

# User guide:

Assume a windows machine is being used with no prior software installed other than Microsoft Excel 2016 and a web browser.

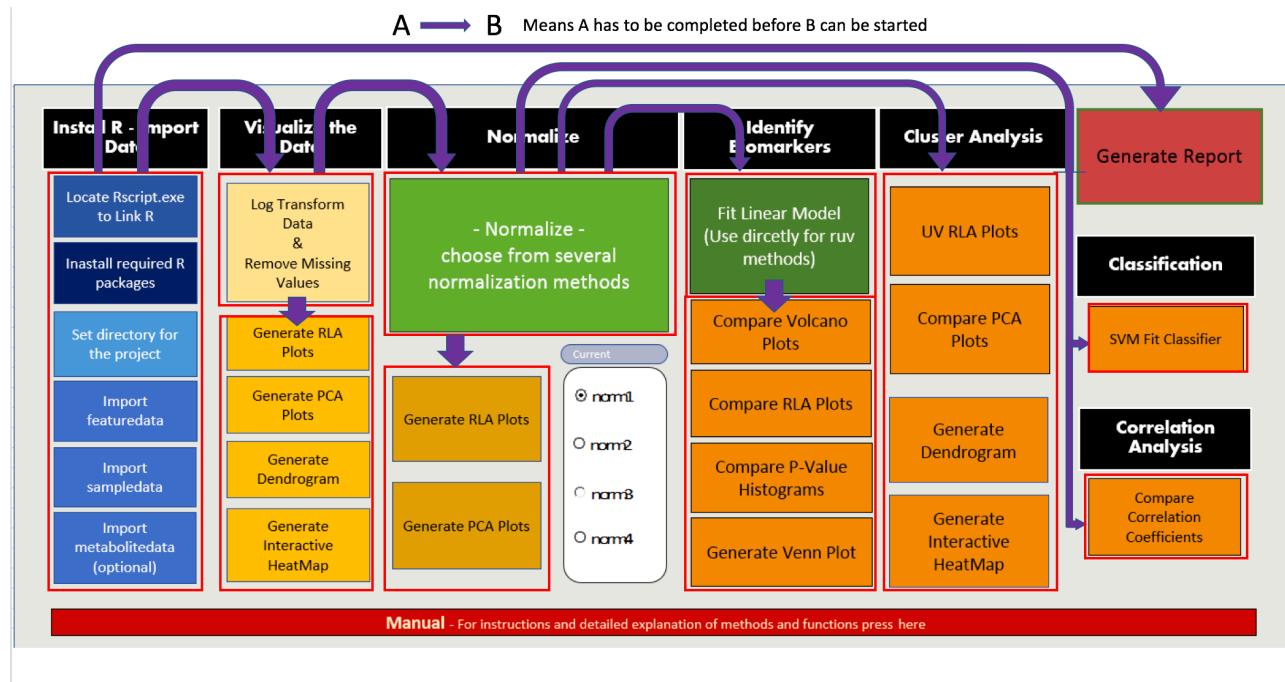
## Getting the files ready:

To get the required files, go to <https://github.com/metabolomicstats/NormalizeMets> and download the latest version of the *ExNormalizeMetsSetup.zip*. After downloading unzip the folder, it should contain the following files:

----- Insert Folder Image once final name is decided upon -----

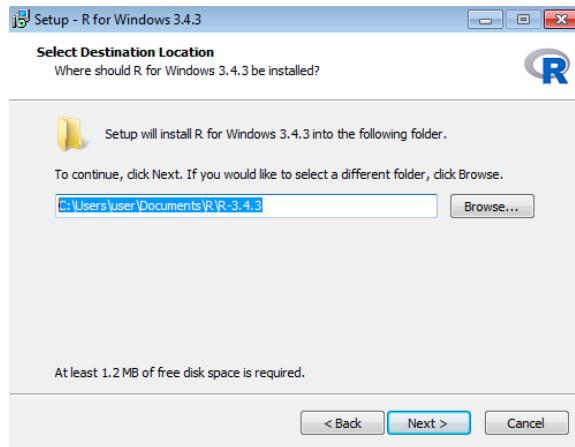
They are all the files needed for installing R and loading NormalizeMets through excel.

## **General Workflow:**



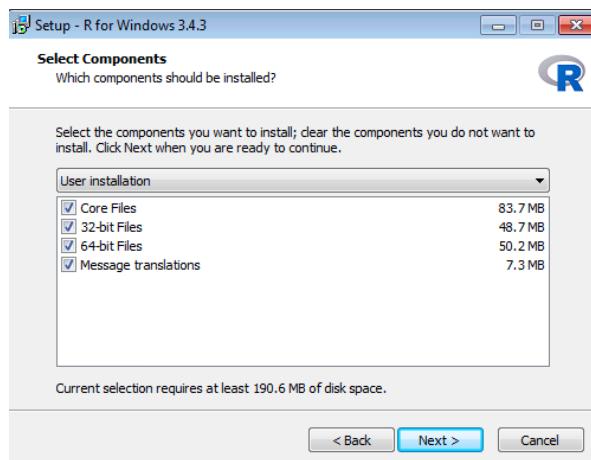
## Installing R:

1. Install R from the extracted ExNormalizeMetsSetup folder by running the *R-3.4.3-win* file.
2. Follow the installation instructions and choose the location where R should be installed.

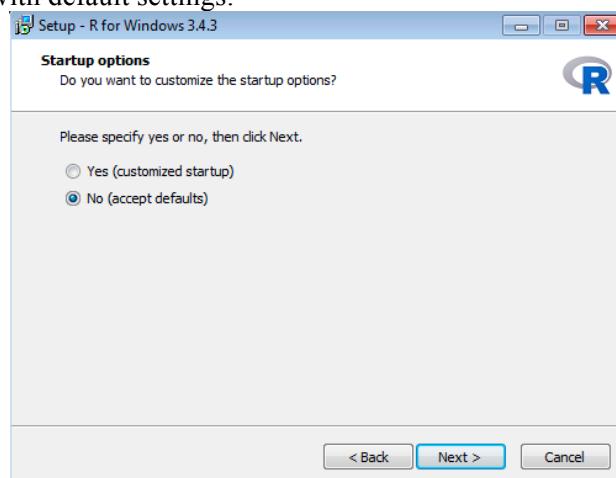


Click *Next* for default location (Recommended) or choose location manually. Make sure you **remember where R is installed** as you will need to locate this folder later to link R to Excel.

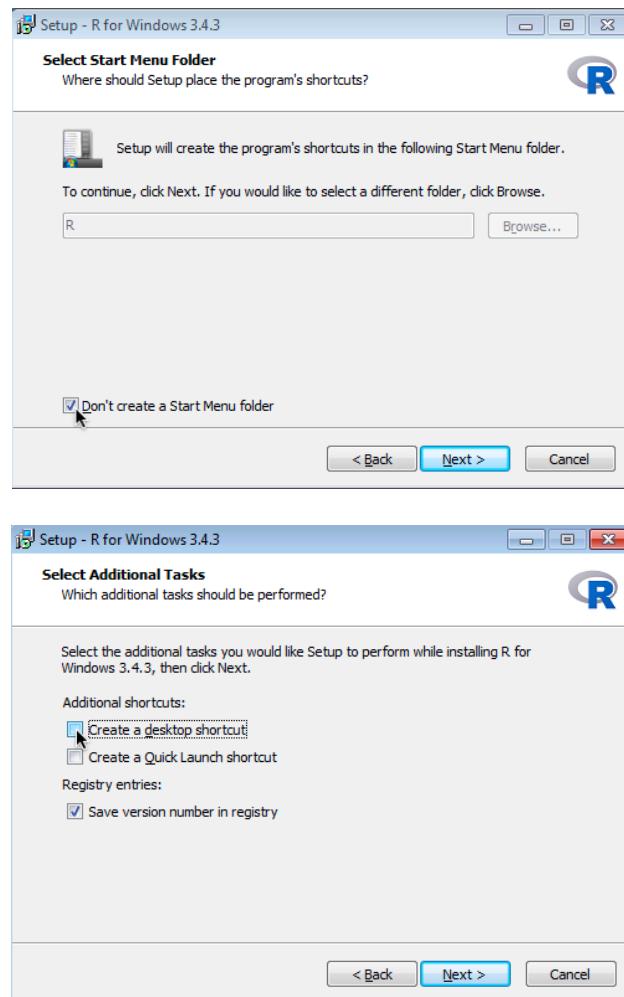
3. Click *Next* to install with the different required settings:



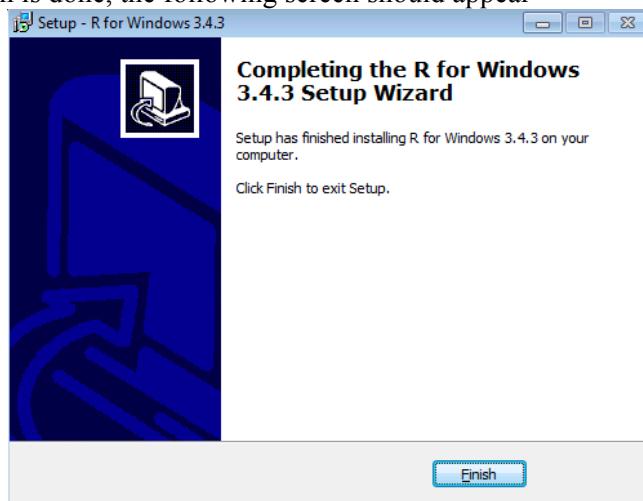
4. Click *Next* to install with default settings:



5. In the next screen check the *Don't create a Start Menu folder* icon if you don't intend to use R by itself and click next. Also uncheck the *Create a desktop shortcut* in the next screen:

