

MetaboAnalystR 2.0 Workflow: From Raw Spectra to Biological Insights

Jasmine Chong, Jeff Xia

2019-03-06

1. Introduction

We present MetaboAnalystR 2.0, which aims to address two important gaps left in its previous version. First, raw spectral processing - the previous version offered very limited support for raw spectra processing and peak annotation. Therefore, we have implemented comprehensive support for raw LC-MS spectral data processing including peak picking, peak alignment and peak annotations. Second, we have enhanced support for functional interpretation directly from m/z peaks. In addition to an efficient implementation of the mummichog algorithm (PMID: 23861661), we have added a new method to support pathway activity prediction based on the well-established GSEA algorithm (PMID: 16199517). In this tutorial, we showcase how to utilize MetaboAnalyst 2.0 to perform a comprehensive end-to-end metabolomics data workflow. The dataset consists of a subset of pediatric IBD stool samples obtained from the Integrative Human Microbiome Project Consortium (<https://ibdmdb.org/>).

2. Raw MS Data Preprocessing

Three main wrapper functions have been implemented for metabolomics data processing based on XCMS (PMIDs: 16448051, 19040729, and 20671148; version 3.4.4) and CAMERA (PMID: 22111785; version 1.38.1) including: (i) the *ImportRawMSData* function for reading in raw data files, (ii) the *PerformPeakProfiling* function for peak picking and alignment, and (iii) the *PerformPeakAnnotation* function for peak annotation. These functions are described below in further detail.

The *ImportRawMSData* function reads in raw MS data files and saves it as an OnDiskMSnExp object. The function outputs two plots - the Total Ion Chromatogram (TIC) which provides an overview of the entire spectra, and the Base Peak Chromatogram (BPC) which is a cleaner profile of the spectra based on the most abundant signals. These plots are useful to inform the setting of parameters downstream. For users who wish to view a peak of interest, an Extracted Ion Chromatogram (EIC) can be generated using the *PlotEIC* function.

The *PerformPeakProfiling* function is a wrapper of several XCMS R functions that performs peak detection, alignment, and grouping in a single step. The resulting peaks are outputted as a XCMSnExp object. The function also generates two diagnostic plots including a retention time adjustment map, and a PCA plot showing the overall sample clustering prior to data cleaning and statistical analysis. Users can specify several parameters such as the mass accuracy, peak width, and retention time range using the *SetPeakParam* function to optimize the peak picking function.

The *PerformPeakAnnotation* function annotates isotope and adduct peaks using the CAMERA R package (PMID: 22111785). It outputs the result as a CSV file ("annotated_peaklist.csv") and saves the annotated peaks as an xsAnnotate object. Finally, the peak list is formatted to the correct structure for MetaboAnalystR and filtered based upon user's specifications using the *FormatPeakList* function. This function permits the filtering of adducts (i.e. removal of all adducts except for [M+H]⁺/[M-H]⁻) and filtering of isotopes (i.e. removal of all isotopes except for monoisotopic peaks). The goal of filtering peaks is to remove degenerative signals and reduce the file size.

1.1 Example data download

Prior to starting this tutorial, please download the 12 example mzML files, containing fecal samples from 6 pediatric IBD patients and 6 healthy controls. These patients are all male, aged 11-16 years old. The data is available for download here: <https://www.dropbox.com/sh/o1nqb7ioyo868nt/AADWDKRsCo6tZkWvDcFBGGhJa?dl=0>. Ensure that you unzip these files into an empty folder and set your working directory to this new folder. If you have issues downloading/unzipping the zipped folder, try downloading these 6 samples individually.

1.2 Preprocessing of example data

We will first import the raw MS data using the *ImportRawMSData* function. Use the *list.files()* functions to view the order of the files within the folder. As you can see, the first 6 samples are the controls, and the last 6 samples are IBD patients - denoted with a CD at the beginning of their file names. Therefore, within *ImportRawMSData* we set *grpA* to “Healthy”, *numA* to 6, *grpB* to “IBD”, and *numB* to 6. The *format* refers to the format of the resulting TIC and BPC plots (i.e. “png”, “jpg”), the *dpi* refers to the resolution (typically 300 dpi is required for publication quality images) and *width* refers to the size of the plot. All plots generated in *MetaboAnalystR* can be found within your current working directory.

```
# Set working directory to 12 clinical samples
setwd("~/Desktop/iHMP/")

# Load the MetaboAnalystR package
library("MetaboAnalystR")

# Import mzML files
rawData <- ImportRawMSData(grpA = "Healthy", numA = 6, grpB = "IBD", numB = 6,
                           format = "png", dpi = 72, width = 9)
```

