

Metabolomic Data Analysis with MetaboAnalyst 6.0

Name: guest5646536859053858679

March 20, 2024

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 Binned Spectral Data

The binned spectra data 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 50 (samples) by 200 (spectra bins) 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 was performed.

Table 1: Summary of data processing results

	Features (positive)	Missing/Zero	Features (processed)
C002	194	6	200
C004	189	11	200
C005	191	9	200
C006	195	5	200
C007	200	0	200
C009	186	14	200
C010	196	4	200
C011	177	23	200
C012	189	11	200
C015	188	12	200
C016	188	12	200
C017	198	2	200
C019	181	19	200
C020	184	16	200
C021	187	13	200
C022	191	9	200
C024	190	10	200
C026	195	5	200
C028	196	4	200
C029	192	8	200
C030	182	18	200
C031	179	21	200
C032	191	9	200
C033	189	11	200
C034	199	1	200
P002	195	5	200
P012	187	13	200
P014	200	0	200
P027	200	0	200
P034	198	2	200
P037	187	13	200
P038	195	5	200
P041	178	22	200
P042	198	2	200
P049	189	11	200
P056	190	10	200
P058	179	21	200
P060	190	10	200
P064	200	0	200
P065	198	2	200
P070	190	10	200
P080	196	4	200
P085	200	0	200
P086	193	7	200
P089	199	1	200
P092	191	9	200
P099	190	10	200
P113	152	48	200
P013b	191	9	200
P100b	199	1	200

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.

Row-wise normalization: Normalization to sample median; Data transformation: Log10 Normalization; Data scaling: Autoscaling.

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

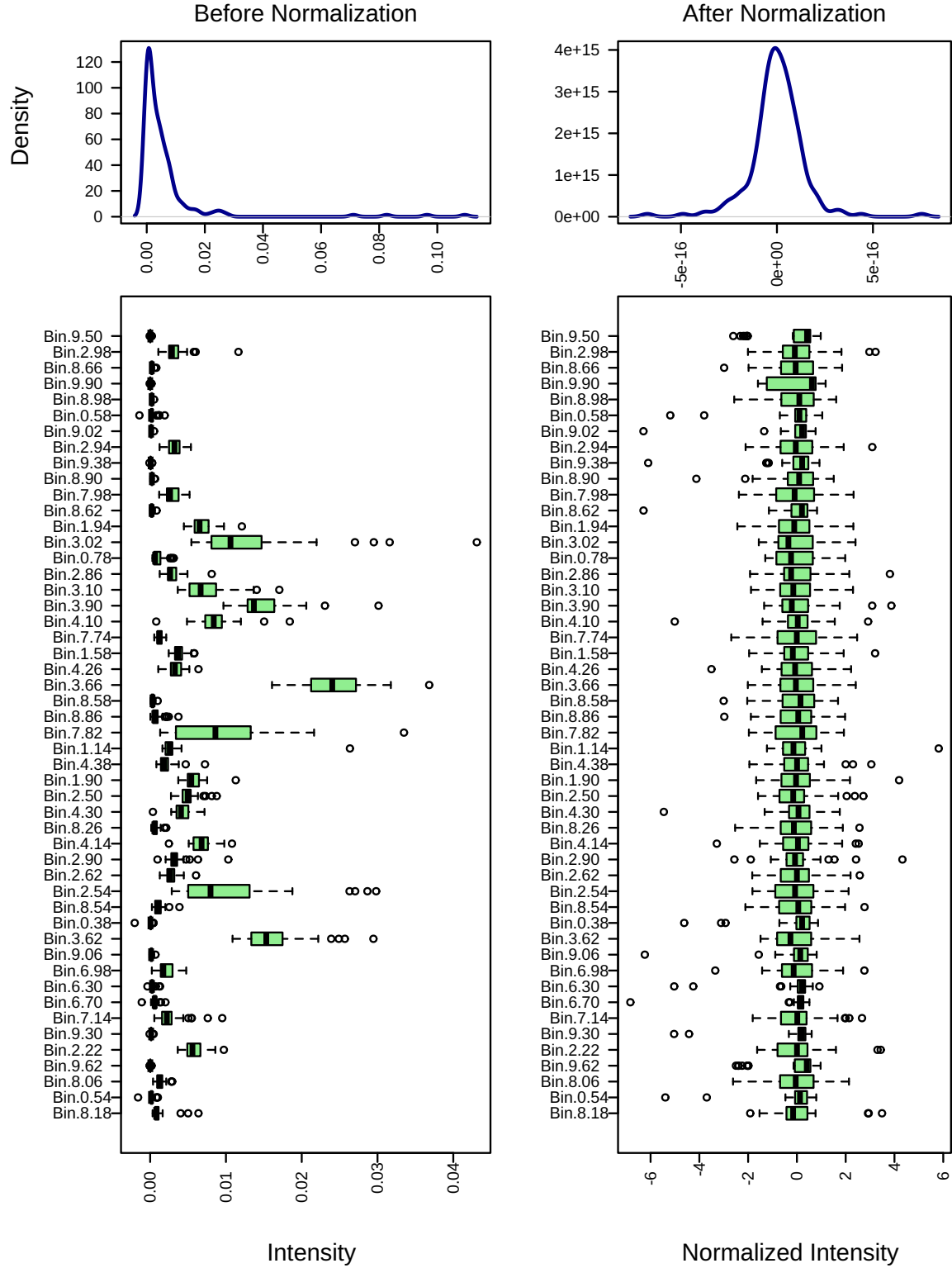


Figure 1: Box plots and kernel density plots before and after normalization. The boxplots show at most 50 features due to space limit. The density plots are based on all samples.

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 Univariate Analysis

Univariate analysis methods are the most common methods used for exploratory data analysis. For two-group data, MetaboAnalyst provides Fold Change (FC) analysis, t-tests, and volcano plot which is a combination of the first two methods. All three these methods support both unpaired and paired analyses. For multi-group analysis, MetaboAnalyst provides two types of analysis - one-way analysis of variance (ANOVA) with associated post-hoc analyses, and correlation analysis to identify significant compounds that follow a given pattern. The univariate analyses provide a preliminary overview about features that are potentially significant in discriminating the conditions under study.

For paired fold change analysis, the algorithm first counts the total number of pairs with fold changes that are consistently above/below the specified FC threshold for each variable. A variable will be reported as significant if this number is above a given count threshold (default > 75% of pairs/variable)

Figure 2 shows the important features identified by t-tests. Table 2 shows the details of these features;

Please note, the purpose of fold change is to compare absolute value changes between two group means. Therefore, the data before column normalization will be used instead. Also note, the result is plotted in log2 scale, so that same fold change (up/down regulated) will have the same distance to the zero baseline.

