

flagme: Fragment-level analysis of GC-MS-based metabolomics data

Mark Robinson

`mrobinson@wehi.edu.au`

Riccardo Romoli

`riccardo.romoli@unifi.it`

April 25, 2023

1 Introduction

This document gives a brief introduction to the **flagme** package, which is designed to process, visualise and statistically analyze sets of GC-MS samples. The ideas discussed here were originally designed with GC-MS-based metabolomics in mind, but indeed some of the methods and visualizations could be useful for LC-MS data sets. The *fragment-level analysis* though, takes advantage of the rich fragmentation patterns observed from electron interaction (EI) ionization.

There are many aspects of data processing for GC-MS data. Generally, algorithms are run separately on each sample to detect features, or *peaks* (e.g. AMDIS). Due to retention time shifts from run-to-run, an alignment algorithm is employed to allow the matching of the same feature across multiple samples. Alternatively, if known standards are introduced to the samples, retention *indices* can be computed for each peak and used for alignment. After peaks are matched across all samples, further processing steps are employed to create a matrix of abundances, leading into detecting differences in abundance.

Many of these data processing steps are prone to errors and they often tend to be black boxes. But, with effective exploratory data analysis, many of the pitfalls can be avoided and any problems can be fixed before proceeding to the downstream statistical analysis. The package provides various visualizations to ensure the methods applied are not black boxes.

The **flagme** package gives a complete suite of methods to go through all common stages of data processing. In addition, R is especially well suited to the downstream data analysis tasks since it is very rich in analysis tools and has excellent visualization capabilities. In addition, it is freely available (www.r-project.org), extensible and there is a growing community of users and developers. For routine analyses, graphical user interfaces could be designed.

2 Reading and visualizing GC-MS data

To run these examples, you must have the **gcspikelite** package installed. This data package contains several GC-MS samples from a spike-in experiment we designed to interrogate data processing methods. So, first, we load the packages:

To load the data and corresponding peak detection results, we simply create vectors of the file-names and create a **peakDataset** object. Note that we can speed up the import time by setting the retention time range to a subset of the elution, as below:

```
> gcmsPath <- paste(find.package("gcspikelite"), "data", sep="/")
> data(targets)
> cdfFiles <- paste(gcmsPath, targets$FileName, sep="/")
```

```

> eluFiles <- gsub("CDF", "ELU", cdfFiles)
> pd <- peaksDataset(cdfFiles, mz=seq(50,550), rtrange=c(7.5,8.5))

Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_468.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_474.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_475.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_485.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_493.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_496.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_470.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_471.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_479.CDF

> pd <- addAMDISPeaks(pd, eluFiles)

Reading retention time range: 7.500133 8.499917
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_468.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_474.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_475.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_485.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_493.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_496.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_470.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_471.ELU ... Done.
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_479.ELU ... Done.

> pd

An object of class "peaksDataset"
9 samples: 0709_468 0709_474 0709_475 0709_485 0709_493 0709_496 0709_470 0709_471 0709_479
501 m/z bins - range: ( 50 550 )
scans: 175 175 175 175 175 174 175 175 175
peaks: 24 23 26 20 27 24 24 25 21

```

Here, we have added peaks from AMDIS, a well known and mature algorithm for deconvolution of GC-MS data. For GC-TOF-MS data, we have implemented a parser for the **ChromaTOF** output (see the analogous **addChromaTOFPeaks** function). The **addXCMSPeaks** allows to use all the XCMS peak-picking algorithms; using this approach it is also possible to elaborate the raw data file from within R instead of using an external software. In particular the function reads the raw data using XCMS, group each extracted ion according to their retention time using CAMERA and attaches them to an already created **peaksDataset** object:

```

> pd.2 <- peaksDataset(cdfFiles[1:3], mz = seq(50, 550), rtrange = c(7.5, 8.5))

Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_468.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_474.CDF
Reading F:/biocbuild/bbs-3.17-bioc/R/library/gcspikelite/data/0709_475.CDF

> mfp <- xcms::MatchedFilterParam(fwhm = 10, snthresh = 5)
> pd.2 <- addXCMSPeaks(cdfFiles[1:3], pd.2, settings = mfp)

Start grouping after retention time.
Created 164 pseudospectra.
Start grouping after correlation.