Next, we will set the parameters for peak picking using the *SetPeakParam* function. This function sets all the parameters used for downstream pre-processing of user's raw MS data such as the mass error (ppm), signal to noise threshold, and bandwidth. Use *?SetPeakParam* for further information of all available parameters. In this case, we set the mass error to 5. Following the setting of parameters, we will use the *PerformPeakProfiling* function to perform peak picking, grouping, retention time correction, and alignment. The subset of 10 samples used in this vignette were run on a Dell XPS 15 9570 laptop with a dual-booted Ubuntu 16.04 and 16 GB of ram. When this function was run, 13.6/16 GB of memory were used during peak picking, causing the function to fail. Therefore for users who may also run into memory issues on their desktop/laptop, either limit the number of cores used, as *xcms* uses multi-core parallel processing - with: *register(bpstart(MulticoreParam(3)))* ~ which for example, sets the number of cores used to 3, or use a smaller subset of IBD samples to demo the workflow.

```
# Set parameters for peak profiling
peakParams <- SetPeakParam(ppm = 5)

# Perform peak profiling
extPeaks <- PerformPeakProfiling(rawData, peakParams, rtPlot = TRUE, pcaPlot = TRUE, labels = TRUE,
                                format = "png", dpi = 300, width = 9)
```

The *PerformPeakProfiling* function outputs a PCA plot of the detected peaks per sample. From this plot, we can see an ok separation between the Healthy and IBD patients. Remember that this workflow utilizes only a subset of patients, and for good results, the number of samples should be maximized.

Following peak profiling we will perform peak annotation. First, set the parameters for peak annotation using the *SetAnnotationParam* function. In this case, specify that the samples were collected in negative ion mode - “negative”, and the allowed variance for the search (for adduct annotation) to 0.005. Next, use the

PerformPeakAnnotation function to annotate the detected peaks to potential isotopes and adducts. This will output a csv file named “annotated_peaklist.csv”, which contains the entire list of all identified m/z features, their retention time, and their annotations. In this case, we identified ~11000 features. Finally, use the *FormatPeakList* function to format this annotated_peaklist into the correct format to be used by MetaboAnalystR (and the web version). In this function, users can choose to filter the peaks, i.e. remove features that are missing in more than 75% of features per group and filter out all isotopes. The function will output a file named “metaboanalyst_input.csv”, which will be used the package to read in the peak table.

```
# Set parameters for peak annotation
annParams <- SetAnnotationParam(polarity = "negative", mz_abs_add = 0.005)

# Perform peak annotation
annotPeaks <- PerformPeakAnnotation(extPeaks, annParams)

# Format and filter the peak list for MetaboAnalystR
maPeaks <- FormatPeakList(annotPeaks, annParams, filtIso = TRUE, filtAdducts = FALSE,
                          missPercent = 0.75)
```

3. Data Processing and Statistical Analysis

Remaining in the same working directory, we will read in the filtered peak table. Then, we will replace missing values with a small value (half of the minimum positive value within the original data). Next we will filter variables out non-informative signals, then normalize the samples to their median and perform log transformation. Finally, we will perform t-test analysis to identify differentially enriched features, setting the FDR-adjusted p-value threshold to 0.25.

```
# Initialize the mSetObj
mSet<-InitDataObjects("pktable", "stat", FALSE)

# Read in the filtered peak list
mSet<-Read.TextData(mSet, "metaboanalyst_input.csv", "colu", "disc")

# Perform data processing
mSet<-SanityCheckData(mSet)
mSet<-ReplaceMin(mSet);
mSet<-FilterVariable(mSet, "iqr", "F", 25)
mSet<-PreparePrenormData(mSet)
mSet<-Normalization(mSet, "MedianNorm", "LogNorm", "NULL", ratio=FALSE, ratioNum=20)

# Perform t-test
mSet<-Ttests.Anal(mSet, F, 0.25, FALSE, TRUE)
```