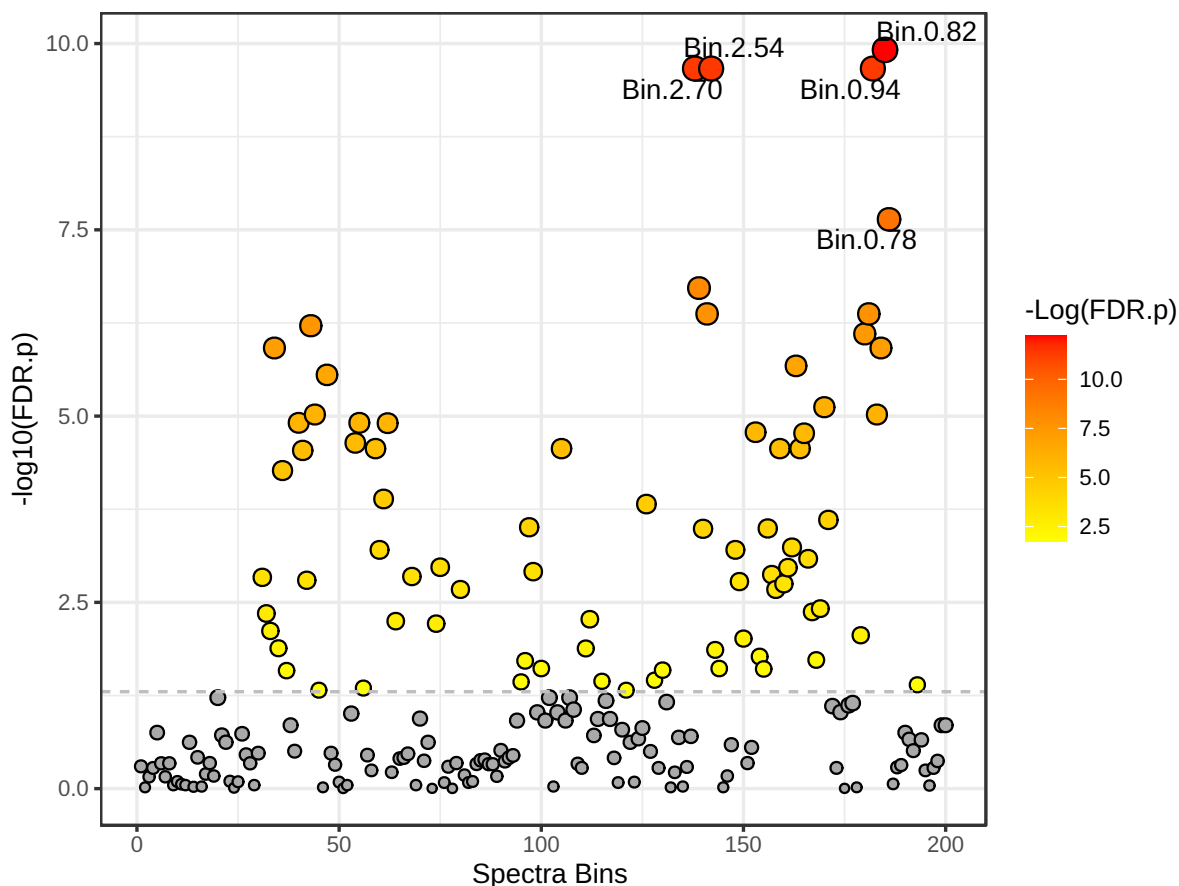


Figure 2: Important features selected by t-tests with threshold 0.05. The red circles represent features above the threshold. Note the p values are transformed by $-\log_{10}$ so that the more significant features (with smaller p values) will be plotted higher on the graph.

Table 2: Top 50 features identified by t-tests

	Spectra Bins	t.stat	p.value	-log10(p)	FDR
1	Bin.0.82	-9.7339	6.0795e-13	12.216	1.2159e-10
2	Bin.2.54	9.1867	3.7747e-12	11.423	2.1653e-10
3	Bin.2.70	9.1659	4.0479e-12	11.393	2.1653e-10
4	Bin.0.94	-9.1459	4.3307e-12	11.363	2.1653e-10
5	Bin.0.78	-7.7267	5.7137e-10	9.2431	2.2855e-08
6	Bin.2.66	7.0703	5.7385e-09	8.2412	1.9128e-07
7	Bin.2.58	6.78	1.5993e-08	7.7961	4.2491e-07
8	Bin.0.98	-6.7627	1.6996e-08	7.7696	4.2491e-07
9	Bin.8.30	-6.6261	2.7548e-08	7.5599	6.1218e-07
10	Bin.1.02	-6.5246	3.9423e-08	7.4042	7.8847e-07
11	Bin.8.66	-6.3556	7.1612e-08	7.145	1.2189e-06
12	Bin.0.86	-6.3496	7.3136e-08	7.1359	1.2189e-06
13	Bin.1.70	-6.1702	1.3771e-07	6.861	2.1186e-06
14	Bin.8.14	-6.0703	1.9581e-07	6.7082	2.7972e-06
15	Bin.1.42	-5.7662	5.6951e-07	6.2445	7.5935e-06
16	Bin.8.26	-5.667	8.0564e-07	6.0939	9.5236e-06
17	Bin.0.90	-5.6656	8.0951e-07	6.0918	9.5236e-06
18	Bin.7.82	5.5737	1.1151e-06	5.9527	1.2309e-05
19	Bin.8.42	-5.5601	1.1693e-06	5.9321	1.2309e-05
20	Bin.7.54	5.5423	1.2441e-06	5.9051	1.2441e-05
21	Bin.2.10	-5.4476	1.7291e-06	5.7622	1.6468e-05
22	Bin.1.62	-5.425	1.8698e-06	5.7282	1.6998e-05
23	Bin.7.86	5.3256	2.6374e-06	5.5788	2.2934e-05
24	Bin.7.66	5.2358	3.5947e-06	5.4443	2.7317e-05
25	Bin.1.86	-5.2299	3.6689e-06	5.4355	2.7317e-05
26	Bin.1.66	-5.2285	3.6868e-06	5.4333	2.7317e-05
27	Bin.3.98	5.2284	3.6878e-06	5.4332	2.7317e-05
28	Bin.8.38	-5.2016	4.0432e-06	5.3933	2.888e-05
29	Bin.8.58	-5.0096	7.803e-06	5.1077	5.3814e-05
30	Bin.7.58	4.74	1.9424e-05	4.7117	0.00012949
31	Bin.3.14	4.6831	2.3506e-05	4.6288	0.00015165
32	Bin.1.38	-4.5265	3.958e-05	4.4025	0.00024738
33	Bin.4.26	-4.4486	5.1192e-05	4.2908	0.00031026
34	Bin.1.98	-4.4282	5.4736e-05	4.2617	0.00032198
35	Bin.2.62	4.416	5.6986e-05	4.2442	0.00032564
36	Bin.1.74	-4.2307	0.00010423	3.982	0.00057905
37	Bin.2.30	-4.1972	0.00011611	3.9351	0.00062411
38	Bin.7.62	4.1907	0.00011858	3.926	0.00062411
39	Bin.1.58	-4.0967	0.00016031	3.795	0.00082213
40	Bin.7.02	-4.0066	0.00021348	3.6706	0.0010674
41	Bin.1.78	-3.9947	0.0002217	3.6542	0.0010814
42	Bin.4.22	-3.9479	0.00025698	3.5901	0.0012237
43	Bin.1.94	-3.9113	0.0002883	3.5401	0.001341
44	Bin.7.30	3.8851	0.00031294	3.5045	0.0014225
45	Bin.8.78	-3.8695	0.00032849	3.4835	0.0014599
46	Bin.8.34	-3.8337	0.00036729	3.435	0.0015969
47	Bin.2.26	-3.8136	0.00039099	3.4078	0.0016638
48	Bin.1.82	-3.7847	0.00042758	3.369	0.0017816
49	Bin.6.82	-3.7191	0.00052323	3.2813	0.002122
50	Bin.1.90	-3.7146	0.0005305	3.2753	0.002122

2.2 Correlation Analysis

Correlation analysis can be used to visualize the overall correlations between different features. It can also be used to identify which features are correlated with a feature of interest. Correlation analysis can also be used to identify if certain features show particular patterns under different conditions. Users first need to define a pattern in the form of a series of hyphenated numbers. For example, in a time-series study with four time points, a pattern of 1-2-3-4 is used to search compounds with increasing the concentration as time changes; while a pattern of 3-2-1-3 can be used to search compounds that decrease at first, then bounce back to the original level.

