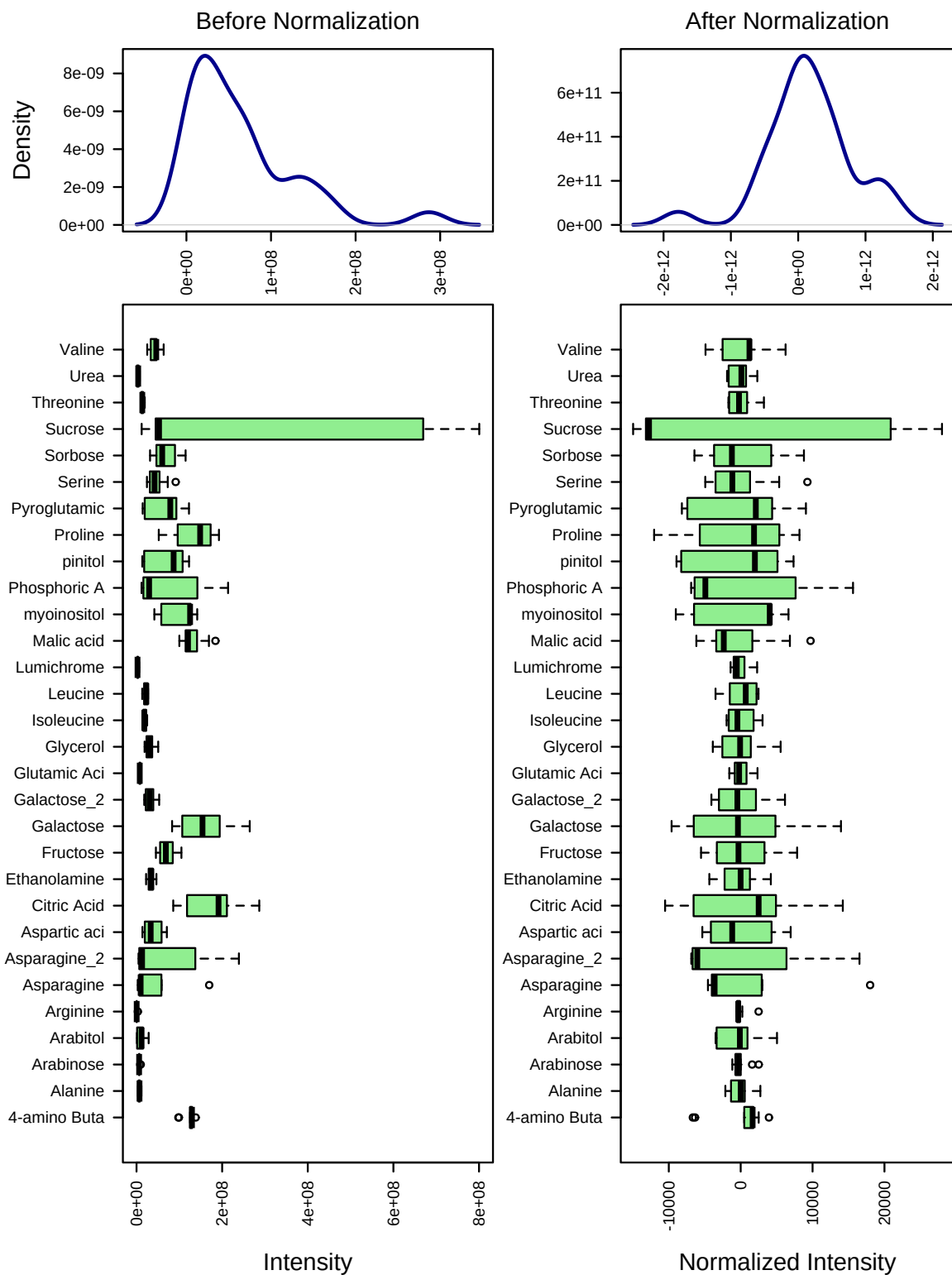


Figure 1: Box plots and kernel density plots before and after normalization. The boxplots show at most 50 features due to space limit. The density plots are based on all samples. Selected methods : Row-wise normalization: N/A; Data transformation: N/A; Data scaling: Pareto Scaling.

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.