4. From Raw Peaks to Biological Insights

Even with high mass accuracy afforded by the current high-resolution MS platforms, it is often impossible to uniquely identify a given peak based on its mass alone. To get around this issue, a key concept is to shift the unit of analysis from individual compounds to individual pathways or a group of functionally related compounds (i.e. metabolite sets (PMID: 20457745)). The general assumption is that the collective behavior of a group is more robust against a certain degree of random errors of individuals. The mummichog algorithm is the first implementation of this concept to infer pathway activities from a ranked MS peaks identified in untargeted metabolomics. This algorithm is available in MetaboAnalystR (the *PerformMummichog* function). The original algorithm implements an over-representation analysis (ORA) method to evaluate pathway-level enrichment based on significant features. Users need to specify a pre-defined cutoff based on either t-statistics

or fold changes. A complementary approach is the Gene Set Enrichment Analysis (GSEA) method, a widely used method to extract biological meaning from a ranked gene list (PMID: 16199517). Unlike ORA, this method considers the overall ranks without using a significance cutoff. It is able to detect subtle and consistent changes which can be missed from ORA methods. The MetaboAnalystR 2.0 offers GSEA (the PerformGSEA function) using the high-performance fgsea R package (<https://www.biorxiv.org/content/10.1101/060012v1>).

To perform these analyses, users must upload a table containing three columns - m/z features, p-values, and t-scores or fold-change values. In this case, as we performed t-test analysis on the peak table, we can use the new function *Convert2Mummichog*, which converts the results of the t-test to the proper format for mummichog analysis. The output of this function is a file titled “mummichog_input_XX.txt”, with XX replaced by the current date.

Next, we follow the typical MS Peaks to Pathways workflow, setting the instrument options to “negative” and the p-value cut-off to 0.25. Then, we will perform the original mummichog algorithm, using the *PerformMummichog* function. We will select the hsa_mfn library, and set the permutations to 1000. This may take some time.

```
# Convert the output of the t-test to the proper input for mummichog analysis
mSet<-Convert2Mummichog(mSet)
```

```
# Read in the ranked peak list (don't forget to change the file name!)
mSet<-Read.PeakListData(mSet, "mummichog_input_2019-03-06.txt");
mSet<-UpdateMummichogParameters(mSet, "5", "negative", 0.25);
mSet<-SanityCheckMummichogData(mSet)
```

```
# Perform original mummichog algorithm
```

```
mSet<- PerformMummichog(mSet, "hsa_mfn", permNum = 1000)
```

```
# View the top pathway hits using head(mSet$mummi.resmat)
```

```
# > head(mSet$mummi.resmat)
```

#	Pathway	total	Hits.total	Hits.sig	EASE
# Vitamin E metabolism		54	33	4	0.0054183
# Linoleate metabolism		46	20	1	1.0000000
# Carnitine shuttle		72	20	1	1.0000000
# Androgen and estrogen biosynthesis and metabolism		95	65	1	1.0000000
# Drug metabolism - cytochrome P450		53	42	1	1.0000000
# Ascorbate (Vitamin C) and Aldarate Metabolism		29	19	1	1.0000000

Now we will perform GSEA, using the *PerformGSEA* function. This function is much faster than the original implementation. As you can see, Vitamin E metabolism is consistently identified between the two algorithms. Vitamin E has been shown to be decreased in patients with Crohn’s disease (PMID: 22488830) and were shown to have anti-inflammatory properties (PMID: 30237488). Finally, we can visualize the comparison of the pathway analysis results of the two algorithms with a scatter plot using the *PlotIntegPaths* function. This function will check that you performed both methods (original mummichog + GSEA), and that both methods used the same selected pathway library. By adding labels = TRUE to the function, we can generate a plot with all of the dots annotated - this may be cluttered for those with many pathways identified.

```
# Perform GSEA
```

```
mSet<- PerformGSEA(mSet, "hsa_mfn", permNum = 1000)
```

```
# View the top pathway hits using head(mSet$mummi.gsea.resmat)
```

```
# > head(mSet$mummi.gsea.resmat)
```

#	Pathway	Total	Hits	P_val	P_adj	NES
# Linoleate metabolism	"46"		"1"	"0.998"	"0.998"	"-0.4419"
# Carnitine shuttle	"72"		"1"	"0.998"	"0.998"	"-0.2882"
# Vitamin E metabolism	"54"		"4"	"0.8219"	"0.998"	"-0.7053"

```

# Androgen and estrogen biosynthesis and metabolism "95" "1" "0.998" "0.998" "0.3674"
# Drug metabolism - cytochrome P450 "53" "1" "0.998" "0.998" "0.276"
# Ascorbate (Vitamin C) and Aldarate Metabolism "29" "1" "0.998" "0.998" "-0.6419"

# Plot the integration of the two algorithms
mSet <- PlotIntegPaths(mSet, dpi = 300, width = 10, format = "jpg")

# Plot with all dots annotated
mSet <- PlotIntegPaths(mSet, dpi = 300, width = 10, format = "jpg", labels = TRUE)