Figure 3 shows the overall correlation heatmap.

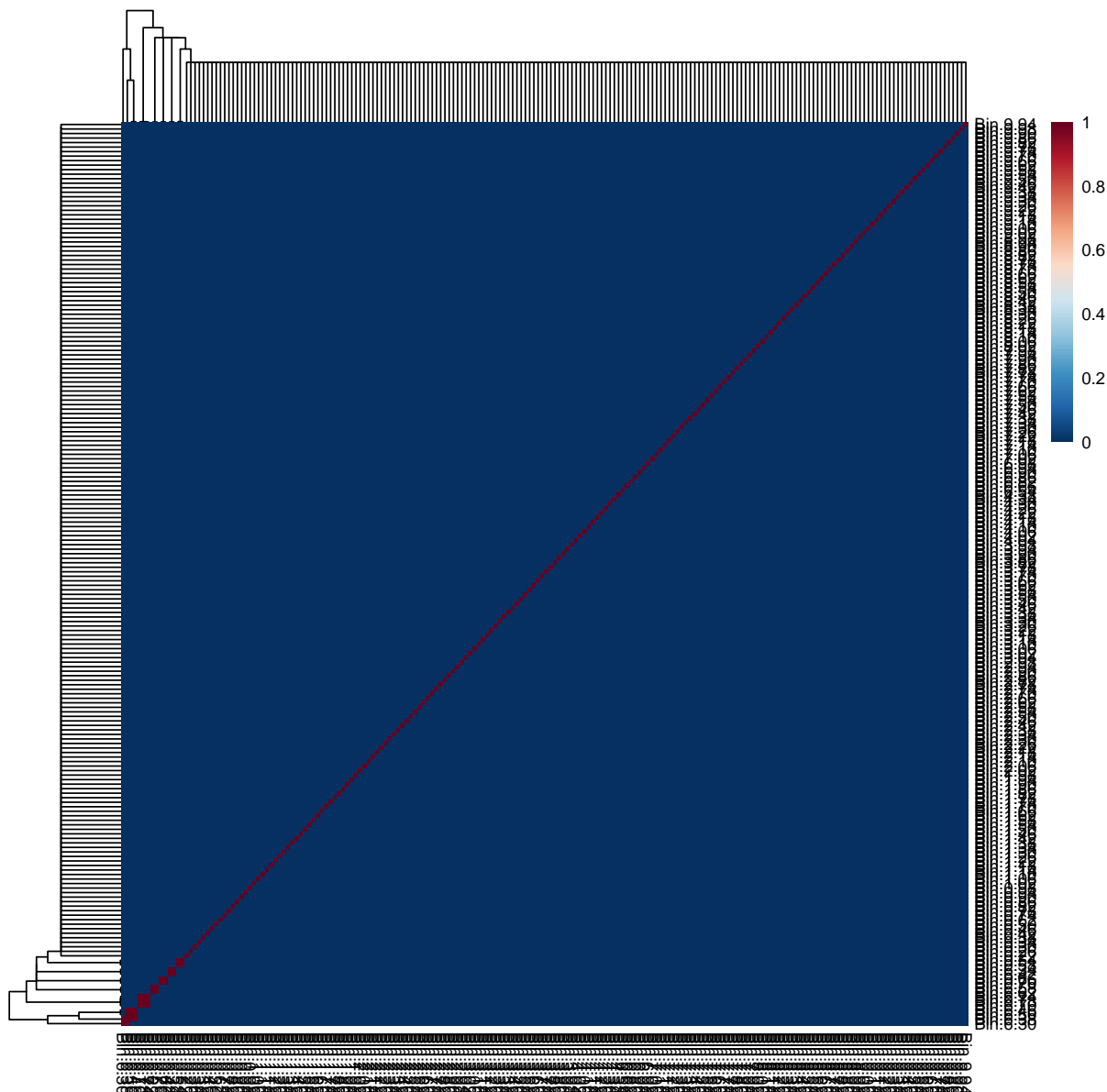


Figure 3: Correlation Heatmaps

2.3 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 4 is pairwise score plots providing an overview of the various separation patterns among the most significant PCs; Figure 5 is the scree plot showing the variances explained by the selected PCs; Figure 6 shows the 2-D scores plot between selected PCs; Figure 7 shows the biplot between the selected PCs. Interactive 3-D scores plots are not included here and can be directly downloaded from website.

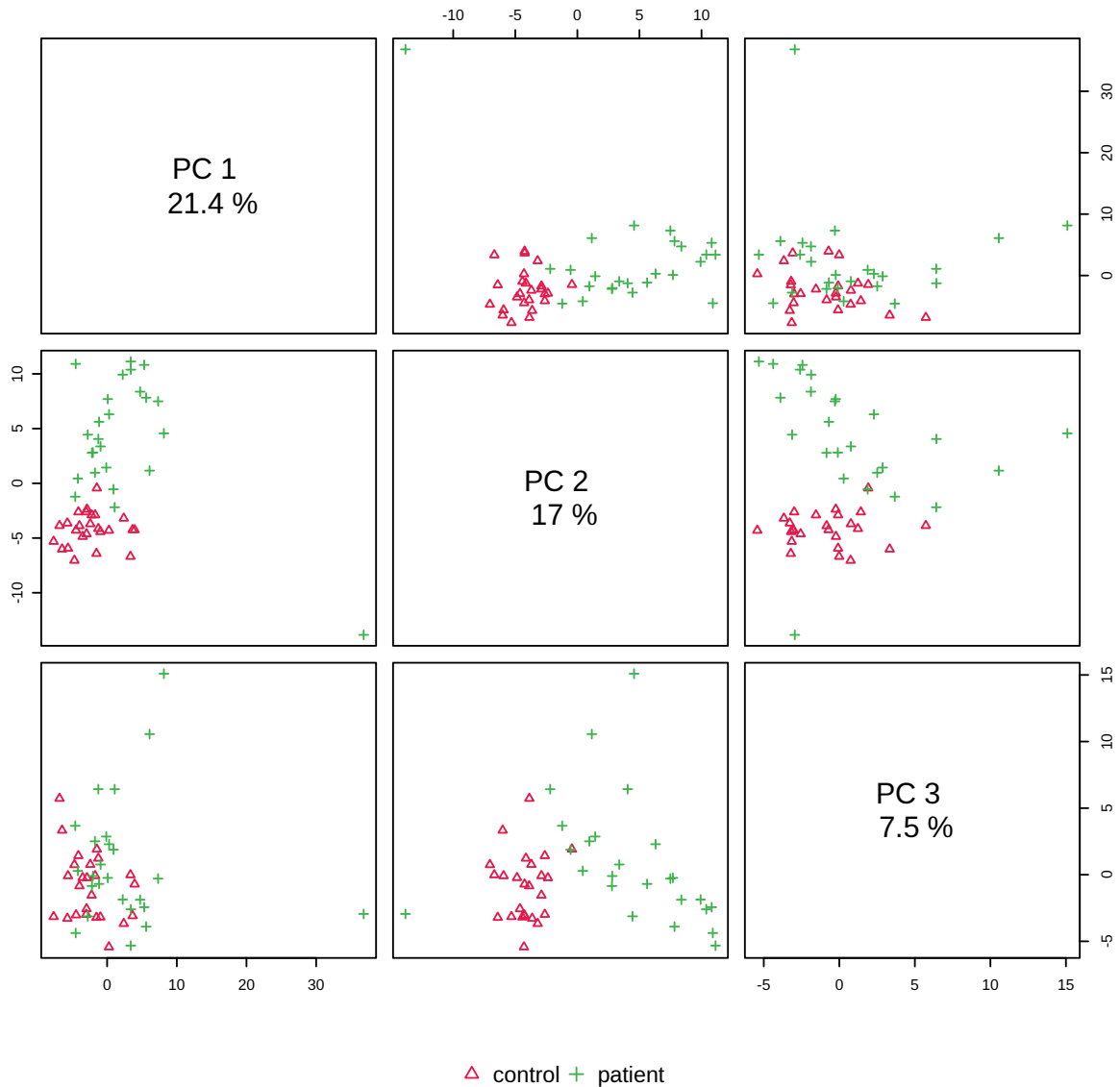


Figure 4: Pairwise score plots between the selected PCs. The explained variance of each PC is shown in the corresponding diagonal cell.