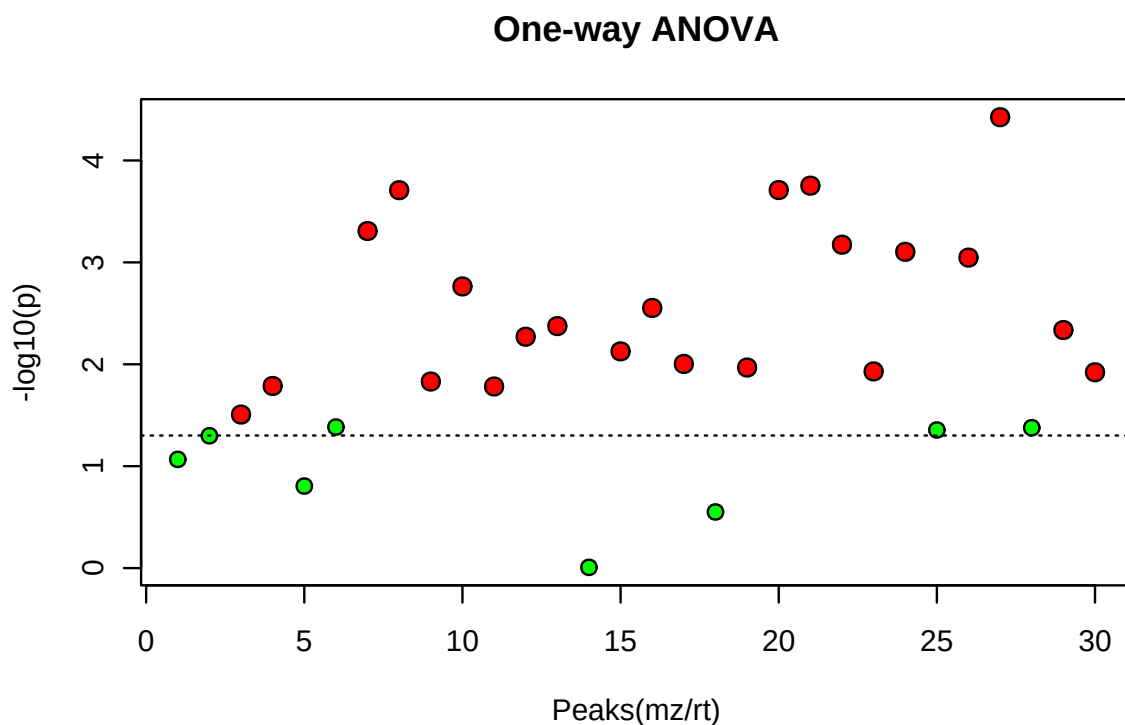


Figure 2: Important features selected by ANOVA plot with p value threshold 0.05.

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

	Peaks(mz/rt)	f.value	p.value	-log10(p)	FDR	Fisher's LSD
1	Sucrose	86.54	0.00	4.42	0.00	OR - CO; OR - KY
2	Phosphoric Acid	50.46	0.00	3.75	0.00	KY - CO; KY - OR
3	myoinositol	48.73	0.00	3.71	0.00	CO - KY; OR - KY
4	Aspartic acid	48.67	0.00	3.71	0.00	KY - CO; CO - OR; KY - OR
5	Asparagine_2	34.99	0.00	3.31	0.00	KY - CO; KY - OR
6	pinitol	31.27	0.00	3.17	0.00	CO - OR; KY - OR
7	Pyroglutamic Acid	29.48	0.00	3.10	0.00	CO - OR; KY - OR
8	Sorbose	28.11	0.00	3.05	0.00	CO - KY; CO - OR; OR - KY
9	Ethanolamine	22.03	0.00	2.76	0.01	KY - CO; CO - OR; KY - OR
10	Isoleucine	18.29	0.00	2.55	0.01	KY - CO; KY - OR
11	Galactose_2	15.57	0.00	2.37	0.01	CO - KY; CO - OR; OR - KY
12	Urea	15.02	0.00	2.34	0.01	CO - OR; KY - OR
13	Galactose	14.14	0.01	2.27	0.01	CO - KY; CO - OR; OR - KY
14	Glycerol	12.35	0.01	2.13	0.02	CO - KY; CO - OR
15	Leucine	10.96	0.01	2.00	0.02	KY - CO; OR - CO
16	Malic acid	10.59	0.01	1.97	0.02	KY - CO; KY - OR
17	Proline	10.20	0.01	1.93	0.02	CO - OR; KY - OR
18	Valine	10.11	0.01	1.92	0.02	CO - OR; KY - OR
19	Citric Acid	9.23	0.01	1.83	0.02	CO - OR; KY - OR
20	Arabitol	8.83	0.02	1.79	0.02	CO - OR; KY - OR
21	Fructose	8.78	0.02	1.78	0.02	CO - KY; CO - OR
22	Arabinose	6.53	0.03	1.51	0.04	CO - KY; CO - OR

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 3-D scores plot between selected PCs; Figure 7 shows the loadings plot between the selected PCs; Figure 8 shows the biplot between the selected PCs.