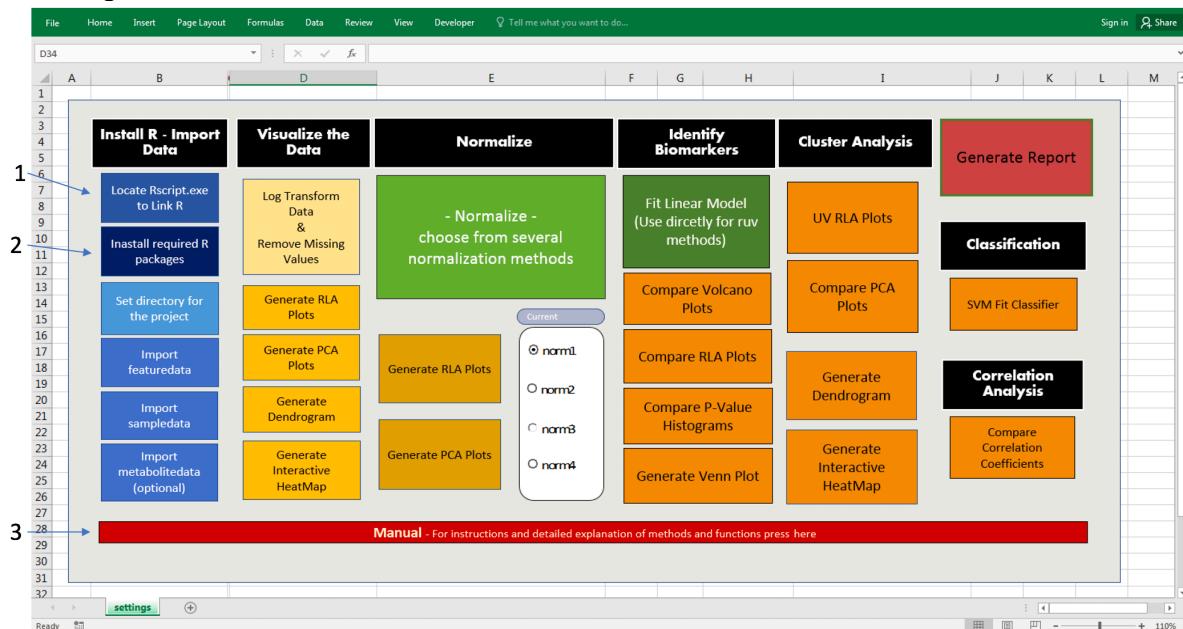


6. Wait until installation is done, the following screen should appear



**Linking R to excel and installing the required packages (first use only):**

After installing R open the excel file ExNormalizeMets.xlsx, this will open the Excel interface onto the settings sheet:

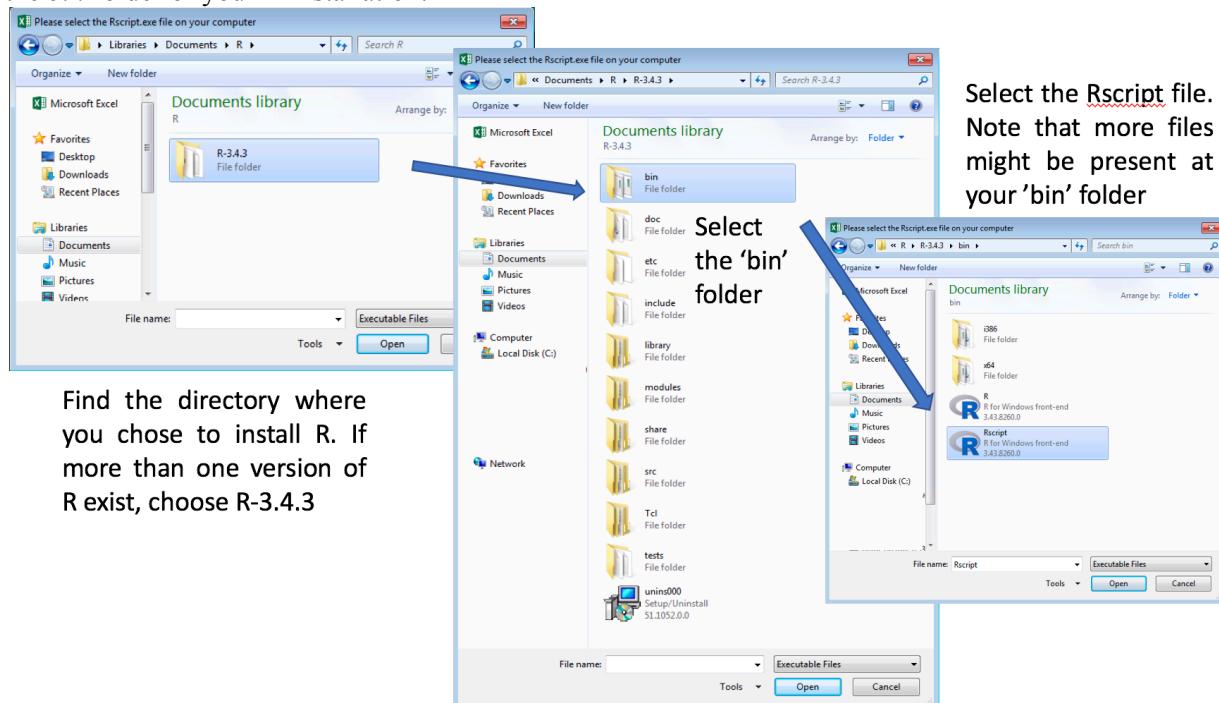


This is your ‘Control’, any function you want to run, from importing data, Normalizing, viewing results and opening the manual can be done from here.

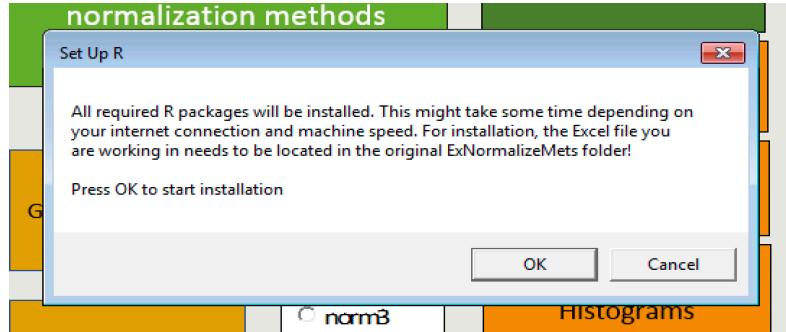
This manual can be accessed at any stage by clicking on (3) but more on this later. For now, first locate the Rscript.exe file so that excel will know how to run R commands it generates.

### Locating Rscript:

Press (1) to locate the Rscript file in the window that opens up. Make sure to select the *Rscript* file in the *bin* folder of your R installation:

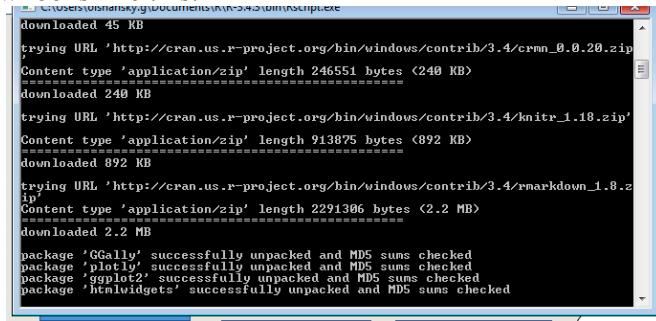


After selecting the file, to install all the required R packages and set up the needed dependencies, press (2).



Pressing ok will start the installation, this might take a few minutes if you are using R for the first time as many of the base packages will need to be installed.

The installation window looks like this:



When the installation is done, a window with the message *Done!* will appear.

Now that the installation is complete, NormalizeMets is ready for use!

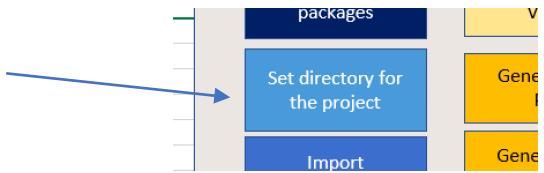
## Using NormalizeMets

Data used in the following examples is provided with NormalizeMets (alldata\_eg in R), it is located in the *ExampleData* folder in your downloaded *ExNormalizeMetsSetup* file. Future references in this guide refer to this data by default.

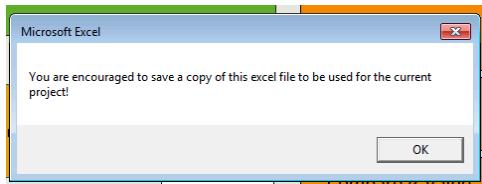
An example excel document containing the data used for the tutorial (MyFirstNormalizeMetsProject\_example.xlsx), with all settings identical to those in the tutorial is provided.

### Starting a new project:

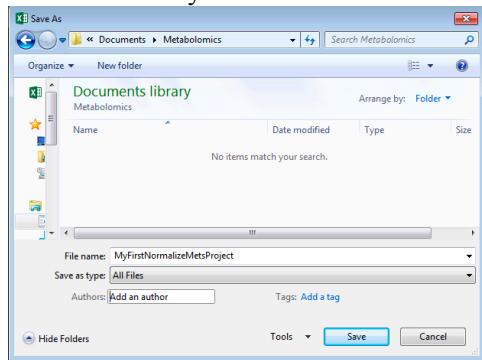
To start a new project, you will need to set up or use an existing directory where the project will ‘live’, data and plots generated by excel will all be saved to that folder.