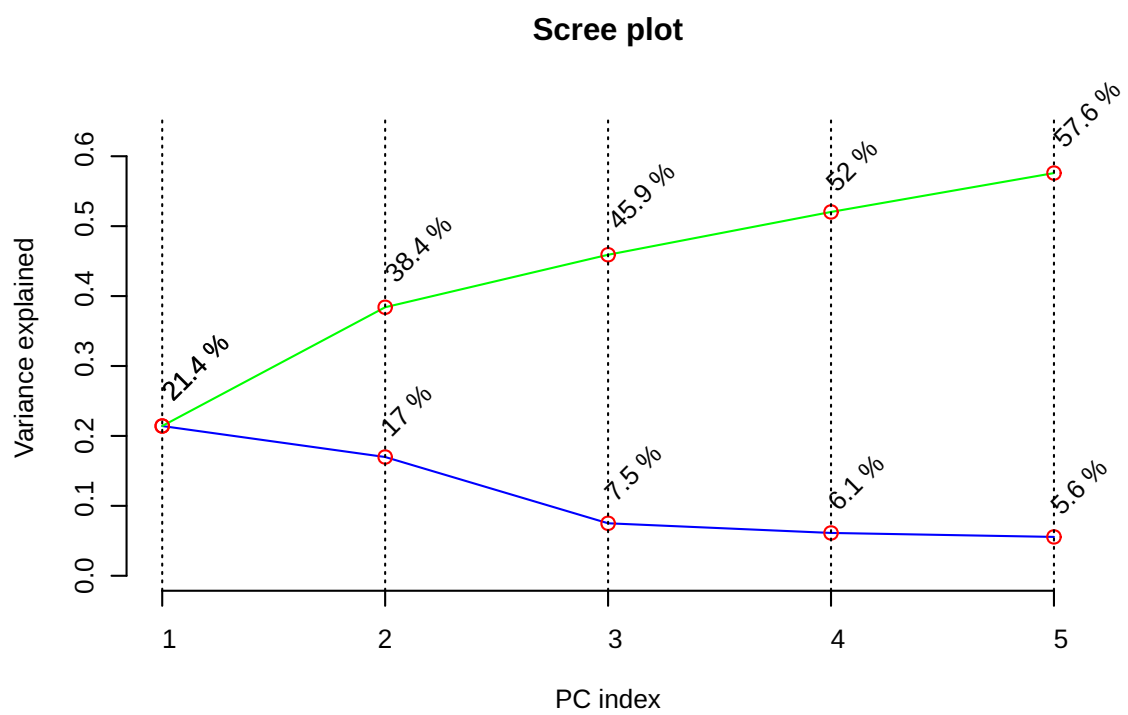


Figure 5: Scree plot shows the variance explained by PCs. The green line on top shows the accumulated variance explained; the blue line underneath shows the variance explained by individual PC.

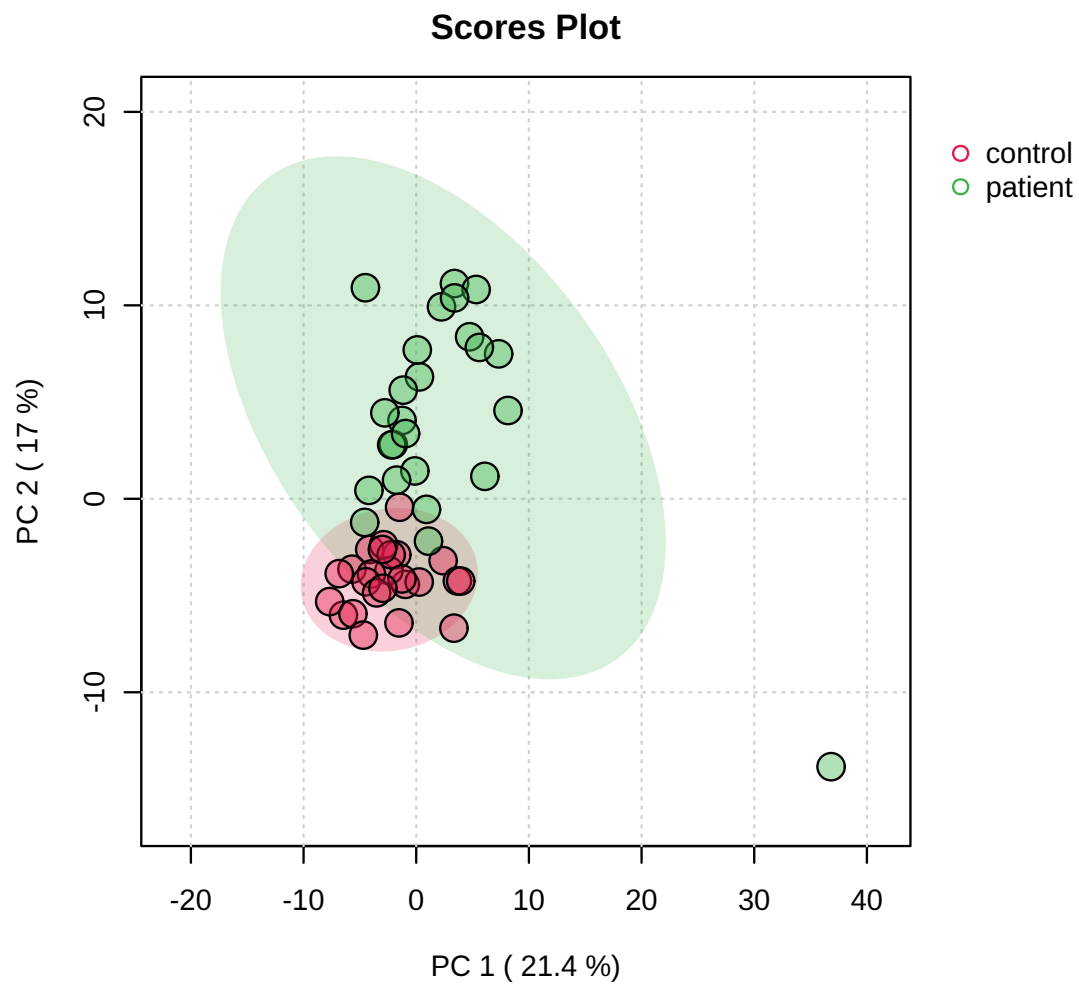


Figure 6: Scores plot between the selected PCs. The explained variances are shown in brackets.

2.4 Partial Least Squares - Discriminant Analysis (PLS-DA)

PLS is a supervised method that uses multivariate regression techniques to extract via linear combination of original variables (X) the information that can predict the class membership (Y). The PLS regression is performed using the `pls` function provided by R `pls` package⁴. The classification and cross-validation are performed using the corresponding wrapper function offered by the `caret` package⁵.