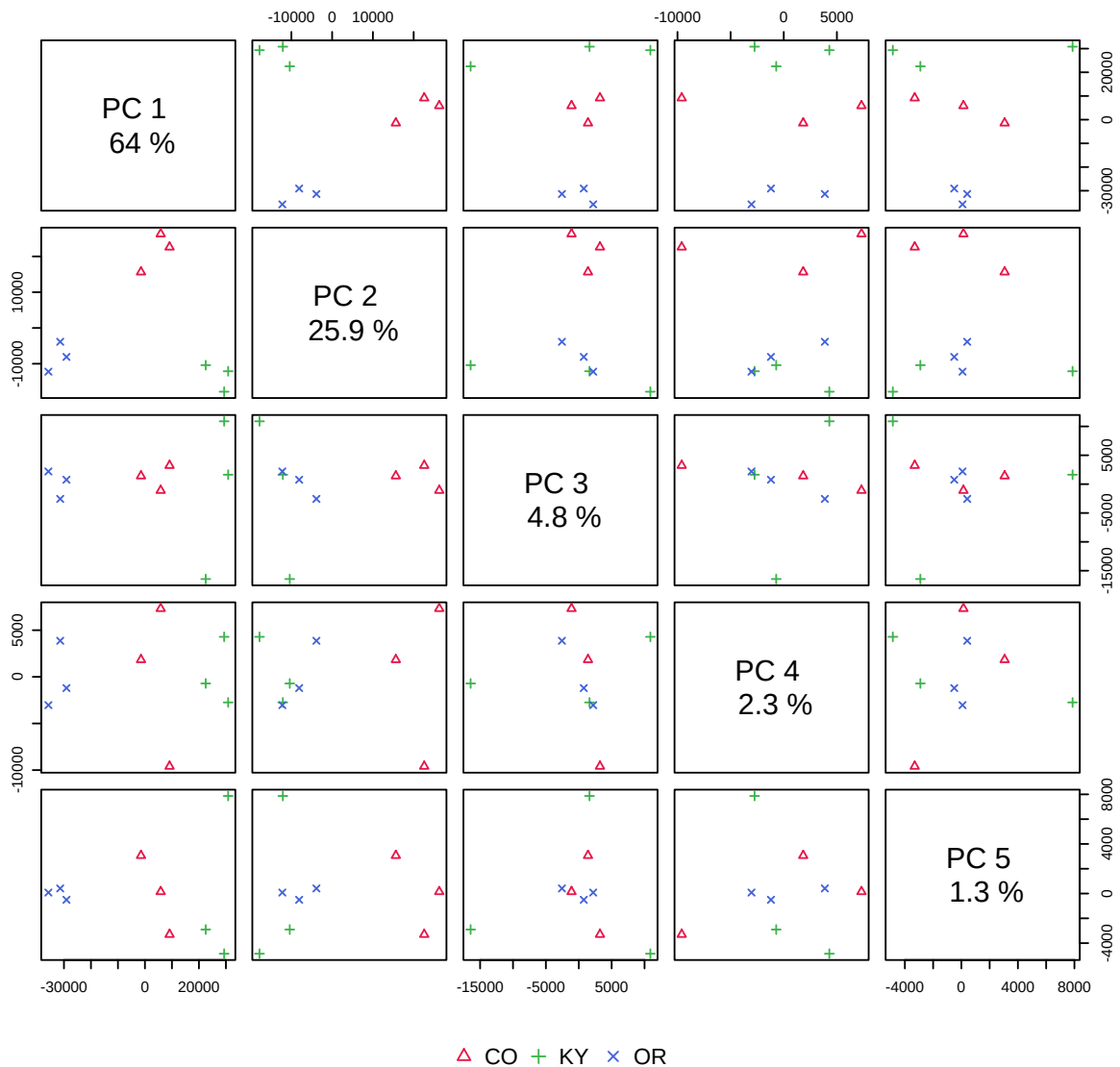


Figure 3: Pairwise score plots between the selected PCs. The explained variance of each PC is shown in the corresponding diagonal cell.

