

Meat comparison report

Lisa Schneider

April 10, 2020

Installation and namespaces

The files with the self-written functions for this workflow are loaded with `source()`. The first source-file “install_packages_lipidome_comparison.R” contains a function that installs all required packages. The installed packages are loaded into the namespace with `library()`.

```
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
```

Set directorys

This is set for my system. The syntax for windows is slightly different. Making this work on all systems, maybe from command line is planned for later.

```
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
```

A lot of data preprocessing

The meat data was used without preprocessing, except exproting the excel sheet to csv. The data was read from csv to the script and stripped of meta data. The string “N/F” for not identified lipids was replaced with NA and the data was separated in target data and internal standard data. The target data was then re-formatted, with the compounds as the columns and the sample IDs as the rows. The used function flip_df() is from the “lipidome_comparison_dataTransformaions.R” file. Metadata was extracted from the sample IDs and the data set was separated by extraction method (AS and N).

```
## knitr
## dplyr
## stringr
## data.table
```

```
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
```

Data imputation

It can be assumed that most of the missing values are due to the concentration being below the detection limit. Therefore the missing values are considered non-random and left-censored. Imputation only works up to a certain percentage of missing values, without changing the results. I found 20% missing values in the literature (Gurke, 2019, doi: 10.3389/fpsy.2019.00041), therefore all compounds with more than 20% missing values were filtered. The remaining missing values were imputed using the QRILC (Quantile Regression Imputation of Left-Censored data) procedure. This procedure has the disadvantage of producing negative values (because it performs random draws from a distribution), which is a problem for the calculation of the log2 foldchange (I'm working on that).

```
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
```

```
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
```

Reducing the replicates

To not artificially produce more samples than we have and to reduce variation in preprocessing and measurement the means for each sample or biological replicate are calculated. The idea is, to continue working with these values, but in most cases this reduces the sample number below a point where most procedures work ($n[\text{fish}] = 1$). So in the for now I continued working with the non reduced data.

The self-written functions in this part are in “lipidome_comparison_dataTransformations.R” and “lipidome_comparison_EDA.R”.

```
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
```

```
## heatmaply
## gplots
## plotly
## dendextend
```

Exploratory data analysis

Exploratory data analysis was performed graphically and using the shapiro-wilk test. The graphical data analysis (using qqplot, histogram and boxplot) turned out not to be very practical due to the number of compounds, but it gives a good overview over the distributions of the variables. It showed, that some of the variables had normal distribution, while others were multi modal. Therefore the results of the following tests for normal distribution and correlation have to be handled carefully. All the graphs are in a separate zip-folder due to their size and runtime.

```
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
```

Principal component analysis

The selection of the number of principal components depends on the explained variance displayed by `get_eigenvalue()` and the scree plots. There are two different representations of the scree plots, because I find both of them useful. Unfortunately in our data each PC only explains a small percentage of variance, PC1 and PC2 together explain only about 36%, therefore this model is not ideal. The problem with the separation is also visible in the biplots. While the fish-values are close together (due to them being replicates of the same sample) and well separated from the other groups, there is almost no separation between meat and game. The last two plots show which compounds have the highest contribution to the principal components. The matrix gives an overview over all five principal components and all variables. Due to the number of components, the elements of the plot are very small. A better representation are the barcharts for the top ten contributions to a selected number of principal components (I plotted PC1 and PC2). With this information the number of variables can be reduced, for example for a parallel plot, a spider plot or for hypothesis generation and testing.

```
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
## knitr
## dplyr
## stringr
## data.table
## tibble
```

```
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
```

```

## gplots
## plotly
## dendextend
## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend

```

Hierarchical clustering

The functions `hclust_performance_table` and `hclust_performance_plot` both use the `dend_expend` function to find the best performing distance and hierarchical clustering methods. The barplot displays the values from the “optim” column graphically. In hierarchical clustering, game is put in one cluster and meat and fish are clustered together. In the heatmap it is visible, that there are two different groups of game. In the zip-folder with the EDA-plots there is also a html-file for the interactive heatmap, because printing it into pdf did not work.

```

## knitr
## dplyr
## stringr
## data.table
## tibble

```



```

## imputeLCMD
## impute
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend

```

Hypothesis testing and volcano plot

To find variables with a significant difference in at least one group, Kruskal-Wallis-test (because not all variables were normally distributed) was performed column-wise. The resulting table can be used to select the significant variables of the original data frame for further analysis.

The volcano plots were performed for two groups each. The significance was assessed using the Wilcoxon rank sum test (again because of the missing normal distribution in some variables). In addition a multiple testing correction using FDR was performed. Log2-foldchange was calculated by performing log2 for all variables in both groups and then subtracting the test-group from the control-group values. There is a problem with the negative values imputed by impute.QRILC (I am working on finding out how to constrain the imputation method so it does not impute negative values): The log2-function produces NaNs from negative values. Until the imputation problem is solved, the variables with NaNs are excluded.

The volcano_plot function prints a dotplot with the log2-foldchange on the x-axis and the negative log10 of the p-value on the y-axis. The significant up- or down regulated lipids are colored and labelled. Because the number of variables is relatively low for a volcano plot (85 variables), the shape does not resemble a volcano, this should be improved with an increasing number of identified lipids.

```

## knitr
## dplyr
## stringr
## data.table
## tibble
## imputeLCMD
## impute

```

```
## ggplot2
## viridis
## GGally
## fmsb
## scales
## ggrepel
## ggpubr
## FactoMineR
## factoextra
## ggfortify
## corrplot
## heatmaply
## gplots
## plotly
## dendextend
```

Outlook

- Fix imputation of negative values
- Play around with normalization methods
- Supervised learning with random forest:
 - Makes sense for establishing if and how the kind of meat can be predicted, because we already have training data.
 - Variable selection is already implemented in the randomForest r-ackage.
 - The only problem is that there is only one small data set, which might be too small to produce training and test data.
- Add option to easily run the workflow (either as a shiny app or commandline interface)
- Prepare installation script for dependencies that works on windows as well