To assess the significance of class discrimination, a permutation test was performed. In each permutation, a PLS-DA model was built between the data (X) and the permuted class labels (Y) using the optimal number of components determined by cross validation for the model based on the original class assignment. MetaboAnalyst supports two types of test statistics for measuring the class discrimination. The first one is based on prediction accuracy during training. The second one is separation distance based on the ratio of the between group sum of the squares and the within group sum of squares (B/W-ratio). If the observed test statistic is part of the distribution based on the permuted class assignments, the class discrimination cannot be considered significant from a statistical point of view.⁶

There are two variable importance measures in PLS-DA. The first, Variable Importance in Projection (VIP) is a weighted sum of squares of the PLS loadings taking into account the amount of explained Y-variation in each dimension. Please note, VIP scores are calculated for each components. When more than components are used to calculate the feature importance, the average of the VIP scores are used. The other importance measure is based on the weighted sum of PLS-regression. The weights are a function of the reduction of the sums of squares across the number of PLS components. Please note, for multiple-group (more than two) analysis, the same number of predictors will be built for each group. Therefore, the coefficient of each feature will be different depending on which group you want to predict. The average of the feature coefficients are used to indicate the overall coefficient-based importance.

Figure 8 shows the overview of scores plots; Figure 9 shows the 2-D scores plot between selected components; Figure 10 shows the 3-D scores plot between selected components; Figure 11 shows the loading plot between the selected components; Figure 12 shows the classification performance with different number of components; Figure 13 shows the results of permutation test for model validation; Figure 14 shows important features identified by PLS-DA.

⁴Ron Wehrens and Bjorn-Helge Mevik. *pls: Partial Least Squares Regression (PLSR) and Principal Component Regression (PCR)*, 2007, R package version 2.1-0

⁵Max Kuhn. Contributions from Jed Wing and Steve Weston and Andre Williams. *caret: Classification and Regression Training*, 2008, R package version 3.45

⁶Bijlsma et al. *Large-Scale Human Metabolomics Studies: A Strategy for Data (Pre-) Processing and Validation*, Anal Chem. 2006, 78 567 - 574

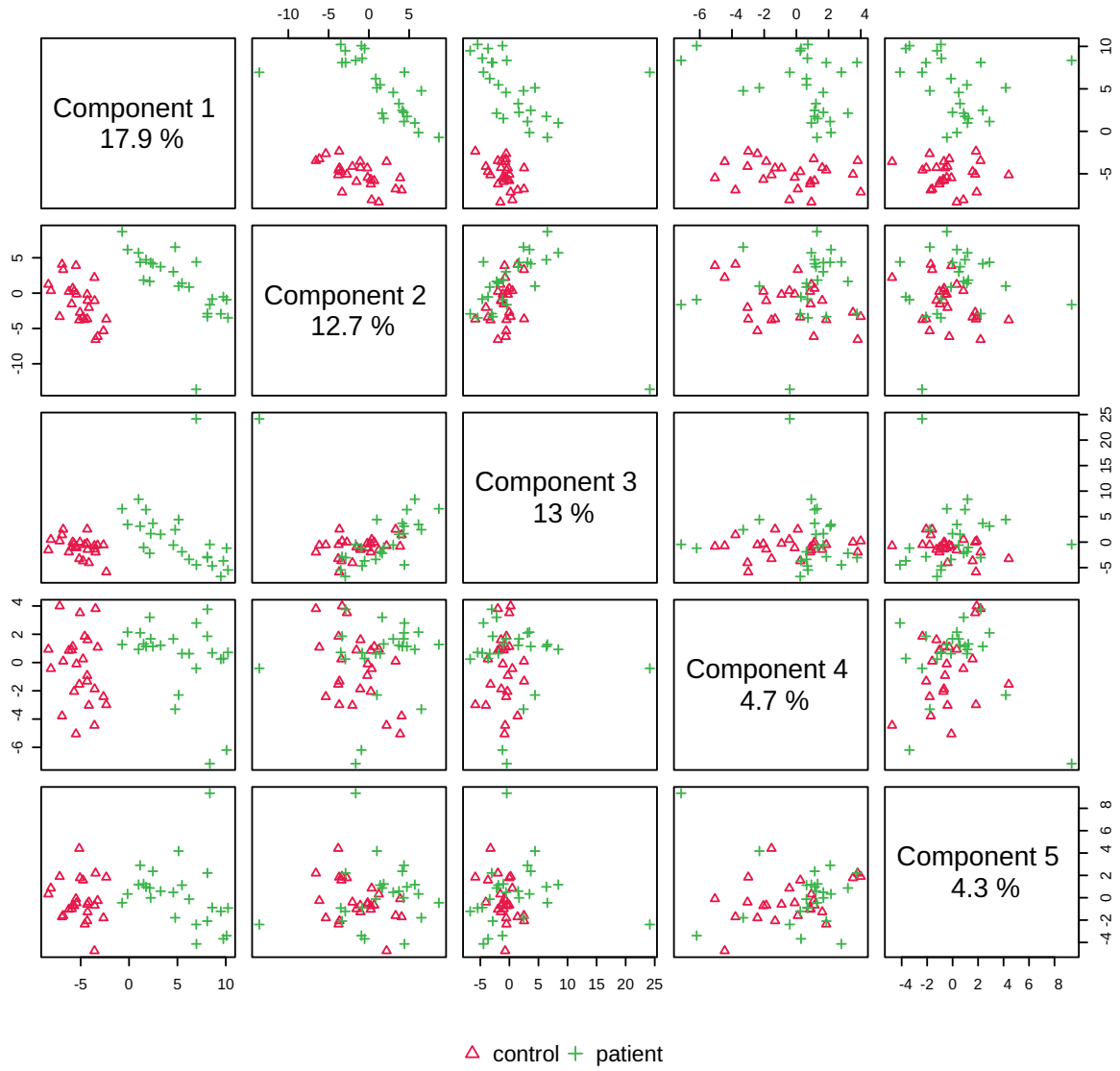
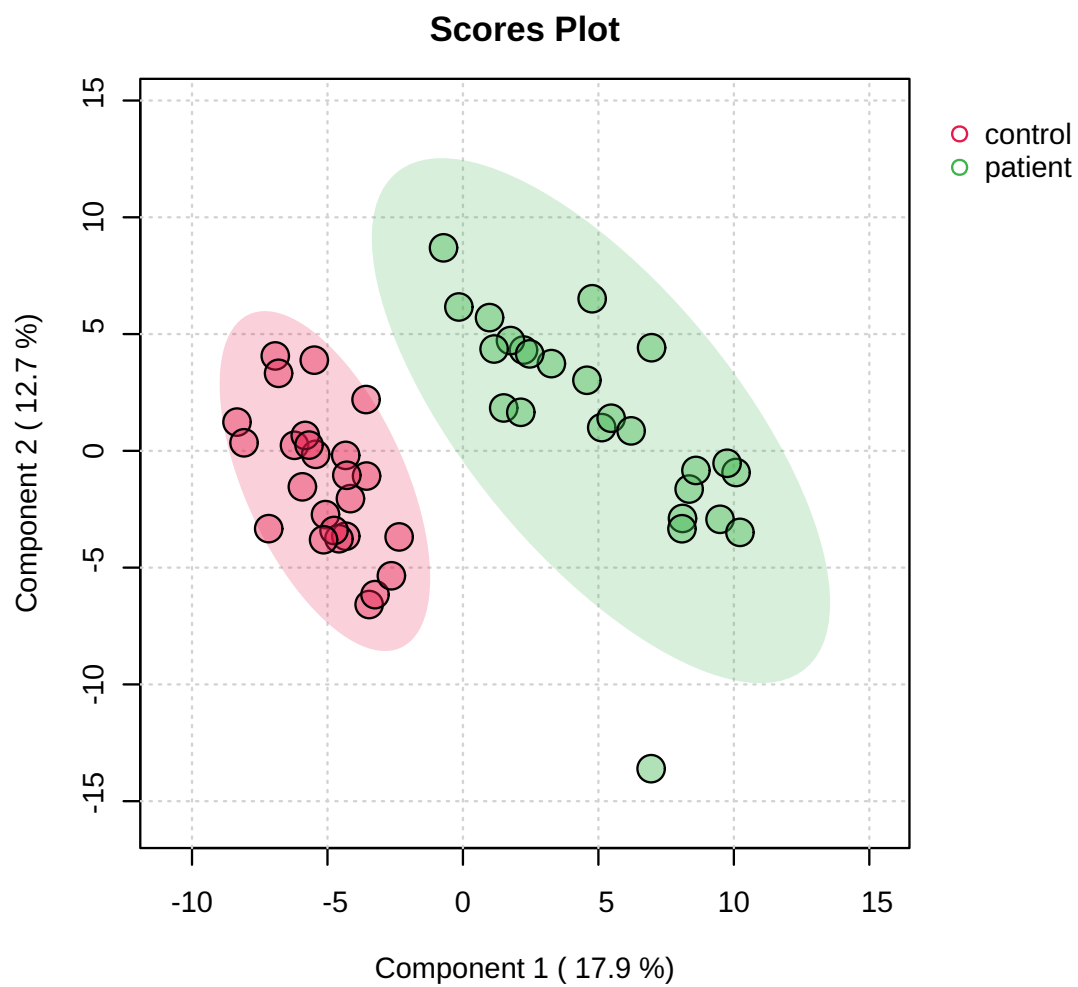