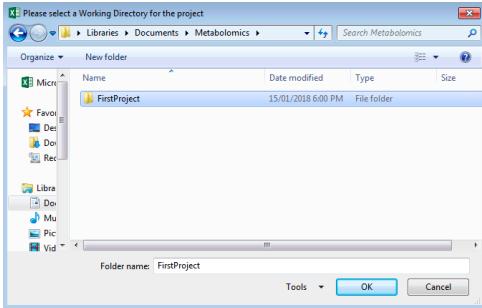
Before selecting the folder, you will be prompted to save a copy of the current version of the excel file, it is recommended to save it with a new name to make sure a ‘clean’ version always stays in your ExNormalizeMetsSetup folder.



After saving the workbook under the name of your choice:

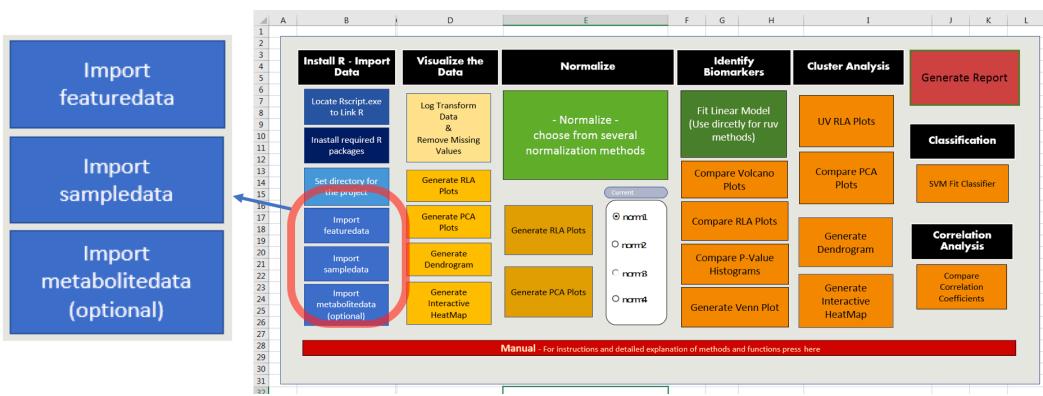


Select the working directory where all files should be generated. We recommend making a new folder for each project.



## **Loading data:**

To load data for the project, in turn click on the following to load the relevant data:



Loaded data needs to be in .csv format. After clicking on the required file, it will open a new sheet, showing the imported data. Select the setting sheet import more data and get back to the options.

For *featureddata*, set metabolites in columns and samples in rows. Unique sample names should be provided as row names.

---

*sampedata* should have sample information matching *featuredata* (samples in rows).

Optional *metabolitedata* should have metabolite information matching featuredata with metabolite names in rows.

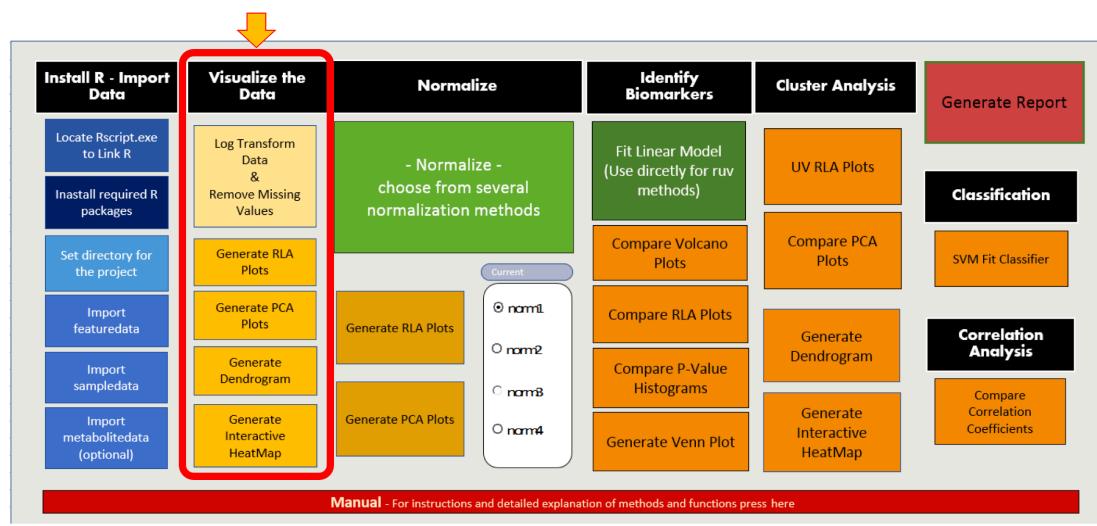
*metabolitedata* can include any metabolite information such as grouping structures, internal standard metabolites, negative control and positive control metabolites.

After the data is loaded, you are ready to proceed to analyse the data!

## NormalizeMets Workflow:

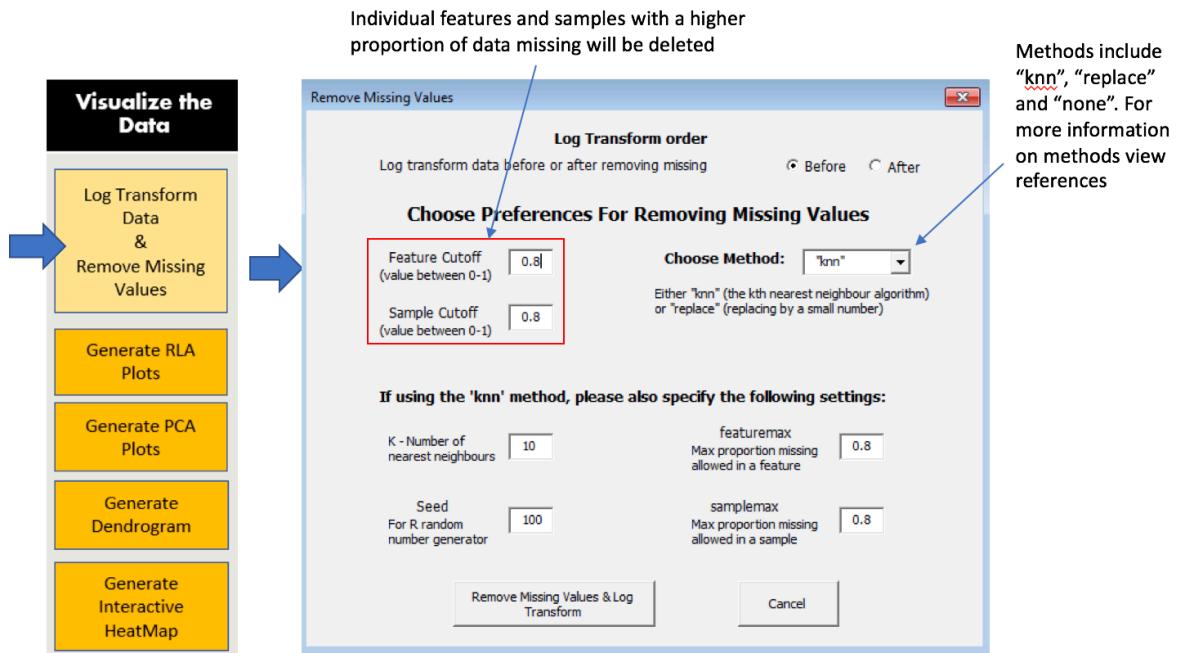
### Visualize the Data

The following section refers to the visualize part:



### **Log Transforming the data and removing missing value (mandatory):**

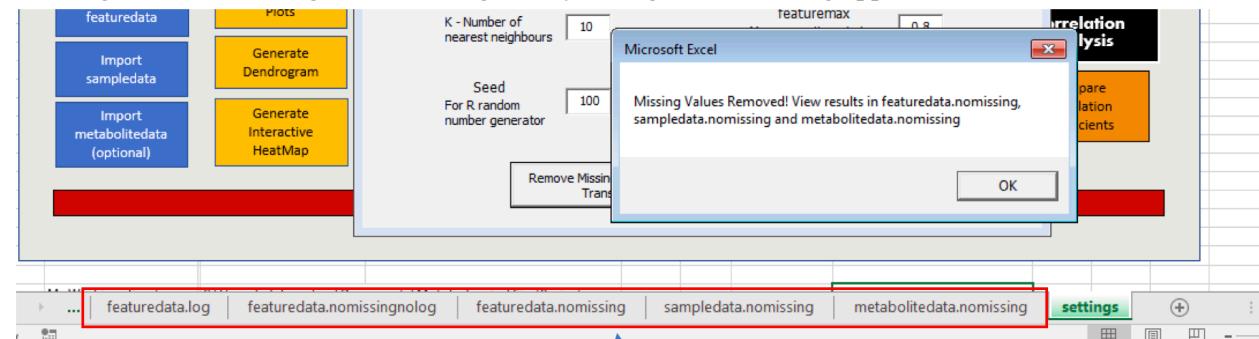
A frequent issue in metabolomics data sets is the occurrence of missing values. It is important to reduce the number of missing values as much as possible by using an effective pre-processing procedure. For example, a secondary peak picking method can be used for LC-MS data to fill in missing peaks which are not detected and aligned.



“knn” – use k nearest neighbours method to replace missing values.

“replace” – replaces missing values by half the minimum value in featuredata.