```

```
Calculating peak correlations in 164 Groups...
% finished: 10 20 30 40 50 60 70 80 90 100
```

```
Calculating graph cross linking in 164 Groups...
% finished: 10 20 30 40 50 60 70 80 90 100
New number of ps-groups: 250
xsAnnotate has now 250 groups, instead of 164
Start grouping after retention time.
Created 153 pseudospectra.
Start grouping after correlation.
```

```
Calculating peak correlations in 153 Groups...
% finished: 10 20 30 40 50 60 70 80 90 100
```

```
Calculating graph cross linking in 153 Groups...
% finished: 10 20 30 40 50 60 70 80 90 100
New number of ps-groups: 245
xsAnnotate has now 245 groups, instead of 153
Start grouping after retention time.
Created 163 pseudospectra.
Start grouping after correlation.
```

```
Calculating peak correlations in 163 Groups...
% finished: 10 20 30 40 50 60 70 80 90 100
```

```
Calculating graph cross linking in 163 Groups...
% finished: 10 20 30 40 50 60 70 80 90 100
New number of ps-groups: 245
xsAnnotate has now 245 groups, instead of 163
```

```
> pd.2
```

```
An object of class "peaksDataset"
3 samples: 0709_468 0709_474 0709_475
418 m/z bins - range: ( 51 468 )
scans: 175 175 175
peaks: 47 45 44
```

The possibility to work using computer cluster will be added in the future.

Regardless of platform and peak detection algorithm, a useful visualization of a set of samples is the set of total ion currents (TIC), or extracted ion currents (XICs). To view TICs, you can call:

```
> plotChrom(pd, rtrange=c(7.5,8.5), plotPeaks=TRUE, plotPeakLabels=TRUE,
+           max.near=8, how.near=0.5, col=rep(c("blue","red","black"), each=3))
```

Note here the little *hashes* represent the detected peaks and are labelled with an integer index. One of the main challenges is to match these peak detections across several samples, given that they appear at slightly different times in different runs.

For XICs, you need to give the indices (of `pd@mz`, the grid of mass-to-charge values) that you want to plot through the `mzind` argument. This could be a single ion (e.g. `mzind=24`) or could be a range of indices if multiple ions are of interest (e.g. `mzind=c(24,25,98,99)`).

There are several other features within the `plot` command on `peaksDataset` objects that can be useful. See `?plot` (and select the `flagme` version) for full details.

Another useful visualization, at least for individual samples, is a 2D heatmap of intensity. Such plots can be enlightening, especially when peak detection results are overlaid. For example (with detected fragment peaks from AMDIS shown in white):

```
> r <- 1
> plotImage(pd, run=r, rtrange=c(7.5,8.5), main="")
> v <- which(pd@peaksdata[[r]] > 0, arr.ind=TRUE) # find detected peaks
> abline(v=pd@peaksrt[[r]])
> points(pd@peaksrt[[r]][v[,2]], pd@mz[v[,1]], pch=19, cex=.6, col="white")
```

3 Pairwise alignment with dynamic programming algorithm

One of the first challenges of GC-MS data is the matching of detected peaks (i.e. metabolites) across several samples. Although gas chromatography is quite robust, there can be some drift in the elution time of the same analyte from run to run. We have devised a strategy, based on dynamic programming, that takes into account both the similarity in spectrum (at the apex of the called peak) and the similarity in retention time, without requiring the identity of each peak; this matching uses the data alone. If each sample gives a ‘peak list’ of the detected peaks (such as that from AMDIS that we have attached to our `peaksDataset` object), the challenge is to introduce gaps into these lists such that they are best aligned. From this a matrix of retention times or a matrix of peak abundances can be extracted for further statistical analysis, visualization and interpretation. For this matching, we created a procedure analogous to a multiple *sequence* alignment.