Figure 8: Pairwise scores plots between the selected components. The explained variance of each component is shown in the corresponding diagonal cell.



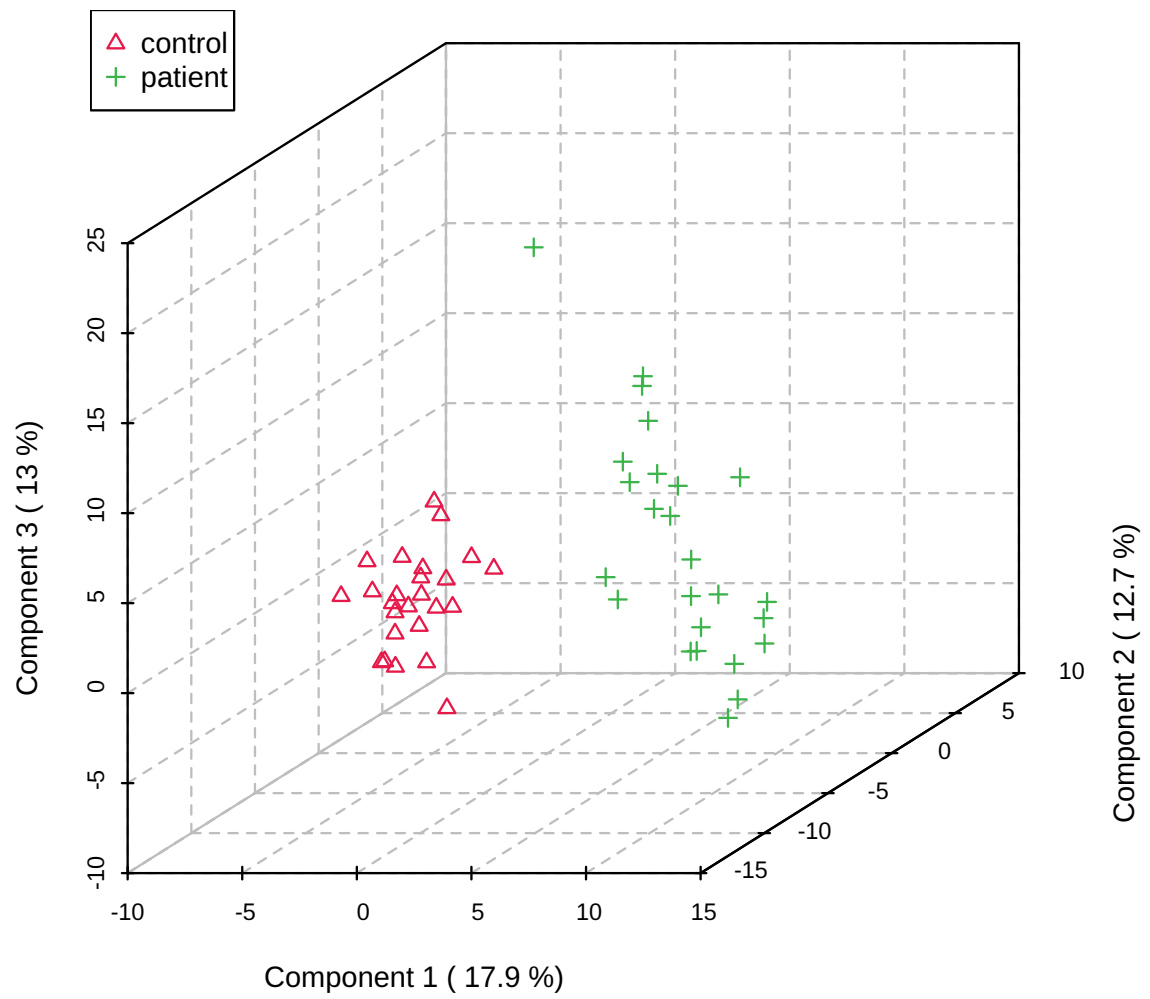
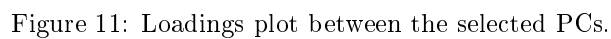


Figure 10: 3D scores plot between the selected PCs. The explained variances are shown in brackets.