Clicking ‘Remove Missing Values & Log Transforming’ the following appears:



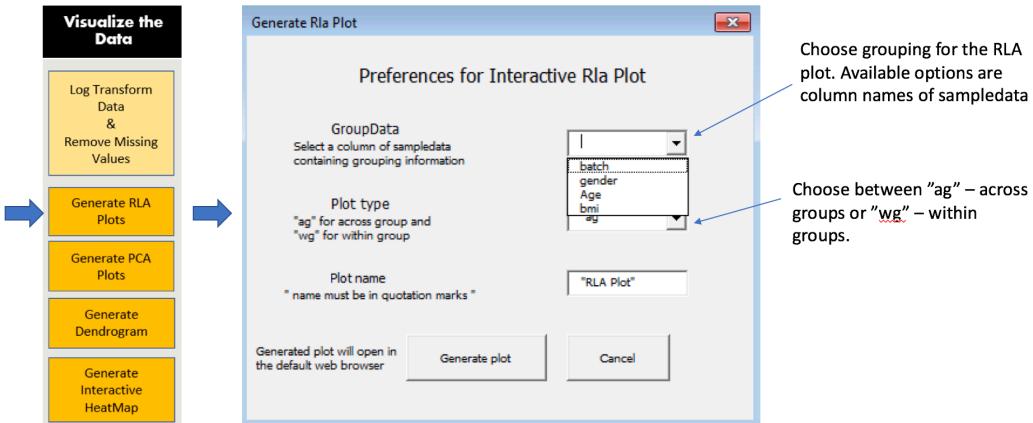
Note the new sheets that appeared, they have respectively the regular log transformed featuredata, the log transformed featuredata with missing values removed, featuredata with missing values removed and without log transformation, sampledata with rows removed corresponding to featuredata.nomissing, metabolitedata with rows removed corresponding to featuredata.nomissing.

Unless you are interested to view or copy any of this data, those sheets are only going to be used for further internal functions.

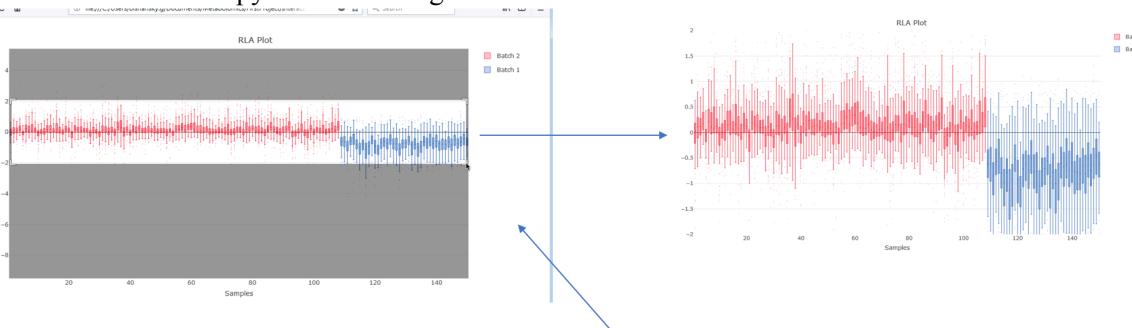
Now the plots in *Visualize the Data* can be generated! The data generated is also going to be used for the Normalization section.

## RLA plots

One way of visualising the log transformed metabolomics data is the use of *across group* or *within group* relative log abundance (RLA) plots (De Livera et al. 2012 De Livera et al. (2015)).

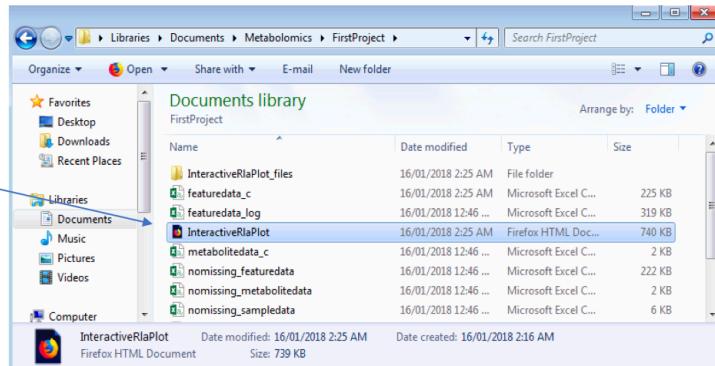


Setting groupdata to *batch* and selecting *Generate plot* opens the interactive plot in the default web browser and saves a copy in the working folder.



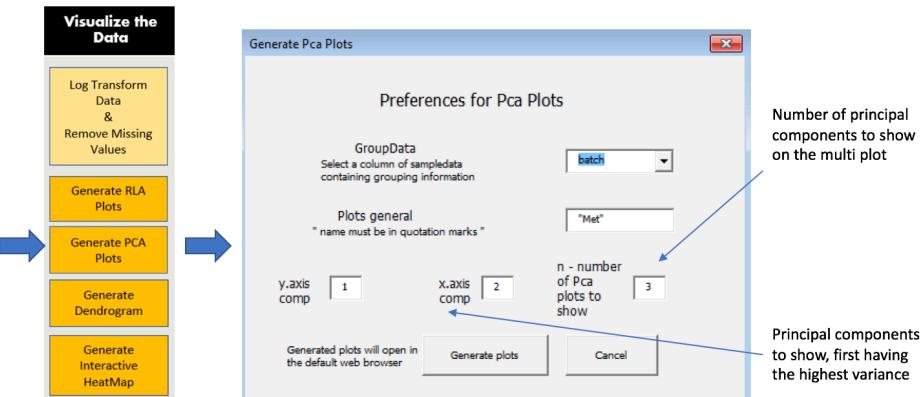
Interactive plot opens in the default browser, to zoom in, simple select the required part

Copy of all plots together with all other generated files are stored in the working directory



## PCA plots

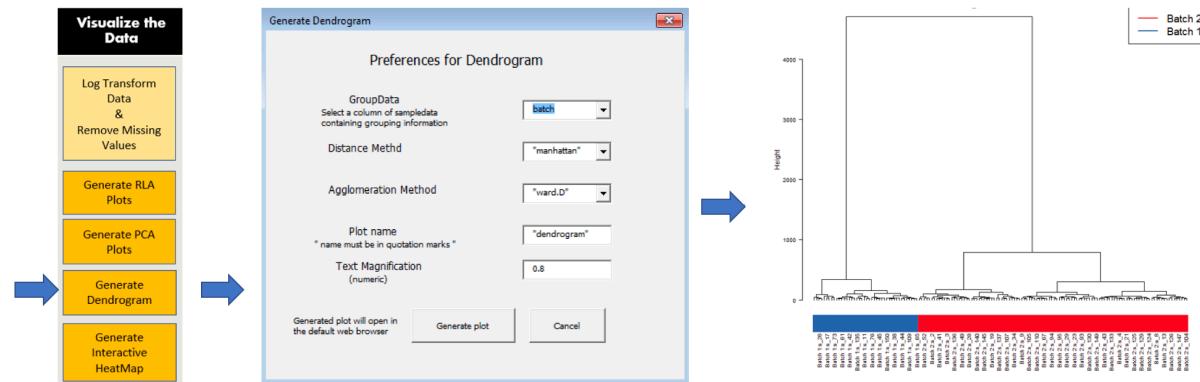
The following function can be used to obtain multiple plots for exploration of the principal components of the *featuredata* matrix: a bar plot indicating the variance explained by each principal component, scores and loading plots with specified axes (interactive and non-interactive), and a pairs plot of the first *n* principal components. These plots are useful in identifying any outlying samples and getting a preliminary understanding of the structure of the data.



All produced plots are stored in the working directory, with interactive plots opened in the browser and static plots located in the new *plots* sheet.

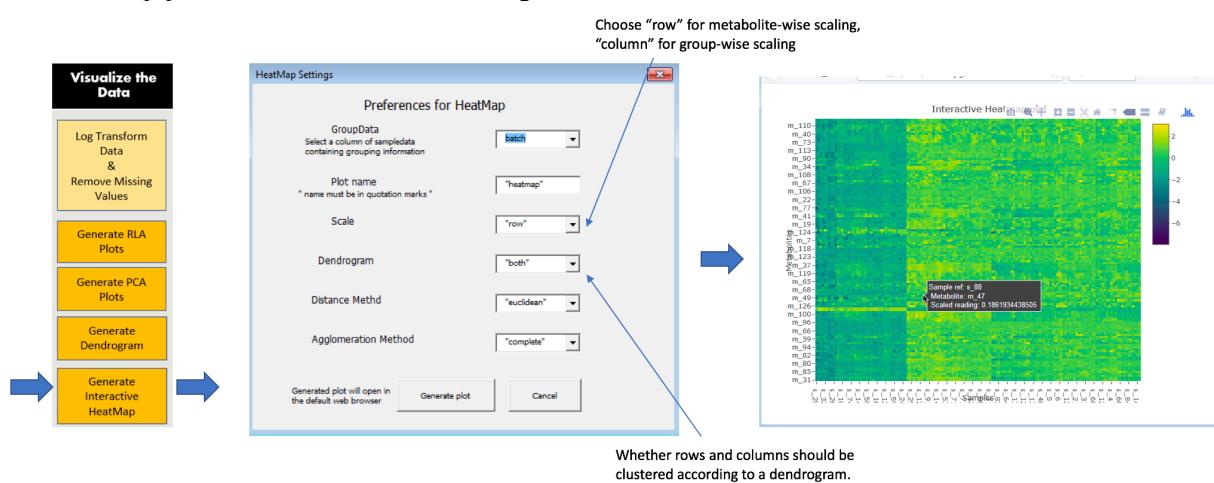
## Dendrogram

Generates a dendrogram to visualise clustering structures in the data, many different methods are available.