To highlight the dynamic programming-based alignment strategy, we first illustrate a pairwise alignment of two peak lists. This example also illustrates the selection of parameters necessary for the alignment. From the data read in previously, let us consider the alignment of two samples, denoted 0709_468 and 0709_474. First, a similarity matrix for two samples is calculated. This is calculated based on a scoring function and takes into account the similarity in retention time and in the similarity of the apex spectra, according to:

$$S_{ij}(D) = \frac{\sum_{k=1}^K x_{ik}y_{jk}}{\sqrt{\sum_{k=1}^K x_{ik}^2 \cdot \sum_{k=1}^K y_{jk}^2}} \cdot \exp\left(-\frac{1}{2} \frac{(t_i - t_j)^2}{D^2}\right)$$

where i is the index of the peak in the first sample and j is the index of the peak in the second sample, \mathbf{x}_i and \mathbf{y}_j are the spectra vectors and t_i and t_j are their respective retention times. As you can see, there are two components to the similarity: spectra similarity (left term) and similarity in retention time (right term). Of course, other metrics for spectra similarity are feasible. Ask the author if you want to see other metrics implemented. We have some non-optimized code for a few alternative metrics.

The peak alignment algorithm, much like sequence alignments, requires a `gap` parameter to be set, here a number between 0 and 1. A high gap penalty discourages gaps when matching the two lists of peaks and a low gap penalty allows gaps at a very low *cost*. We find that a gap penalty in the middle range (0.4-0.6) works well for GC-MS peak matching. Another parameter, `D`, modulates the impact of the difference in retention time penalty. A large value for `D` essentially eliminates the effect. Generally, we set this parameter to be a bit larger than the average width of a peak, allowing a little flexibility for retention time shifts between samples. Keep in mind the `D` parameter should be set on the scale (i.e. seconds or minutes) of the `peaksrt` slot of the `peaksDataset` object. The next example shows the effect of the `gap` and `D` penalty on the matching of a small ranges of peaks.

```
> Ds <- c(0.1, 10, 0.1, 0.1)
> gaps <- c(0.5, 0.5, 0.1, 0.9)
> par(mfrow=c(2,2), mai=c(0.8466,0.4806,0.4806,0.1486))
> for(i in 1:4){
+   pa <- peaksAlignment(pd@peaksdata[[1]], pd@peaksdata[[2]],
```

```

+                 pd@peaksrt[[1]], pd@peaksrt[[2]], D=Ds[i],
+                 gap=gaps[i], metric=1, type=1, compress = FALSE)
+ plotAlignment(pa, xlim=c(0, 17), ylim=c(0, 16), matchCol="yellow",
+              main=paste("D=", Ds[i], " gap=", gaps[i], sep=""))
+ }

[peaksAlignment] Comparing 24 peaks to 23 peaks -- gap= 0.5 D= 0.1 , metric= 1 , type= 1
[peaksAlignment] 21 matched. Similarity= 0.2308905
[peaksAlignment] Comparing 24 peaks to 23 peaks -- gap= 0.5 D= 10 , metric= 1 , type= 1
[peaksAlignment] 21 matched. Similarity= 0.2180835
[peaksAlignment] Comparing 24 peaks to 23 peaks -- gap= 0.1 D= 0.1 , metric= 1 , type= 1
[peaksAlignment] 15 matched. Similarity= 0.01170268
[peaksAlignment] Comparing 24 peaks to 23 peaks -- gap= 0.9 D= 0.1 , metric= 1 , type= 1
[peaksAlignment] 22 matched. Similarity= 0.2785359

```

You might ask: is the `flagme` package useful without peak detection results? Possibly. There have been some developments in alignment (generally on LC-MS proteomics experiments) without peak/feature detection, such as Prince et al. 2006, where a very similar dynamic programming is used for a pairwise alignment. We have experimented with alignments without using the peaks, but do not have any convincing results. It does introduce a new set of challenges in terms of highlighting differentially abundant metabolites. However, in the `peaksAlignment` routine above (and those mentioned below), you can set `usePeaks=FALSE` in order to do *scan*-based alignments instead of *peak*-based alignments. In addition, the `flagme` package may be useful simply for its bare-bones dynamic programming algorithm.

3.1 Normalizing retention time score to drift estimates

In what is mentioned above for pairwise alignments, we are penalizing for differences in retention times that are non-zero. But, as you can see from the TICs, some differences in retention time are consistent. For example, all of the peaks from sample 0709_485 elute at later times than peaks from sample 0709_496. We should be able to estimate this drift and normalize the time penalty to that estimate, instead of penalizing to zero. That is, we should replace $t_i - t_j$ with $t_i - t_j - \hat{d}_{ij}$ where \hat{d}_{ij} is the expected drift between peak i of the first sample and peak j of the second sample.