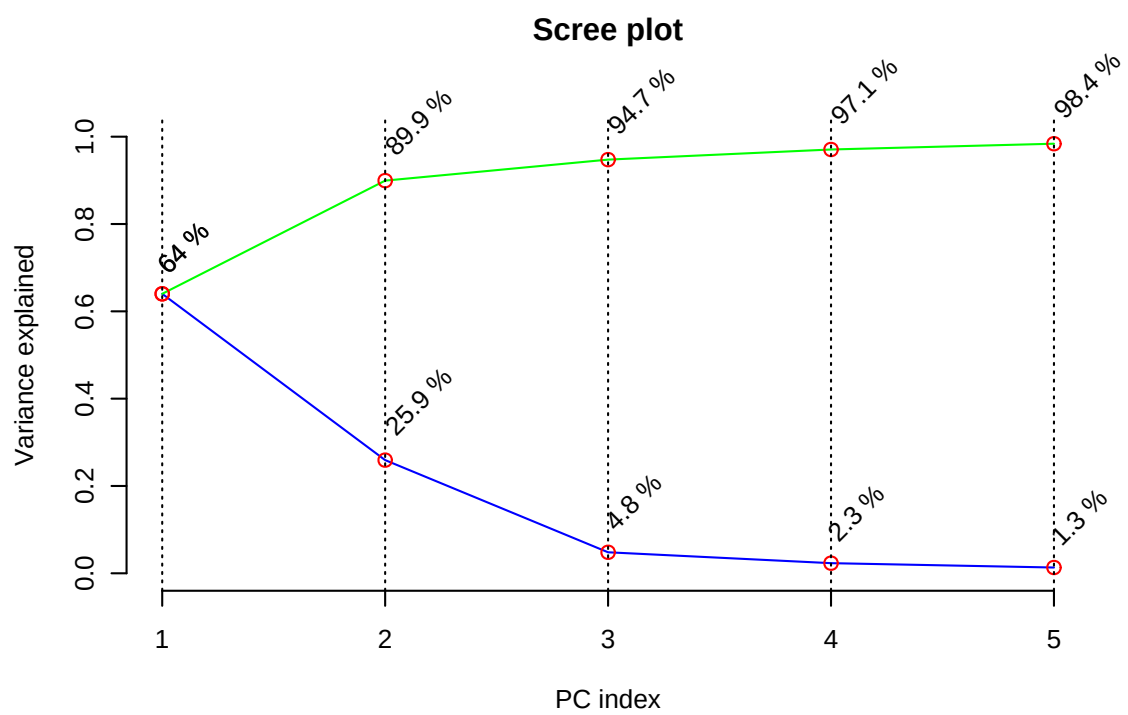


Figure 4: Scree plot shows the variance explained by PCs. The green line on top shows the accumulated variance explained; the blue line underneath shows the variance explained by individual PC.