## HeatMap

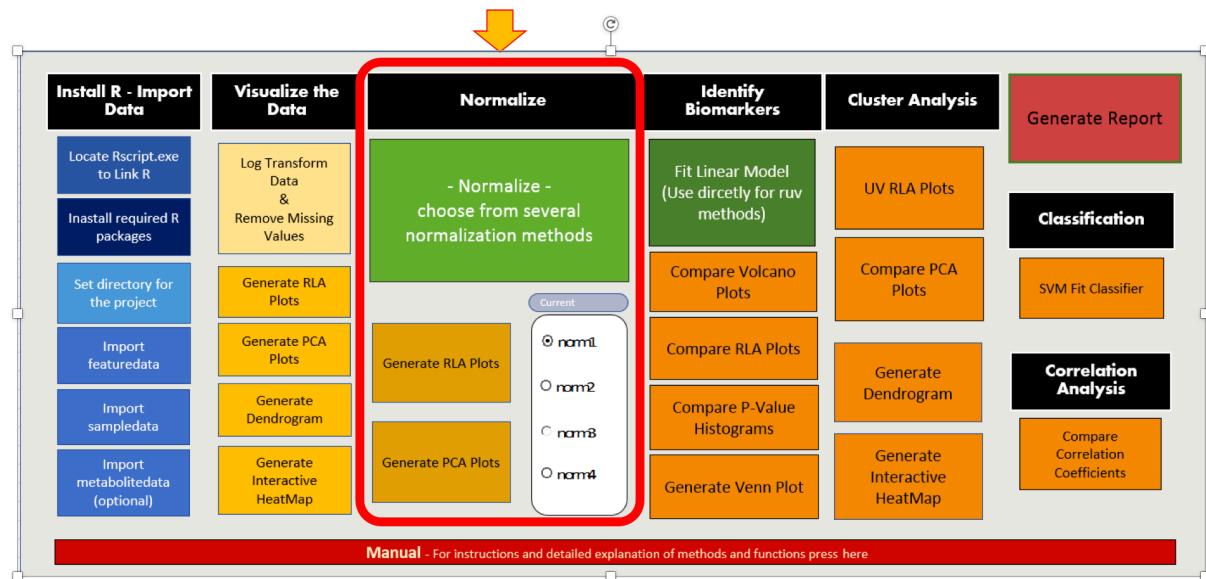
The HeatMap produced can reveal interesting structures in the data.



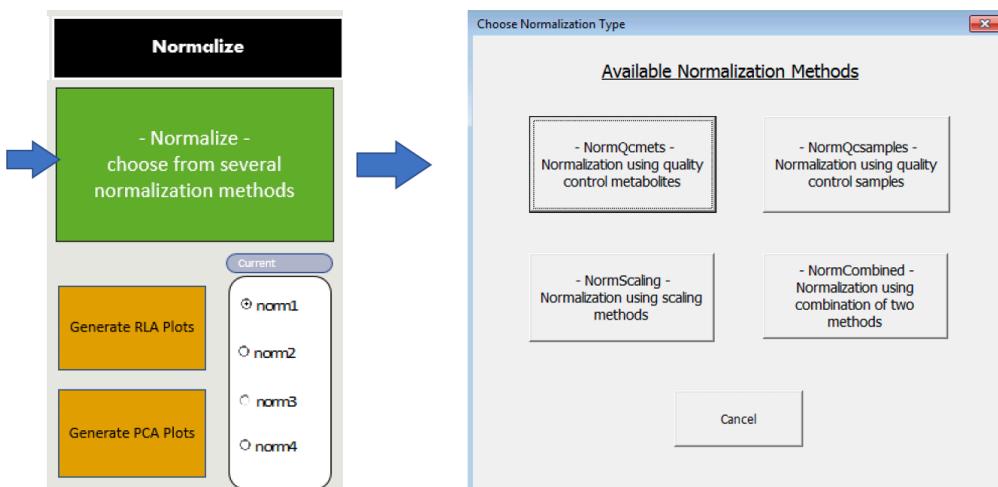
## Normalization

Normalization methods presented in this package are divided into four categories, as those which use (i) internal, external standards and other quality control metabolites (*NormQcmets*) (Systi-Aho et al. 2007, Redestig et al. (2009), De Livera et al. (2012), De Livera et al. (2015), Gullberg et al. (2004)) (ii) quality control samples (*NormQcsamples*) (Dunn et al. 2011), (iii) scaling methods (*NormScaling*) (Scholz et al. 2004, Wang et al. (2003)), and (iv) combined methods (*NormCombined*) (Kirwan and Broadhurst (2013)).

The normalization methods are accessible in the following section:



Clicking on the Normalize button opens the following menu enabling the choice of different normalization methods.

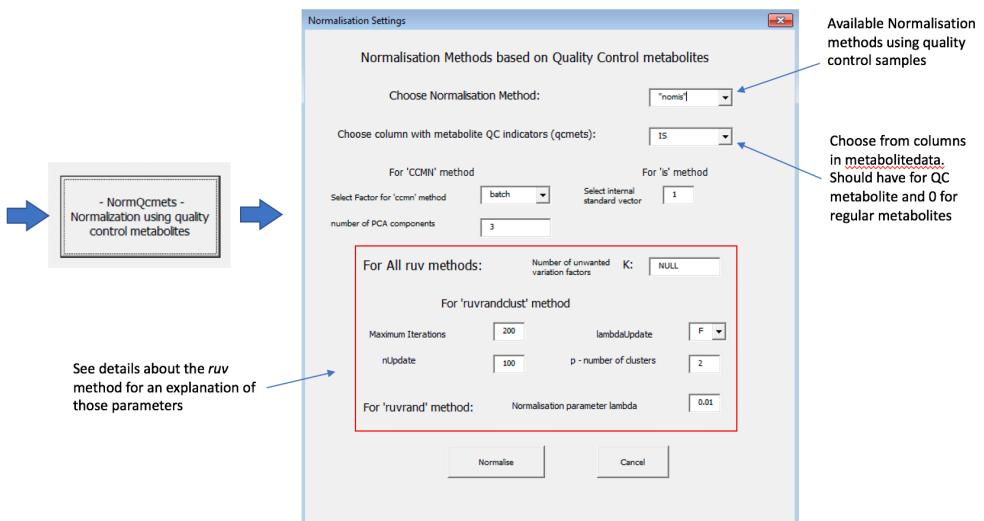


## NormQcmets

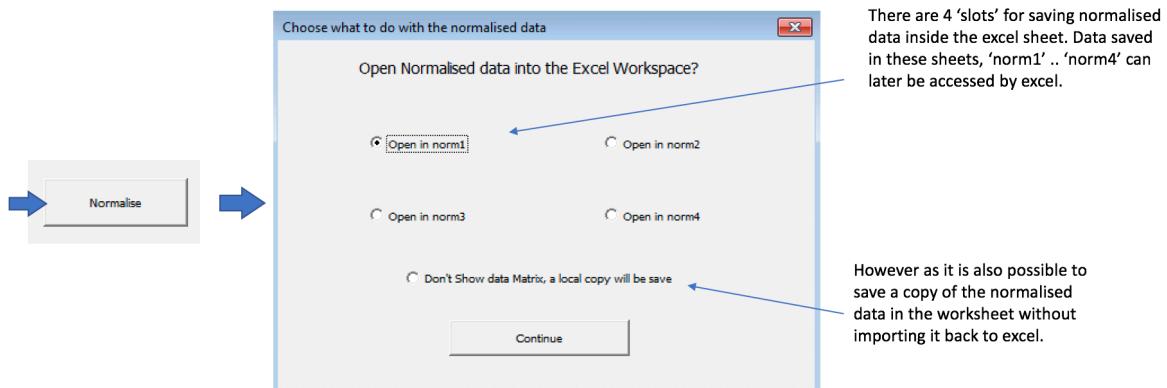
The approaches in *NormQcmets* use internal, external standards and other quality control metabolites. These include the *is* method which uses a single standard (Gullberg et al. 2004), the *ccmn* (cross contribution compensating multiple internal standard) method (Redestig et al. 2009), the *nomis* (normalization using optimal selection of multiple internal standards) method (Systi-Aho et al. 2007), and the remove unwanted variation methods (J. A. Gagnon-Bartsch, Jacob, and Speed 2014) as applied to metabolomics using “ruv2” (De Livera et al. 2012), “ruvrnd” and

“rvrrandclust” (De Livera et al. 2015). Note that *rvv2* is an application specific method designed for identifying biomarkers using a linear model that adjusts for the unwanted variation component.

To Normalize:



After Clicking the Normalise button the screen asking you where the normalized data is to be saved appears.



Upon clicking continue, you will return back to the settings sheet but you can notice some changes:

Shows current normalized data in the excel workbook

Normalisation location is not important, the names norm1, norm2, norm3 and norm4 are solely for convenience to keep track of the location of the normalised data

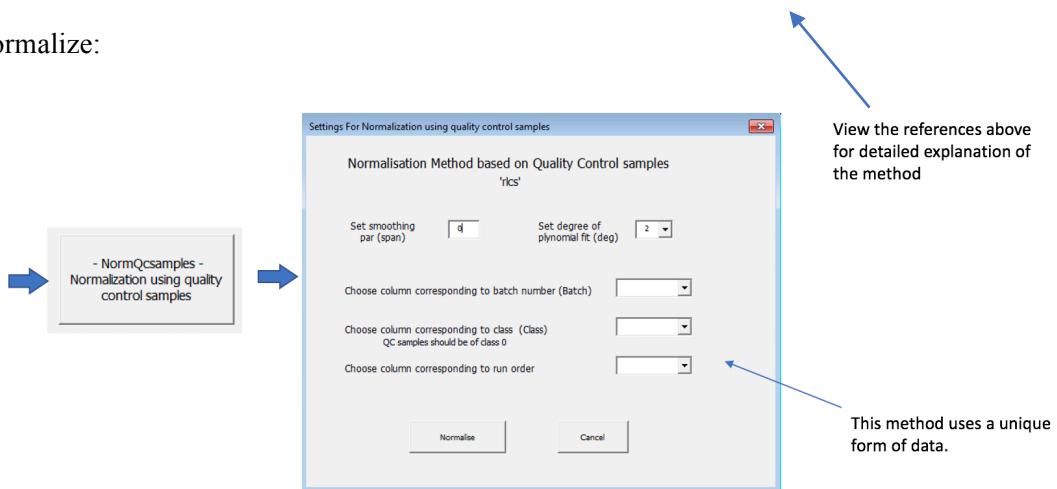
A quick visualization of the normalised data can be viewed by selecting the appropriate method generating the following plots