More details coming soon.

3.2 Imputing location of undetected peaks

One goal of the alignment leading into downstream data analyses is the generation of a table of abundances for each metabolite across all samples. As you can see from the TICs above, there are some low intensity peaks that fall below detection in some but not all samples. Our view is that instead of inserting arbitrary low constants (such as 0 or half the detection limit) or imputing the intensities post-hoc or having missing data in the data matrices, it is best to return to the area of the where the peak should be and give some kind of abundance. The alignments themselves are rich in information with respect to the location of undetected peaks. We feel this is a more conservative and statistically valid approach than introducing arbitrary values.

More details coming soon.

4 Multiple alignment of several experimental groups

Next, we discuss the multiple alignment of peaks across many samples. With replicates, we typically do an alignment within replicates, then combine these together into a summarized form. This cuts down on the computational cost. For example, consider 2 sets of samples, each with 5 replicates. Aligning first within replicates requires 10+10+1 total alignments whereas an all-pairwise alignment requires 45 pairwise alignments. In addition, this allows some

flexibility in setting different gap and distance penalty parameters for the *within* alignment and *between* alignment. An example follows.

```
> print(targets)

      FileName Group
1 0709_468.CDF  mmA
2 0709_474.CDF  mmA
3 0709_475.CDF  mmA
4 0709_485.CDF  mmC
5 0709_493.CDF  mmC
6 0709_496.CDF  mmC
7 0709_470.CDF  mmD
8 0709_471.CDF  mmD
9 0709_479.CDF  mmD

> ma <- multipleAlignment(pd, group=targets$Group, wn.gap=0.5, wn.D=.05,
+                          bw.gap=.6, bw.D=0.05, usePeaks=TRUE, filterMin=1,
+                          df=50, verbose=FALSE, metric = 1, type = 1) # bug

[clusterAlignment] Aligning 0709_468 to 0709_474
[peaksAlignment] Comparing 24 peaks to 23 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 22 matched. Similarity= 0.2816753
[clusterAlignment] Aligning 0709_468 to 0709_475
[peaksAlignment] Comparing 24 peaks to 26 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 20 matched. Similarity= 0.1618763
[clusterAlignment] Aligning 0709_474 to 0709_475
[peaksAlignment] Comparing 23 peaks to 26 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 20 matched. Similarity= 0.1831043
[progressiveAlignment] Doing merge -1 -3
[progressiveAlignment] left.runs: 1 , right.runs: 3
[progressiveAlignment] Doing merge -2 1
[progressiveAlignment] left.runs: 2 , right.runs: 1 3
[progressiveAlignment] (dot=50) going to 23 :
      used (Mb) gc trigger (Mb) max used (Mb)
Ncells 8357279 446.4 12530440 669.2 12530440 669.2
Vcells 15824157 120.8 38264303 292.0 38264303 292.0
[clusterAlignment] Aligning 0709_485 to 0709_493
[peaksAlignment] Comparing 20 peaks to 27 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 20 matched. Similarity= 0.2221932
[clusterAlignment] Aligning 0709_485 to 0709_496
[peaksAlignment] Comparing 20 peaks to 24 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 18 matched. Similarity= 0.2047329
[clusterAlignment] Aligning 0709_493 to 0709_496
[peaksAlignment] Comparing 27 peaks to 24 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 22 matched. Similarity= 0.2446219
[progressiveAlignment] Doing merge -4 -6
[progressiveAlignment] left.runs: 4 , right.runs: 6
[progressiveAlignment] Doing merge -5 1
[progressiveAlignment] left.runs: 5 , right.runs: 4 6
[progressiveAlignment] (dot=50) going to 27 :
      used (Mb) gc trigger (Mb) max used (Mb)
```

```

Ncells 8357377 446.4 12530440 669.2 12530440 669.2
Vcells 15827059 120.8 38264303 292.0 38264303 292.0
[clusterAlignment] Aligning 0709_470 to 0709_471
[peaksAlignment] Comparing 24 peaks to 25 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 23 matched. Similarity= 0.2891437
[clusterAlignment] Aligning 0709_470 to 0709_479
[peaksAlignment] Comparing 24 peaks to 21 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 18 matched. Similarity= 0.1674945
[clusterAlignment] Aligning 0709_471 to 0709_479
[peaksAlignment] Comparing 25 peaks to 21 peaks -- gap= 0.5 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 20 matched. Similarity= 0.165173
[progressiveAlignment] Doing merge -8 -9
[progressiveAlignment] left.runs: 8 , right.runs: 9
[progressiveAlignment] Doing merge -7 1
[progressiveAlignment] left.runs: 7 , right.runs: 8 9
[progressiveAlignment] (dot=50) going to 24 :
      used (Mb) gc trigger (Mb) max used (Mb)
Ncells 8357514 446.4 12530440 669.2 12530440 669.2
Vcells 15829808 120.8 38264303 292.0 38264303 292.0
[clusterAlignment] Aligning to
[peaksAlignment] Comparing 31 peaks to 29 peaks -- gap= 0.6 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 26 matched. Similarity= 0.3217042
[clusterAlignment] Aligning to
[peaksAlignment] Comparing 31 peaks to 28 peaks -- gap= 0.6 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 26 matched. Similarity= 0.2240976
[clusterAlignment] Aligning to
[peaksAlignment] Comparing 29 peaks to 28 peaks -- gap= 0.6 D= 0.05 , metric= 1 , type= 1
[peaksAlignment] 26 matched. Similarity= 0.2825078
[progressiveAlignment] Doing merge -1 -3
[progressiveAlignment] left.runs: 1 , right.runs: 3
[progressiveAlignment] Doing merge -2 1
[progressiveAlignment] left.runs: 2 , right.runs: 1 3
[progressiveAlignment] (dot=50) going to 29 :
      used (Mb) gc trigger (Mb) max used (Mb)
Ncells 8357755 446.4 12530440 669.2 12530440 669.2
Vcells 15878598 121.2 38264303 292.0 38264303 292.0