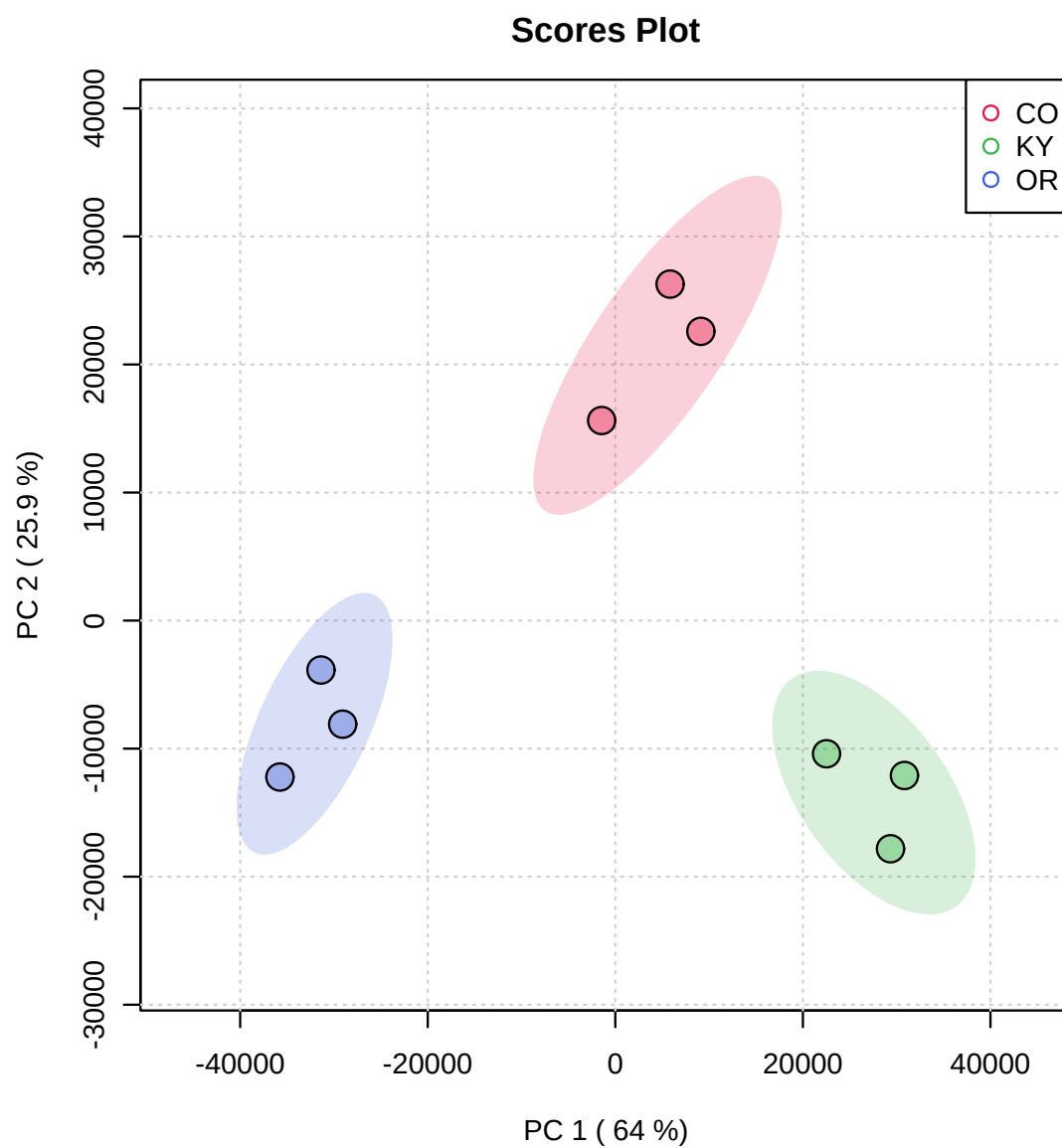


Figure 5: Scores plot between the selected PCs. The explained variances are shown in brackets.

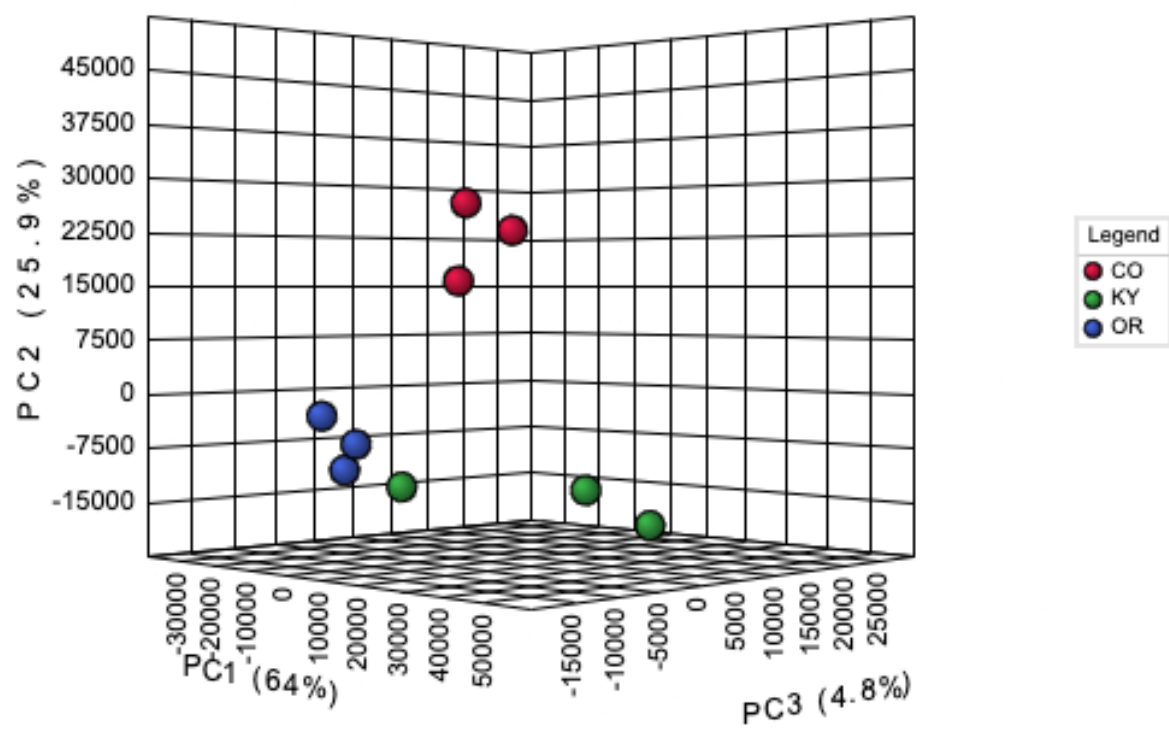


Figure 6: 3D score plot between the selected PCs. The explained variances are shown in brackets.