Location of the normalised featuredata

## NormQcsamples

This function is based on the quality control sample based robust LOESS (locally estimated scatterplot smoothing) signal correction (QC-RLSC) method as described by Dunn et al. (2011) and implemented statTarget (Luan 2017)

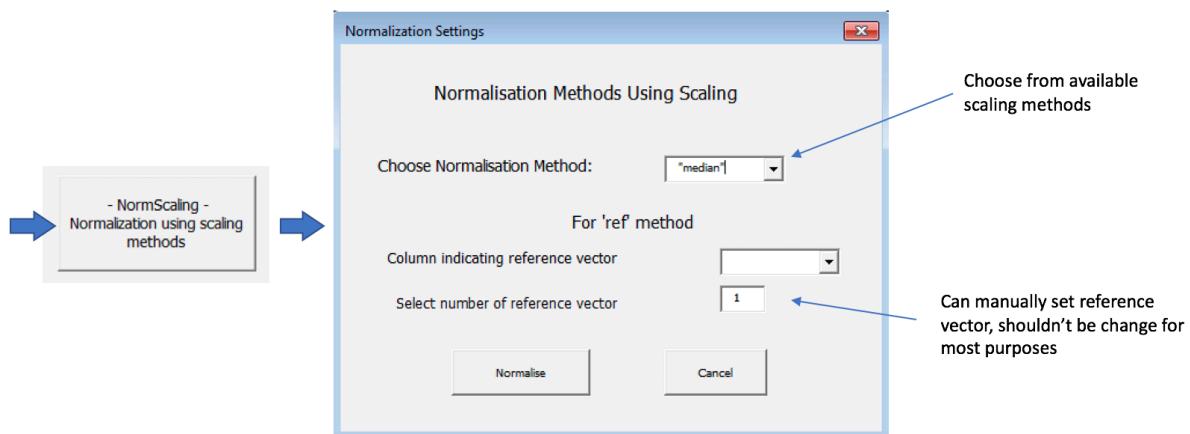
To Normalize:



## NormScaling

The scaling normalization methods (Scholz et al. 2004, Wang et al. (2003)) included in the package are normalization to a total sum, normalisation by the median or mean of each sample, and are denoted by *sum*, *median*, and *mean* respectively. The method *ref* normalises the metabolite abundances to a specific reference vector such as the sample weight or volume.

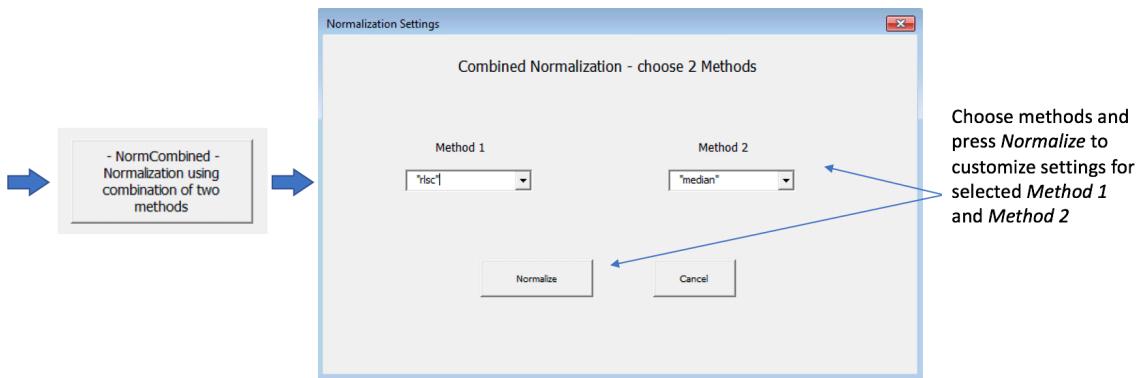
To Normalize:



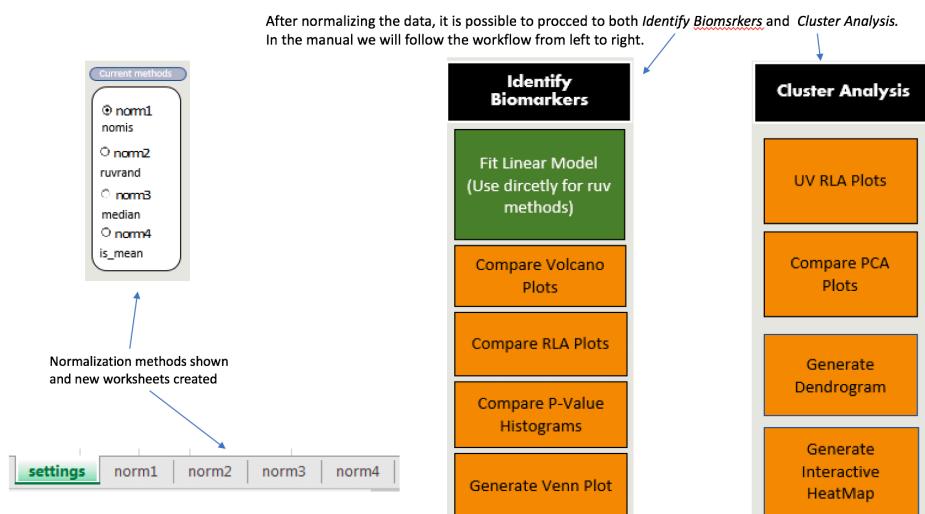
## NormCombined

In some circumstances, researchers use a combination of the above normalizations (i.e., one method followed by another). This can be achieved using the *NormCombined* function. The function defaults to employing 'rlsc' approach followed by the 'median'.

To Normalize:



Note that normalizing the data is not necessary to proceed to fitting a linear model although it is highly recommended to try a few normalization methods when analysis data.



## Assessing and choosing normalization methods

The criteria for assessing and choosing a normalization method implemented NormlizeMets have been described in detail by De Livera et al. (2012), De Livera et al. (2015) and J. A. Gagnon-Bartsch, Jacob, and Speed (2014).

### Identifying Biomarkers

To view and compare the biomarkers identified, first a linear model has to be fitted to the data.

#### **Fit Linear Model**

A linear model has to be fitted for every Normalization method that is to be used down the line for Biomarker identification. Setting from one ‘run’ of the Fit Linear model will be saved for the next.

To Fit Linear Model:

Available normalized data to be used for the linear model

Add factors that should be used when fitting the linear model . Factors can be viewed in the Formula box

Settings when fitting the linear model

Identify Biomarkers

- Fit Linear Model (Use directly for rvu methods)
- Compare Volcano Plots
- Compare RLA Plots
- Compare P-Value Histograms
- Generate Venn Plot

Choose Normalization Source/Target  
If no data present in location chosen, the unadjusted data will be used. For rvu2, the unadjusted data with missing values is used.

Add Factors for the model

Formula: +batch + gender + Age + bmi

For rvu2  
1 rvu2 method removes unwanted variation when fitting the linear model and is not suitable for combined normalisation. Log transformed data is used by default and any data in the current normalization source will be deleted !

Use rvu2 FALSE

K - number of unwanted variation factors NULL

Choose column with metabolite QC indicators (qcmetcs):

Compute moderated statistics FALSE

p-adjustment method "BH"

Significance level for confidence interval (ci\_alpha) 0.05

Fit Linear Model Cancel

My BiostatisticsProject.xlsx

A new sheet is created, it stores the linear model output. This data should not be altered manually as other functions will need to access it.

	A	B	C	D	E	F	G	H	I	J	K	L	M
1													
2	m_1	367.664000	1.86E-09	8.48970511	0.01846057	0.00000011	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	t
3	m_2	1.26E+03	1.86E-09	8.23E-94	0.02051323	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.23123818
4	m_3	39.83815528	1.93E-75	1.67E-25	6.70338880	0.020513485	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	3.70581385 -0.672814342
5	m_4	14.54000231	1.93E-11	3.76E-11	4.93533128	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	7.73517531 0.39043884
6	m_5	74.63010813	2.06E-48	3.23E-48	8.66379030	0.018454175	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	2.712459324 0.770537426
7	m_6	1.66E+02	1.93E-52	2.35E-52	7.80000000	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.346548175 0.33111104
8	m_7	15.47029338	3.49E-12	4.04E-12	8.07250299	0.04888541	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	3.590003368 0.454694912
9	m_8	146.23070707	4.02E-122	2.47E-120	8.10725099	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	-0.705951318 0.844774732
10	m_9	1226.58955	1.05E-116	3.18E-116	7.23E-116	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.87447684 0.188477432
11	m_10	3.00E+02	1.86E-90	8.48970511	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.29119384 0.15089945
12	m_11	309.103018	1.93E-75	8.79E-75	10.00005299	0.085395384	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	39.15770844 0.210375971
13	m_12	89.17162338	1.84E-42	3.23E-42	7.64515919	0.039796591	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.23123818 0.135172979
14	m_13	88.00000000	2.13E-42	3.68E-42	8.13077722	0.039334734	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.346548175 0.135172978
15	m_14	4.47790084	1.55E-57	1.37E-57	5.83E-57	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	13.33810000 0.404968264
16	m_15	45.54797024	4.31E-28	5.96E-28	5.23100527	0.057307074	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	14.0593205 0.17053694
17	m_16	149.098253	2.60E-55	5.81E-55	7.47397000	0.029229971	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	27.01482123 0.276260389
18	m_17	32.00000000	2.29E-22	2.80E-22	2.29E-22	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001
19	m_18	1.70E+00	4.55E-57	1.05E-57	5.74018116	0.112320640	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.73321001 0.31320000
20	m_19	84.1882905	4.09E-41	1.79E-41	5.9933042	0.057424040	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.31320000 0.17332100
21	m_20	13.128383	1.51E-10	1.71E-10	2.81292127	0.173470041	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.38505505 1.690139137
22	m_21	261.120608	1.24E-42	4.80E-42	7.95000000	0.020739347	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.23123818 0.126522002
23	m_22	1.00E+00	4.55E-57	1.05E-57	4.8933047	0.060000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.31320000 0.17332100
24	m_23	190.9438527	6.22E-62	1.59E-61	6.410546103	0.051413949	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	27.01482123 0.241040468
25	m_24	19.24080811	2.00E-06	1.02E-79	9.71953441	0.104950007	0.000919387	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.241040468 -0.778292248
26	m_25	19.24080811	4.24E-62	1.11E-61	8.28708118	0.050522111	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.241040468 0.2420443202
27	m_26	20.26700000	1.55E-19	1.66E-33	9.55E-05	0.200484200	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	1.87494394 1.823812386
28	m_27	60.18175603	1.33E-33	1.66E-33	8.95480878	0.204345414	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	0.00000001	-2.6471454 -0.13551198