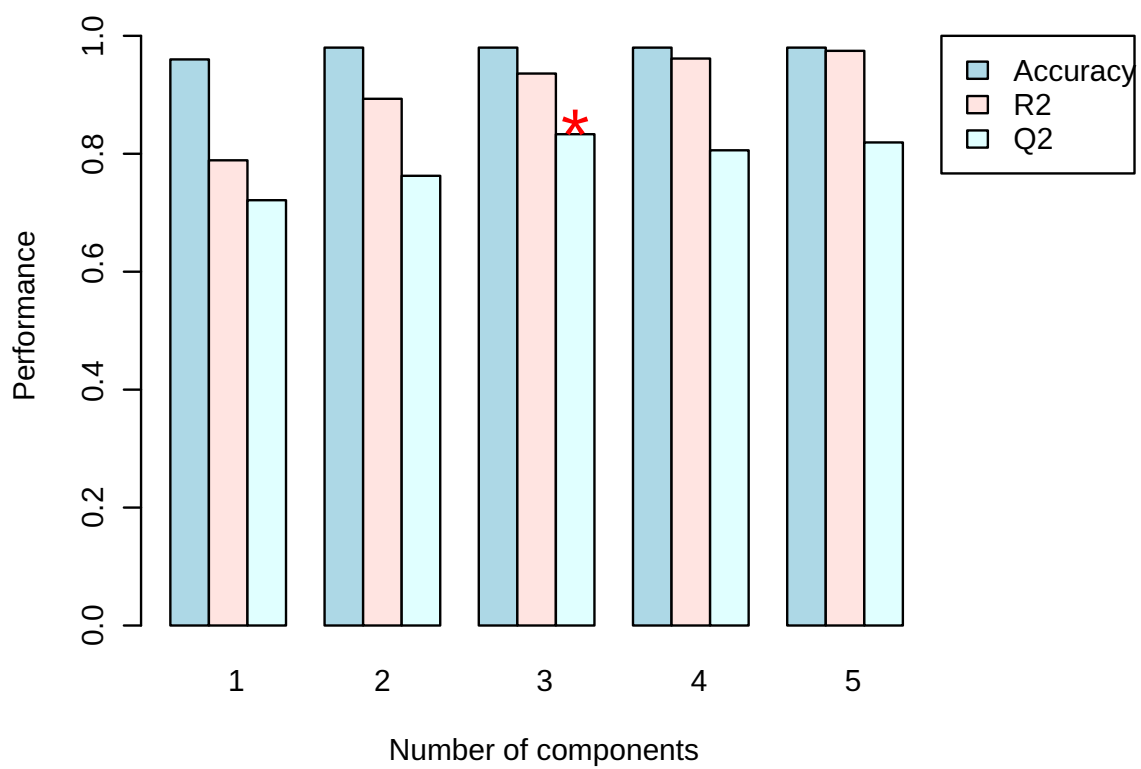


Figure 12: PLS-DA classification using different number of components. The red star indicates the best classifier.

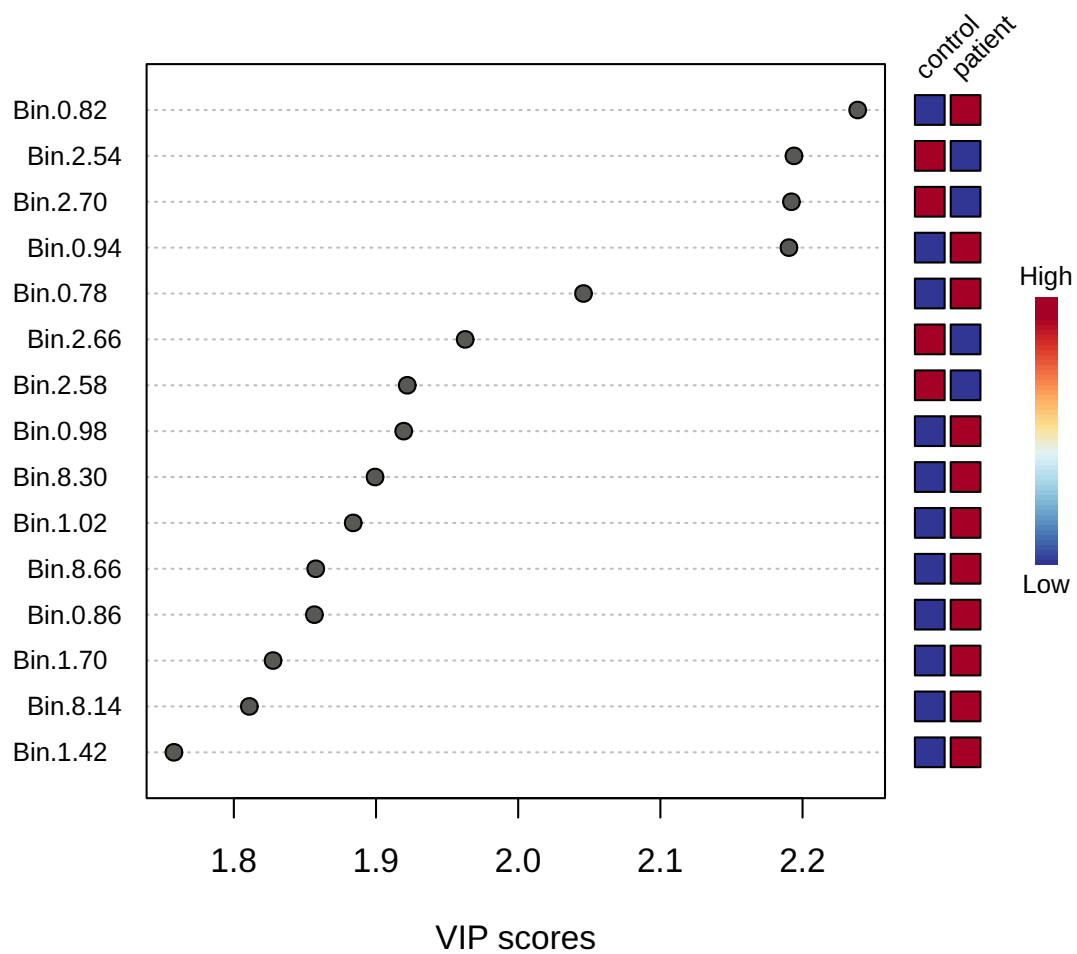


Figure 13: Important features identified by PLS-DA. The colored boxes on the right indicate the relative concentrations of the corresponding metabolite in each group under study.

2.5 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 15 shows the clustering result in the form of a dendrogram. Figure 16 shows the clustering result in the form of a heatmap.