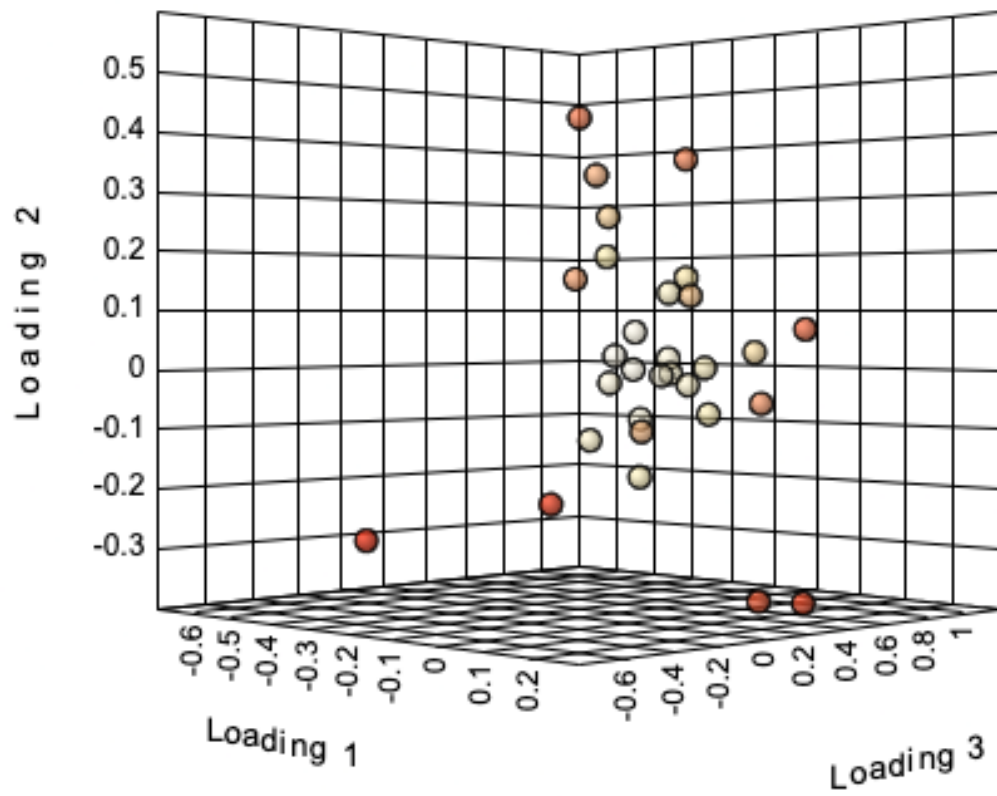


Figure 7: Loadings plot for the selected PCs.