Current methods

- norm1 nomis
- norm2 rvurand
- norm3 median
- norm4 rvu2

Plots

- Compare RLA Plots
- Generate Dendrogram
- Generate Interactive HeatMap

SVM Fit Classifier

- Correlation Analysis
- Compare Correlation Coefficients

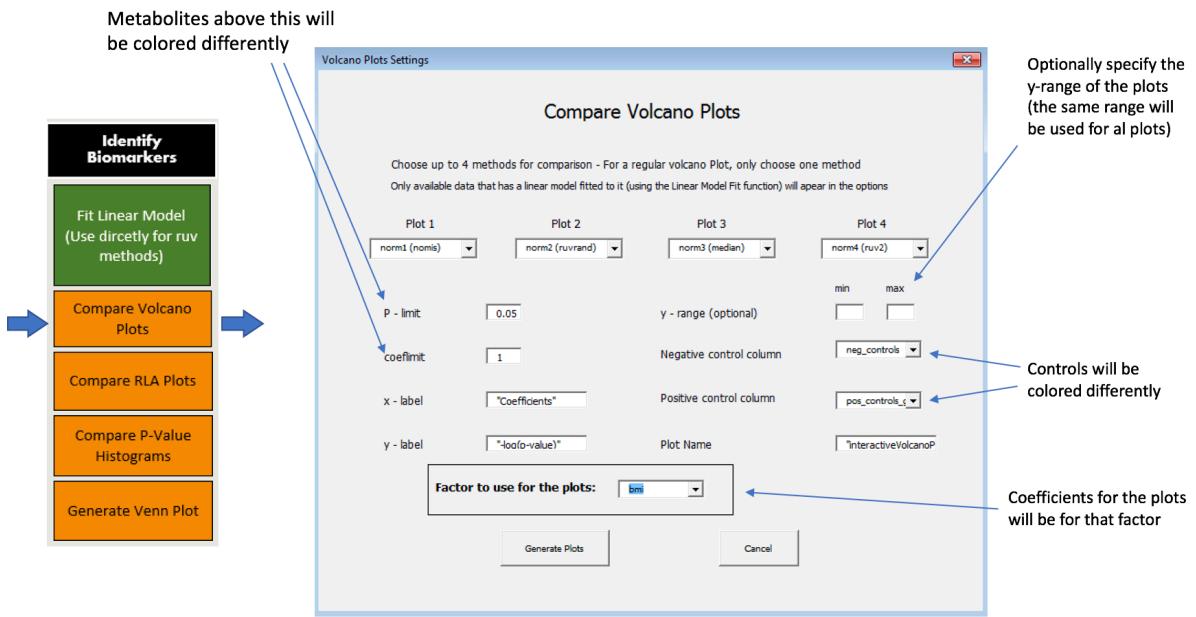
Manual - For instructions and detailed explanation of methods and functions press here

norm1 norm1LMResults norm2 norm2LMResults norm3 norm3LMResults norm4 norm4LMResults

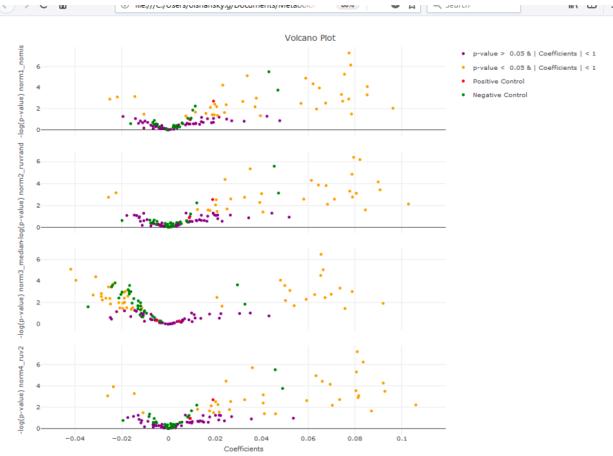
## Volcano Plots

Volcano plots are useful in identifying biomarkers and generally assessing the normalization.

## Compare Volcano Plots:



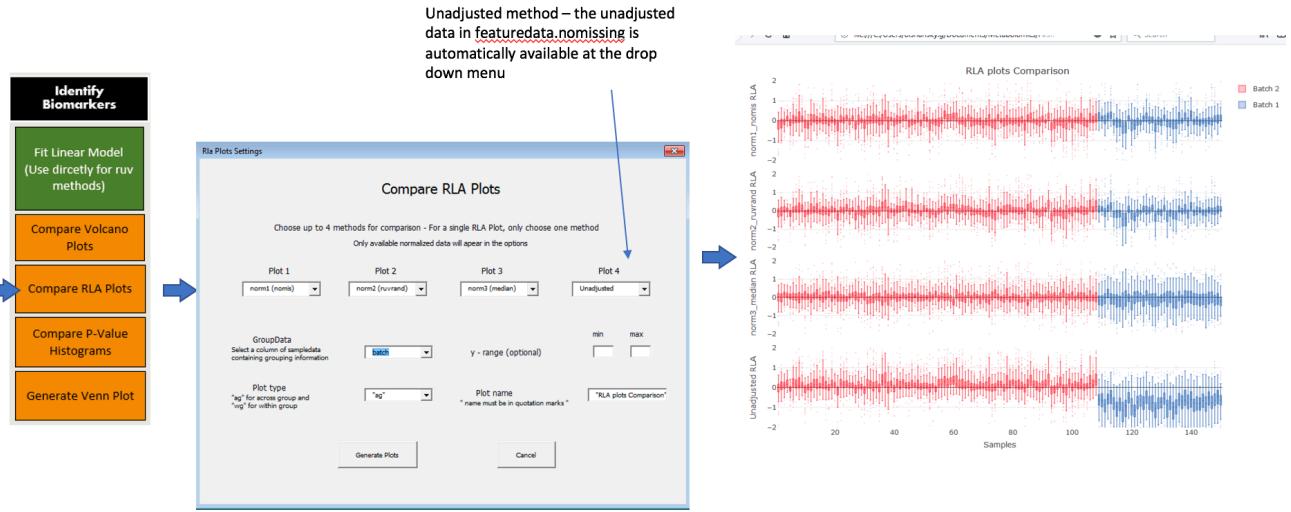
The plot is saved in the working directory and opens in the default browser.



## Compare RLA plots

Used to assess normalization by comparing relative log abundance plots, similar input to the Generate RLA plots function

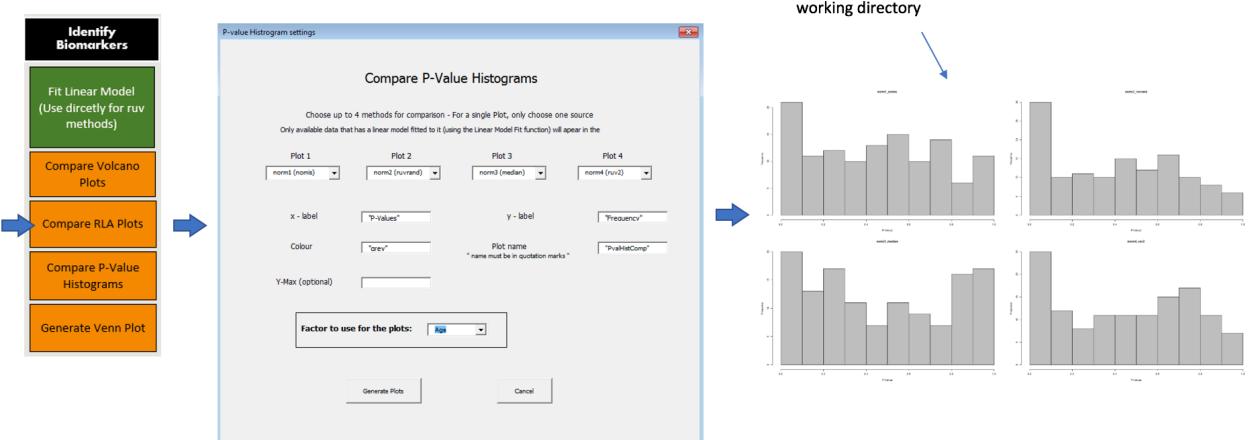
Compare RLA Plots:



## Compare P-Value Histograms

Compare histograms of the coefficient's p-values. The distribution of the p-values should be used to assess the success of the normalization.