```

```
> ma
```

```

An object of class "multipleAlignment"
3 groups: 3 3 3 samples, respectively.
35 merged peaks

```

If you set `verbose=TRUE`, many nitty-gritty details of the alignment procedure are given. Next, we can take the alignment results and overlay it onto the TICs, allowing a visual inspection.

```

> plotChrom(pd, rtrange=c(7.5,8.5), runs=ma@betweenAlignment@runs,
+   mind=ma@betweenAlignment@ind, plotPeaks=TRUE,
+   plotPeakLabels=TRUE, max.near=8, how.near=.5,
+   col=rep(c("blue","red","black"), each=3))
> mp <- correlationAlignment(object=pd.2, thr=0.85, D=20, penalty=0.2,
+   normalize=TRUE, minFilter=1)
> mp

```

4.1 Gathering results

The alignment results can be extracted from the `multipleAlignment` object as:

```
> ma@betweenAlignment@runs
```

```
[1] 5 4 6 2 1 3 7 8 9
```

```
> ma@betweenAlignment@ind
```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]
[1,]	1	1	1	1	1	1	1	1	1
[2,]	2	NA	NA	2	2	2	2	2	NA
[3,]	3	2	2	3	3	3	3	3	2
[4,]	4	3	3	4	4	4	4	4	3
[5,]	5	NA	4	NA	5	NA	5	5	NA
[6,]	6	4	5	5	6	5	6	6	4
[7,]	NA	NA	NA	6	7	NA	NA	NA	NA
[8,]	7	5	6	7	8	6	7	7	5
[9,]	8	6	NA	8	9	7	NA	NA	NA
[10,]	9	7	7	9	10	8	NA	8	6
[11,]	NA	NA	NA	NA	NA	NA	8	NA	7
[12,]	10	8	8	10	11	NA	9	9	8
[13,]	11	NA	NA	NA	NA	NA	NA	NA	NA
[14,]	12	NA	NA	11	12	9	NA	10	9
[15,]	13	9	9	NA	13	NA	NA	11	10
[16,]	14	10	NA	12	14	10	10	12	11
[17,]	15	11	10	13	15	11	11	13	12
[18,]	16	NA	11	NA	16	12	12	14	NA
[19,]	17	NA	12	NA	NA	NA	NA	NA	NA
[20,]	18	12	13	14	NA	13	13	15	13
[21,]	19	NA	14	NA	NA	14	14	16	14
[22,]	20	13	15	15	17	15	15	17	15
[23,]	NA	NA	NA	16	NA	NA	16	NA	NA
[24,]	21	14	16	17	18	16	17	18	16
[25,]	NA	15	17	18	19	17	18	19	17
[26,]	22	NA	18	NA	NA	18	NA	NA	NA
[27,]	23	16	19	19	20	19	19	20	18
[28,]	NA	NA	20	NA	NA	NA	NA	21	19
[29,]	24	17	21	20	NA	20	20	22	NA
[30,]	NA	NA	NA	NA	21	21	21	NA	NA
[31,]	25	18	22	21	22	22	22	23	NA
[32,]	NA	NA	NA	NA	NA	23	NA	NA	NA
[33,]	26	19	23	22	23	24	23	24	20
[34,]	NA	NA	NA	NA	NA	25	NA	NA	NA
[35,]	27	20	24	23	24	26	24	25	21