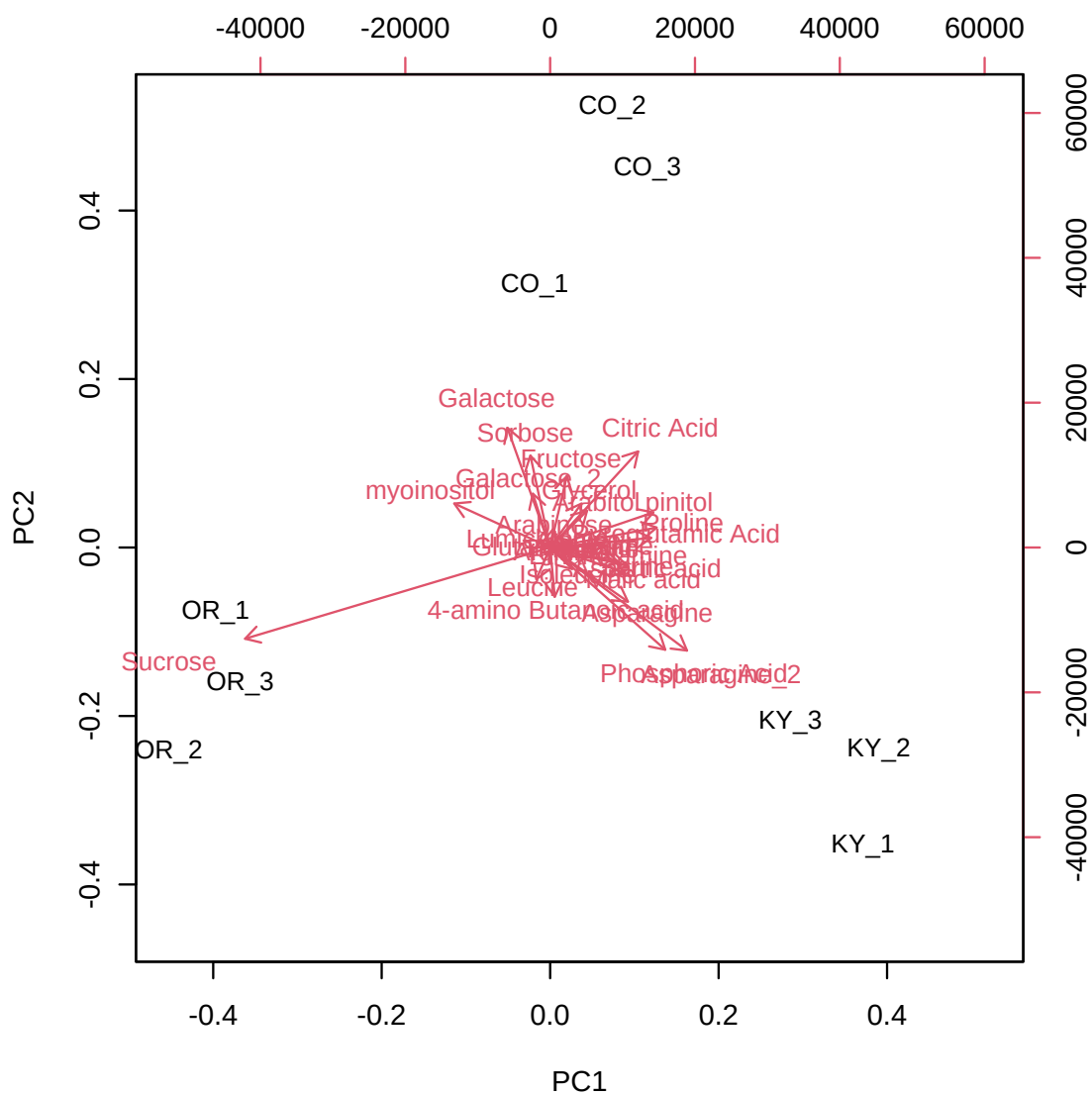


Figure 8: PCA biplot between the selected PCs. Note, you may want to test different centering and scaling normalization methods for the biplot to be displayed properly.

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 9 shows the clustering result in the form of a dendrogram. Figure 10 shows the clustering result in the form of a heatmap.