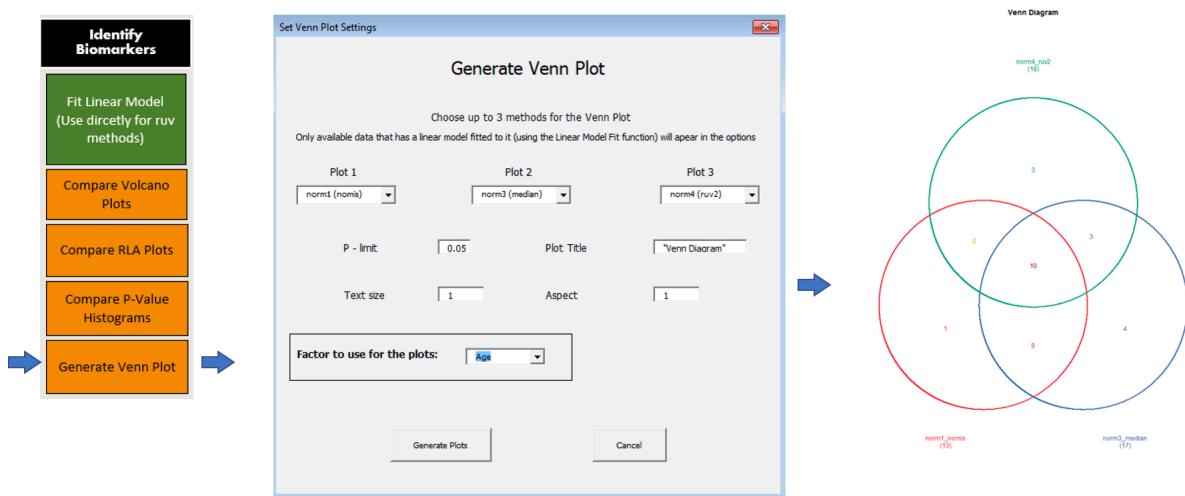
Compare P-Value histograms:



## Generate Venn Plot

Generates a Venn plot that compares the biomarkers identified by the different normalization methods.

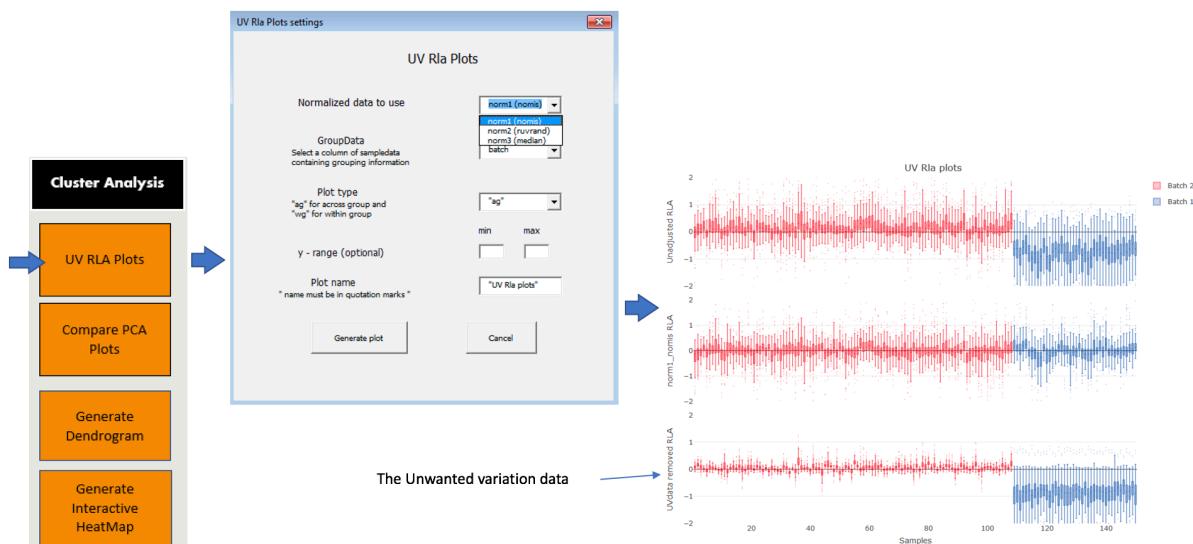
Generate Venn Plot:



## Cluster Analysis

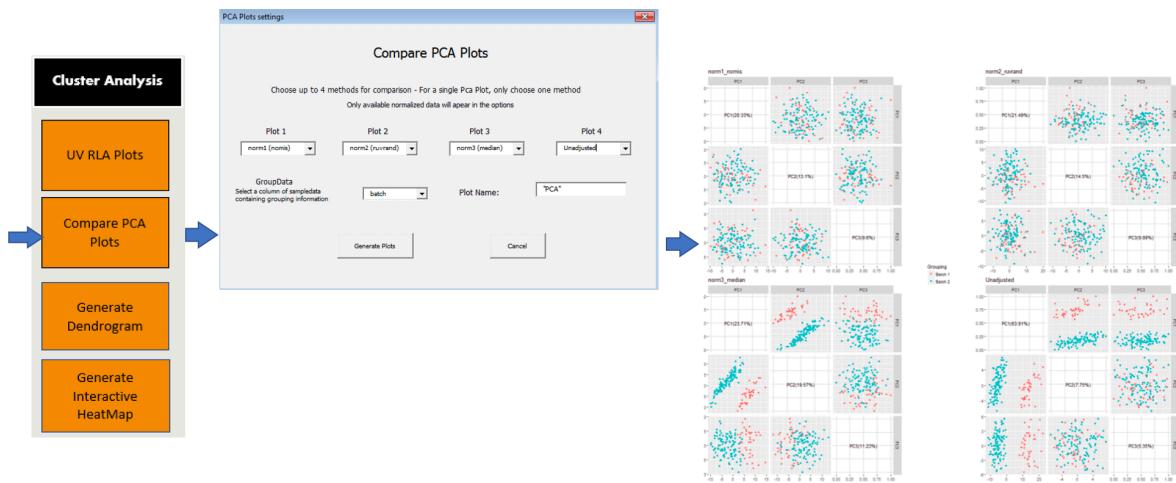
### UV RLA plots

Unwanted Variation relative log abundance plots enable visualisation of the unwanted variation removed by each normalization method.



## Compare PCA Plots

Compare principal component multi-plots for differed normalization methods.



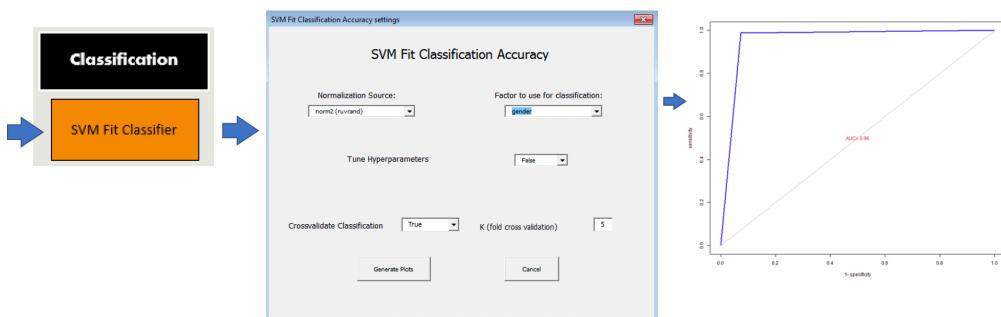
**Generate Dendrogram** and **Generate Interactive HeatMap** are identical to those discussed in the *Visualize the Data* section. The user has to choose the normalized data to be used.

## Classification

Classification accuracy is a good way to assess the success of a given normalization method

### SVM Fit

The Support Vector Machines method is used to classify the data, the classification accuracy is then assessed based on the factor specified.



## Correlation Analysis

It is important to look at correlation coefficients when normalizing data.

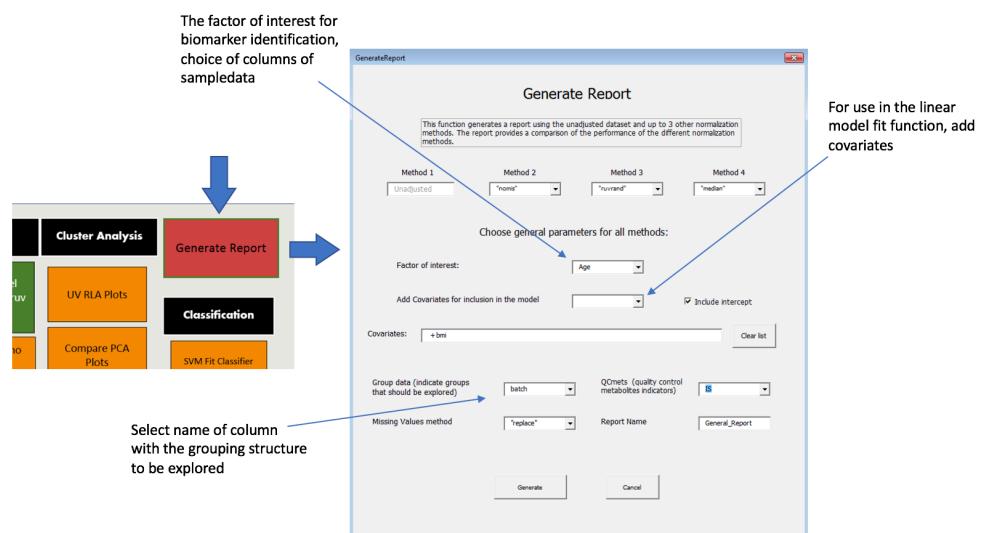
### Compare Correlation coefficients



## Generate Report

The *Generate Report* function generates an interactive report based on basic user input. There is a choice of up to 3 normalization methods to be included together with the unadjusted data. The report includes various plots and diagnostic to assess the normalization. Guidance on interpretation of the various plots, together with notes of what the user should look for when assessing the results is provided in the generated document.

### Generate Report:



An example report is available in your downloaded ExNormalizeMetsSetup.