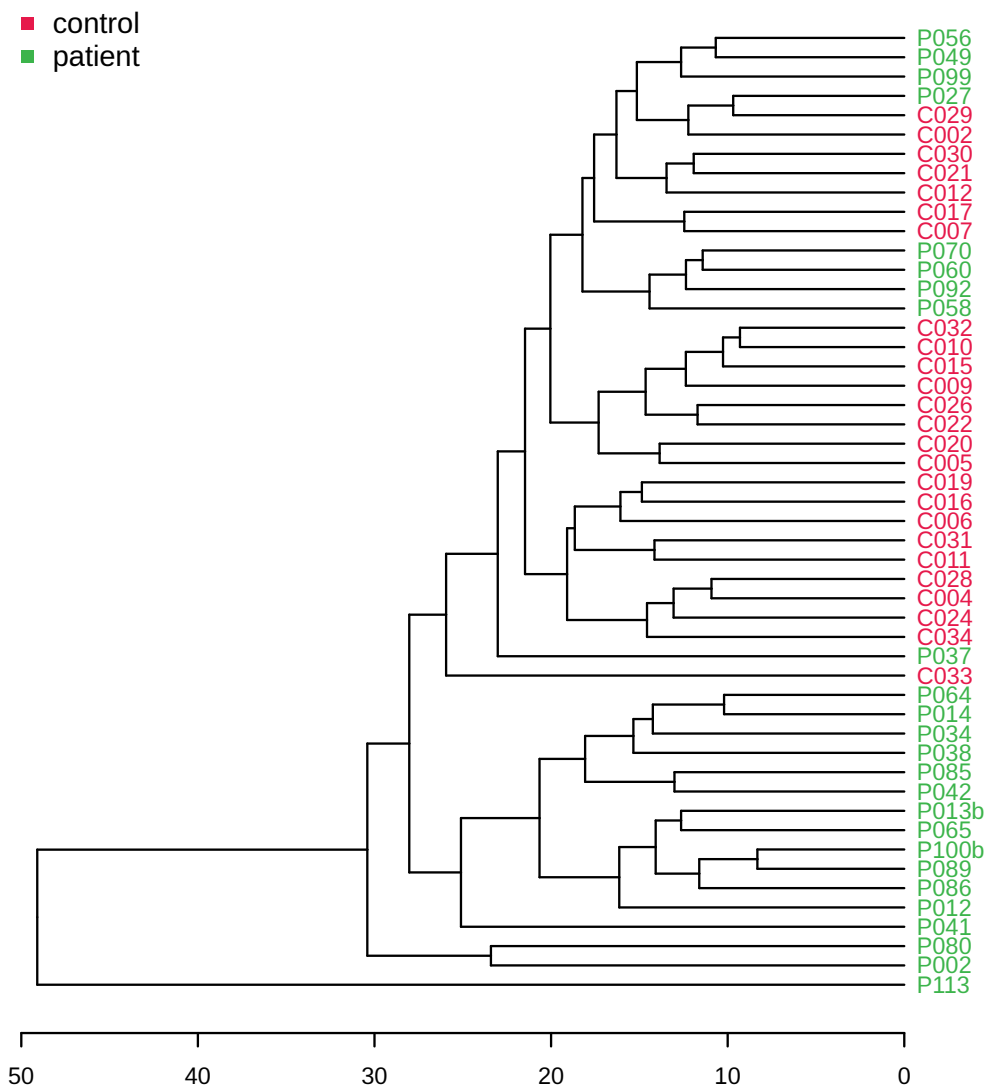


Figure 14: Clustering result shown as dendrogram (distance measure using `euclidean`, and clustering algorithm using `complete`).

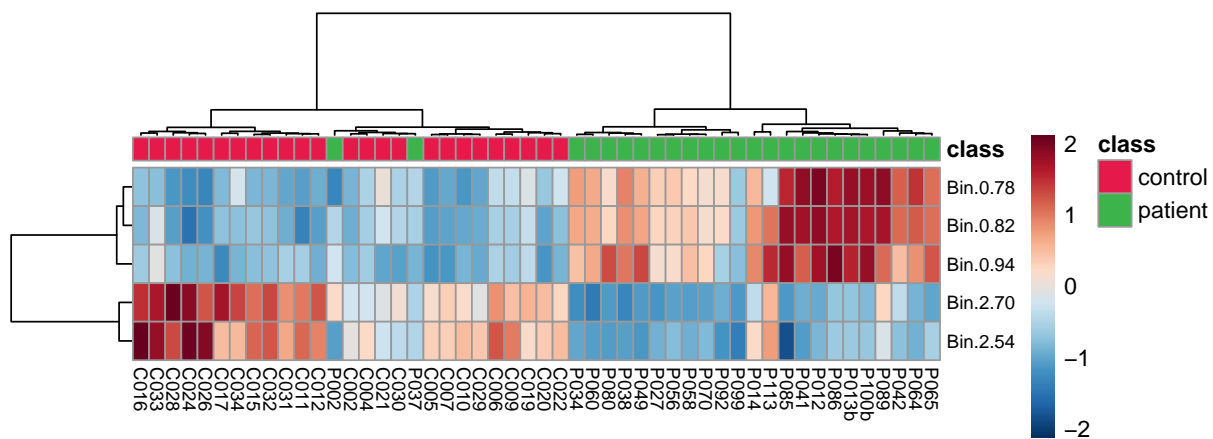


Figure 15: Clustering result shown as heatmap (distance measure using `euclidean`, and clustering algorithm using `ward.D`).

3 Appendix: R Command History

```
[1] "mSet<-InitDataObjects(\"specbin\", \"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<-PreparePrenormData(mSet)"
[7] "mSet<-Normalization(mSet, \"MedianNorm\", \"LogNorm\", \"AutoNorm\", ratio=FALSE, ratioNum=20)"
[8] "mSet<-PlotNormSummary(mSet, \"norm_0_\", \"png\", 72, width=NA)"
[9] "mSet<-PlotSampleNormSummary(mSet, \"snorm_0_\", \"png\", 72, width=NA)"
[10] "mSet<-PlotHCTree(mSet, \"tree_0_\", \"png\", 72, width=NA, \"euclidean\", \"ward.D\")"
[11] "mSet<-PlotHCTree(mSet, \"tree_1_\", \"png\", 72, width=NA, \"euclidean\", \"complete\")"
[12] "mSet<-PlotSubHeatMap(mSet, \"heatmap_1_\", \"png\", 72, width=NA, \"norm\", \"row\", \"euclidean\")"
[13] "mSet<-PCA.Anal(mSet)"
[14] "mSet<-PlotPCAPairSummary(mSet, \"pca_pair_0_\", \"png\", 72, width=NA, 5)"
[15] "mSet<-PlotPCAScree(mSet, \"pca_screes_0_\", \"png\", 72, width=NA, 5)"
[16] "mSet<-PlotPCA2DScore(mSet, \"pca_score2d_0_\", \"png\", 72, width=NA, 1,2,0.95,0,0, \"na\")"
[17] "mSet<-PlotPCALoading(mSet, \"pca_loading_0_\", \"png\", 72, width=NA, 1,2);"
[18] "mSet<-PlotPCABiplot(mSet, \"pca_biplot_0_\", \"png\", 72, width=NA, 1,2)"
[19] "mSet<-PlotPCA3DLoading(mSet, \"pca_loading3d_0_\", \"json\", 1,2,3)"
[20] "mSet<-PlotPCAPairSummary(mSet, \"pca_pair_1_\", \"png\", 72, width=NA, 7)"
[21] "mSet<-PlotPCAPairSummary(mSet, \"pca_pair_2_\", \"png\", 72, width=NA, 3)"
[22] "mSet<-PLSR.Anal(mSet, reg=TRUE)"
[23] "mSet<-PlotPLSPairSummary(mSet, \"pls_pair_0_\", \"png\", 72, width=NA, 5)"
[24] "mSet<-PlotPLS2DScore(mSet, \"pls_score2d_0_\", \"png\", 72, width=NA, 1,2,0.95,0,0, \"na\")"
[25] "mSet<-PlotPLS3DScoreImg(mSet, \"pls_score3d_0_\", \"png\", 72, width=NA, 1,2,3, 40)"
[26] "mSet<-PlotPLSLoading(mSet, \"pls_loading_0_\", \"png\", 72, width=NA, 1, 2);"
[27] "mSet<-PlotPLS3DLoading(mSet, \"pls_loading3d_0_\", \"json\", 1,2,3)"
[28] "mSet<-PlotPLS.Imp(mSet, \"pls_imp_0_\", \"png\", 72, width=NA, \"vip\", \"Comp. 1\", 15,FALSE)"
[29] "mSet<-PLSDA.CV(mSet, \"5\", 5,5, \"Q2\")"
[30] "mSet<-PlotPLS.Classification(mSet, \"pls_cv_1_\", \"png\", 72, width=NA)"
[31] "mSet<-Ttests.Anal(mSet, F, 0.05, FALSE, TRUE, \"fdr\", FALSE)"
[32] "mSet<-PlotTT(mSet, \"tt_0_\", \"png\", 72, width=NA)"
[33] "mSet<-PlotCorrHeatMap(mSet, \"corr_1_\", \"png\", 72, width=NA, \"col\", \"pearson\", \"bwm\", \"bwm\")"
[34] "mSet<-PlotCorrHeatMap(mSet, \"corr_2_\", \"png\", 72, width=NA, \"col\", \"pearson\", \"bwm\", \"bwm\")"
[35] "mSet<-SaveTransformedData(mSet)"
[36] "mSet<-PreparePDFReport(mSet, \"guest5646536859053858679\")\n"
```

The report was generated on Wed Mar 20 17:19:01 2024 with R version 4.3.2 (2023-10-31), OS system: Linux, version: -Ubuntu SMP Tue Jan 9 15:25:40 UTC 2024 .