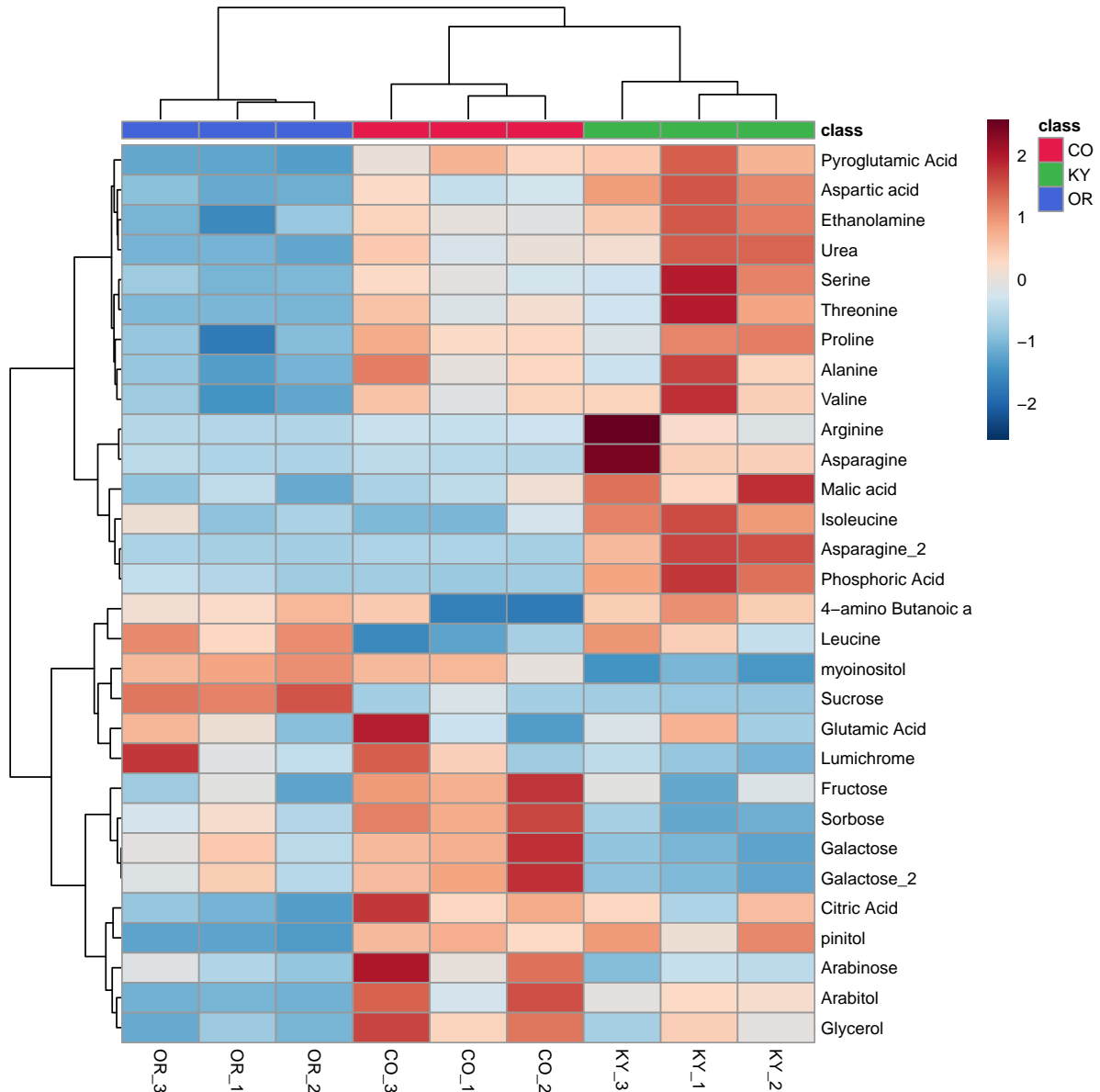


Figure 9: Clustering result shown as heatmap (distance measure using `euclidean`, and clustering algorithm using `ward.D`).

3 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"pktable\", \"stat\", FALSE)"
[2] "mSet<-Read.TextData(mSet, \"Replacing_with_your_file_path\", \"colu\", \"disc\");"
[3] "mSet<-SanityCheckData(mSet)"
[4] "mSet<-ReplaceMin(mSet);"
[5] "mSet<-SanityCheckData(mSet)"
[6] "mSet<-FilterVariable(mSet, \"none\", \"F\", 25)"
[7] "mSet<-PreparePrenormData(mSet)"
[8] "mSet<-Normalization(mSet, \"NULL\", \"NULL\", \"ParetoNorm\", ratio=FALSE, ratioNum=20)"
[9] "mSet<-PlotNormSummary(mSet, \"norm_0_\", \"png\", 72, width=NA)"
[10] "mSet<-PlotSampleNormSummary(mSet, \"snorm_0_\", \"png\", 72, width=NA)"
[11] "mSet<-PCA.Anal(mSet)"
[12] "mSet<-PlotPCAPairSummary(mSet, \"pca_pair_0_\", \"png\", 72, width=NA, 5)"
[13] "mSet<-PlotPCAScree(mSet, \"pca_screes_0_\", \"png\", 72, width=NA, 5)"
[14] "mSet<-PlotPCA2DScore(mSet, \"pca_score2d_0_\", \"png\", 72, width=NA, 1,2,0.95,0,0)"
[15] "mSet<-PlotPCALoading(mSet, \"pca_loading_0_\", \"png\", 72, width=NA, 1,2);"
[16] "mSet<-PlotPCABiplot(mSet, \"pca_biplot_0_\", \"png\", 72, width=NA, 1,2)"
[17] "mSet<-PlotPCA3DLoading(mSet, \"pca_loading3d_0_\", \"json\", 1,2,3)"
[18] "mSet<-ANOVA.Anal(mSet, F, 0.05, \"fisher\", FALSE)"
[19] "mSet<-PlotANOVA(mSet, \"aov_0_\", \"png\", 72, width=NA)"
[20] "mSet<-PlotHeatMap(mSet, \"heatmap_0_\", \"png\", 72, width=NA, \"norm\", \"row\", \"euclidean\"
[21] "mSet<-SaveTransformedData(mSet)"
[22] "mSet<-PreparePDFReport(mSet, \"guest1848273810588172982\")\n"
```

The report was generated on Wed Apr 27 23:37:36 2022 with R version 4.0.2 (2020-06-22).