This table would suggest that matched peak 8 (see numbers below the TICs in the figure above) corresponds to detected peaks 9, 12, 11 in runs 4, 5, 6 and so on, same as shown in the above plot.

In addition, you can gather a list of all the merged peaks with the `gatherInfo` function, giving elements for the retention times, the detected fragment ions and their intensities. The example below also shows the how to construct a table of retention times of the matched peaks (No attempt is made here to adjust retention times onto a common scale. Instead, the peaks are matched to each other on their original scale). For example:


```

> outList <- gatherInfo(pd,ma)
> outList[[8]]

$rt
      mmC.5      mmC.4      mmC.6      mmA.2      mmA.1      mmA.3      mmD.7      mmD.8
7.694983 7.716883 7.694267 7.713550 7.708567 7.648733 7.702967 7.701717
      mmD.9
7.713933

$mz
[1] 58 59 66 72 73 74 75 79 89 104 116 133 147 148 188 204

$data
      mmC.5 mmC.4 mmC.6 mmA.2 mmA.1 mmA.3 mmD.7 mmD.8 mmD.9
[1,]      0      0 16656      0      0      0 14648      0      0
[2,] 5740 4717 5699 5113 4425      0 4855 4557 4728
[3,] 5633 4824 5475 3926 5146      0 4831 3354 4783
[4,] 4909 3931 5123 3680 3910      0 4051 3427 3907
[5,] 75336 69832 77784 61816 65680      0 65912 52848 61560
[6,] 8067 7367 9000 6705 6185      0 6642 6235 7088
[7,] 29936 27440 29256 26376 23328      0 28016 26304 27184
[8,]      0      0      0      0      0      0 38712      0      0
[9,] 6213 5886 7266 5617 5347      0 5173 3946 5659
[10,] 15490 14203 17408 13173 13808      0 12852 9816 12492
[11,] 6447 5440 7809 5417 5525      0 5577 4504 5201
[12,] 3583 3334 4461 3539 3730      0 3599 2436 3893
[13,] 24600 24296 28672 21864 20016      0 21904 17896 22400
[14,] 4473 4299 4918 4413 3430      0 3006 3335 3851
[15,] 17208 16208 18224 14433 14751      0 14731 10680 14061
[16,] 8053 7692 8977 6878 6667      0 6935 5149 6830

> rtmat <- matrix(unlist(lapply(outList,.subset,"rt"), use.names=FALSE),
+                 nr=length(outList), byrow=TRUE)
> colnames(rtmat) <- names(outList[[1]]$rt); rownames(rtmat) <- 1:nrow(rtmat)
> round(rtmat, 3)

      mmC.5 mmC.4 mmC.6 mmA.2 mmA.1 mmA.3 mmD.7 mmD.8 mmD.9
1 7.512 7.534 7.506 7.531 7.526 7.540 7.520 7.519 7.531
2 7.535 NA NA 7.559 7.549 7.557 7.543 7.547 NA
3 7.558 7.580 7.551 7.576 7.566 7.574 7.560 7.565 7.577
4 7.575 7.597 7.569 7.588 7.583 7.592 7.577 7.582 7.594
5 7.586 NA 7.586 NA 7.600 NA 7.594 7.599 NA
6 7.615 7.614 7.614 7.616 7.617 7.614 7.617 7.610 7.617
7 NA NA NA 7.691 7.663 NA NA NA NA
8 7.695 7.717 7.694 7.714 7.709 7.649 7.703 7.702 7.714
9 7.741 7.728 NA 7.736 7.783 7.712 NA NA NA
10 7.804 7.803 7.711 7.799 7.800 7.803 NA 7.805 7.805
11 NA NA NA NA NA NA 7.806 NA 7.817
12 7.809 7.825 7.803 7.828 7.823 NA 7.812 7.816 7.823
13 7.849 NA NA NA NA NA NA NA NA NA
14 7.946 NA NA 7.942 7.880 7.826 NA 7.907 7.874
15 7.958 7.946 7.951 NA 7.943 NA NA 7.936 7.943