```

4.1 MetaboAnalystR 2.0 IBD Case Study

In this vignette, we demonstrated the workflow using a subset of IBD patients and healthy controls. In our manuscript “MetaboAnalystR 2.0: From Raw Spectra to Biological Insights”, we performed this entire workflow using 24 pediatric IBD and 24 healthy controls. The entire R history, resulting annotated peak list, and filtered peak table are available for download from: https://github.com/jsychong/MetaboAnalystR/tree/master/MetaboAnalystR_2_Supplementary_Data

To recreate our findings in the manuscript, using the provided peak table from the GitHub repository, use the following code snippet. Ensure the “iHMP_48_metaboanalyst_input.csv” file is located in your current working directory.

```

mSet<-InitDataObjects("pktable", "stat", FALSE)
mSet<-Read.TextData(mSet, "iHMP_48_metaboanalyst_input.csv", "colu", "disc")
mSet<-SanityCheckData(mSet)
mSet<-ReplaceMin(mSet);
mSet<-FilterVariable(mSet, "iqr", "F", 25)
mSet<-PreparePrenormData(mSet)
mSet<-Normalization(mSet, "MedianNorm", "LogNorm", "NULL", ratio=FALSE, ratioNum=20)
mSet<-Ttests.Anal(mSet, F, 0.05, FALSE, TRUE)
mSet<-Convert2Mummichog(mSet)

# RENAME FILE TO CREATED mummichog_input
mSet<-Read.PeakListData(mSet, "mummichog_input_2019-03-06.txt");
mSet<-UpdateMummichogParameters(mSet, "5", "negative", 0.25);
mSet<-SanityCheckMummichogData(mSet)

# First perform original mummichog algorithm
mSet<- PerformMummichog(mSet, "hsa_mfn", permNum = 1000)

# View the top enriched pathways
# > head(mSet$mummi.resmat)
#
# Pathway total Hits.total Hits.sig EASE FET Ga
# Bile acid biosynthesis 82 52 29 0.0062805 0.0028215 0.00027
# Vitamin E metabolism 54 33 20 0.0094632 0.0035614 0.00028
# Fatty Acid Metabolism 63 11 9 0.0152200 0.0026797 0.00029
# Vitamin D3 (cholecalciferol) metabolism 16 10 8 0.0314750 0.0061618 0.00031
# De novo fatty acid biosynthesis 106 15 10 0.0520170 0.0162020 0.00035
# Fatty acid activation 74 15 10 0.0520170 0.0162020 0.00035

```

Using the original mummichog algorithm, the top enriched pathways are bile acid biosynthesis and Vitamin E metabolism. Both metabolic pathways are known to be perturbed in IBD patients.

```
mSet<- PerformGSEA(mSet, "hsa_mfn", permNum = 1000)
```

```
# > head(mSet$mummi.gsea.resmat)
#
```

	Pathway_Total	Hits	P_val	P_adj	NES
# Bile acid biosynthesis	"82"	"29"	"0.001931"	"0.1506"	"-1.4"
# Vitamin D3 (cholecalciferol) metabolism	"16"	"10"	"0.02495"	"0.7985"	"-1.4"
# Biopterin metabolism	"22"	"7"	"0.03071"	"0.7985"	"-1.4"
# Putative anti-Inflammatory metabolites formation from EPA	"27"	"4"	"0.1062"	"0.8288"	"1.4"
# Androgen and estrogen biosynthesis and metabolism	"95"	"11"	"0.0631"	"0.8288"	"-1.4"
# Arachidonic acid metabolism	"95"	"8"	"0.09623"	"0.8288"	"1.4"

Using the GSEA algorithm, the top enriched pathways are bile acid biosynthesis and Vitamin D3 metabolism. When we plot the results of both pathway analysis results (code below), we can see that by looking at the size and color of the dots (representing each pathway, big + red is considered most significant), bile acid biosynthesis and Vitamin D3 metabolism are consistently shown to be perturbed in these IBD patients. This workflow demonstrates how biological insights can be easily obtained directly from m/z peaks, without the need for a priori metabolite identification.

```
# View plot with no dots labeled
mSet <- PlotIntegPaths(mSet, dpi = 300, width = 10, format = "jpg")

# View plot with all dots labeled
mSet <- PlotIntegPaths(mSet, dpi = 300, width = 10, format = "jpg", labels = TRUE)
```