```

16	7.969	7.974	NA	7.976	7.966	7.975	7.966	7.965	7.977
17	7.986	8.008	7.980	7.999	7.994	7.997	7.989	7.993	8.000
18	8.009	NA	8.003	NA	8.011	8.009	8.012	8.010	NA
19	8.049	NA	8.043	NA	NA	NA	NA	NA	NA
20	8.061	8.077	8.060	8.079	NA	8.077	8.069	8.068	8.080
21	8.107	NA	8.100	NA	NA	8.095	8.086	8.085	8.091
22	8.204	8.111	8.111	8.114	8.109	8.112	8.109	8.108	8.114
23	NA	NA	NA	8.182	NA	NA	8.172	NA	NA
24	8.244	8.254	8.237	8.251	8.246	8.249	8.246	8.245	8.251
25	NA	8.266	8.254	8.285	8.263	8.283	8.280	8.262	8.263
26	8.301	NA	8.294	NA	NA	8.312	NA	NA	NA
27	8.324	8.334	8.323	8.337	8.332	8.335	8.326	8.330	8.337
28	NA	NA	8.329	NA	NA	NA	NA	8.342	8.400
29	8.352	8.363	8.352	8.359	NA	8.357	8.360	8.359	NA
30	NA	NA	NA	NA	8.360	8.375	8.377	NA	NA
31	8.392	8.403	8.386	8.399	8.394	8.403	8.395	8.393	NA
32	NA	NA	NA	NA	NA	8.420	NA	NA	NA
33	8.432	8.437	8.432	8.434	8.434	8.437	8.435	8.433	8.440
34	NA	NA	NA	NA	NA	8.443	NA	NA	NA
35	8.461	8.477	8.460	8.474	8.469	8.472	8.469	8.468	8.474

5 Future improvements and extension

There are many procedures that we have implemented in our investigation of GC-MS data, but have not made part of the package just yet. Some of the most useful procedures will be released, such as:

1. Parsers for other peak detection algorithms (e.g. MzMine) and parsers for other alignment procedures (e.g. SpectConnect) and perhaps retention indices procedures.
2. More convenient access to the alignment information and abundance table.
3. Statistical analysis of differential metabolite abundance.
4. Fragment-level analysis, an alternative method to summarize abundance across all detected fragments of a metabolite peak.

6 References

See the following for further details:

1. Robinson MD. *Methods for the analysis of gas chromatography - mass spectrometry data*. **Ph.D. Thesis**. October 2008. Department of Medical Biology (Walter and Eliza Hall Institute of Medical Research), University of Melbourne.
2. Robinson MD, De Souza DP, Keen WW, Saunders EC, McConville MJ, Speed TP, Likić VA. (2007) *A dynamic programming approach for the alignment of signal peaks in multiple gas chromatography-mass spectrometry experiments*. **BMC Bioinformatics**. 8:419.
3. Prince JT, Marcotte EM (2006) *Chromatographic alignment of ESI-LC-MS proteomics data sets by ordered bijective interpolated warping*. **Anal Chem**. 78(17):6140-52.

7 This vignette was built with/at ...

```
> sessionInfo()
```

```
R version 4.3.0 RC (2023-04-13 r84269 ucrt)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows Server 2022 x64 (build 20348)
```

```
Matrix products: default
```

```
locale:
```

```
[1] LC_COLLATE=C
[2] LC_CTYPE=English_United States.utf8
[3] LC_MONETARY=English_United States.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.utf8
```

```
time zone: America/New_York
```

```
tzcode source: internal
```

```
attached base packages:
```

```
[1] stats4      stats      graphics  grDevices  utils      datasets  methods
[8] base
```

```
other attached packages:
```

```
[1] flagme_1.56.0      CAMERA_1.56.0      xcms_3.22.0
[4] MSnbase_2.26.0     ProtGenerics_1.32.0 S4Vectors_0.38.0
[7] mzR_2.34.0         Rcpp_1.0.10        Biobase_2.60.0
[10] BiocGenerics_0.46.0 BiocParallel_1.34.0 gcspikelite_1.37.0
```

```
loaded via a namespace (and not attached):
```

```
[1] bitops_1.0-7          RBGL_1.76.0
[3] gridExtra_2.3         rlang_1.1.0
[5] magrittr_2.0.3        clue_0.3-64
[7] MassSpecWavelet_1.66.0 matrixStats_0.63.0
[9] compiler_4.3.0        vctrs_0.6.2
[11] stringr_1.5.0         pkgconfig_2.0.3
[13] fastmap_1.1.1         backports_1.4.1
[15] XVector_0.40.0        caTools_1.18.2
[17] utf8_1.2.3            rmarkdown_2.21
[19] graph_1.78.0          preprocessCore_1.62.0
[21] xfun_0.39             zlibbioc_1.46.0
[23] GenomeInfoDb_1.36.0   DelayedArray_0.26.0
[25] parallel_4.3.0        cluster_2.1.4
[27] R6_2.5.1              stringi_1.7.12
[29] RColorBrewer_1.1-3    limma_3.56.0
[31] rpart_4.1.19          GenomicRanges_1.52.0
[33] SummarizedExperiment_1.30.0 iterators_1.0.14
[35] knitr_1.42            snow_0.4-4
[37] base64enc_0.1-3       IRanges_2.34.0
```

[39] igraph_1.4.2	Matrix_1.5-4
[41] splines_4.3.0	nnet_7.3-18
[43] tidyselect_1.2.0	rstudioapi_0.14
[45] gplots_3.1.3	doParallel_1.0.17
[47] codetools_0.2-19	affy_1.78.0
[49] lattice_0.21-8	tibble_3.2.1
[51] plyr_1.8.8	evaluate_0.20
[53] foreign_0.8-84	survival_3.5-5
[55] pillar_1.9.0	affyio_1.70.0
[57] BiocManager_1.30.20	MatrixGenerics_1.12.0
[59] KernSmooth_2.23-20	checkmate_2.1.0
[61] foreach_1.5.2	MALDIquant_1.22.1
[63] ncd4_1.21	generics_0.1.3
[65] RCurl_1.98-1.12	ggplot2_3.4.2
[67] munsell_0.5.0	scales_1.2.1
[69] gtools_3.9.4	glue_1.6.2
[71] MsFeatures_1.8.0	Hmisc_5.0-1
[73] tools_4.3.0	data.table_1.14.8
[75] mzID_1.38.0	robustbase_0.95-1
[77] SparseM_1.81	vsn_3.68.0
[79] RANN_2.6.1	XML_3.99-0.14
[81] grid_4.3.0	impute_1.74.0
[83] MsCoreUtils_1.12.0	colorspace_2.1-0
[85] GenomeInfoDbData_1.2.10	htmlTable_2.4.1
[87] Formula_1.2-5	cli_3.6.1
[89] fansi_1.0.4	dplyr_1.1.2
[91] pcaMethods_1.92.0	gtable_0.3.3
[93] DEoptimR_1.0-12	digest_0.6.31
[95] htmlwidgets_1.6.2	htmltools_0.5.5
[97] multtest_2.56.0	lifecycle_1.0.3
[99] MASS_7.3-59	

> date()

[1] "Tue Apr 25 22:14:33 2023"