

Supplementary data for: ‘MSnbase’, efficient and elegant R-based processing and visualisation of raw mass spectrometry data

Laurent Gatto, Sebastian Gibb and Johannes Rainer

25 June 2020

Contents

1	Introduction	2
2	Data handling and analysis with MSnbase.	2
2.1	Performance of the on-disk backend on large scale data sets	3
2.2	System information	5

1 Introduction

This document describes handling of mass spectrometry data from large experiments using the `MSnbase` package and more specifically its *on-disk* backend. For demonstration purposes, the `MassIVE` data set `MSV000080030` is used. This consists of over 1,000 mzXML files from swab-samples collected from hands and various personal objects of 80 volunteers.

2 Data handling and analysis with `MSnbase`

In this section we demonstrate data handling and access by `MSnbase` on a large experiment consisting of more than 1,000 data files.

To reproduce the analysis in this document, download the `MSV000080030` folder from <ftp://massive.ucsd.edu/MSV000080030/> and place it into the same folder than this document.

Below we load the required libraries and define the files to be analyzed.

```
library(MSnbase)
library(magrittr)
library(pryr)

fls <- dir("MSV000080030/ccms_peak/Forensic_study_80_volunteers/",
          pattern = "mzXML", full.names = TRUE, recursive = TRUE)
```

The data set consists of in total 1182 mzXML files. We next load the data using the two different `MSnbase` backends `"inMemory"` and `"onDisk"`. For the in-memory backend, due to the larger memory requirements, we import the data only from a subset of the files.

```
ms_mem <- readMSData(fls[grepl("Hand", fls)], mode = "inMemory")
```

Next we load data from all mzXML files as an on-disk `MSnExp` object.

```
ms_dsk <- readMSData(fls, mode = "onDisk")
```

Below we count the number of spectra per MS level of the whole experiment.

```
table(msLevel(ms_dsk))
##
##      1      2
## 1173678 4599786
```

Note that the in-memory `MSnExp` object contains only MS2 spectra (in total 2140520) from a subset of data files, still, data import was much slower (over ~ 12 hours for the in-memory backend while creating the on-disk object from the full data data set took ~ 3 hours).

Next we subset the on-disk object to contain the same set of spectra than the in-memory `MSnExp` and compare their memory footprint.

```
ms_dsk_hands <- ms_dsk %>%
  filterFile(grep("Hand", fls)) %>%
  filterMsLevel(2L)

object_size(ms_mem)
```

Supplementary data for: 'MSnbase', efficient and elegant R-based processing and visualisation of raw mass spectrometry data

```
## 21.8 GB
object_size(ms_dsk_hands)
## 617 MB
```

Since the on-disk object stores only spectra metadata in memory it occupies also much less system memory. As a comparison, the on-disk `MSnExp` for the full experiment was still much smaller than the in-memory object:

```
object_size(ms_dsk)
## 1.66 GB
```

2.1 Performance of the on-disk backend on large scale data sets

To demonstrate `MSnbase`'s efficiency in processing large scale experiments we perform some standard subsetting, data access and manipulation operations.

We first compare the performance of the on-disk and in-memory backend on accessing m/z values with the `mz` function on a set of 100 randomly selected spectra. The performance is assessed with the `microbenchmark` function.

```
set.seed(123)
idx <- sample(seq_along(ms_mem), 100)

library(microbenchmark)
microbenchmark(mz(ms_mem[idx]),
               mz(ms_dsk_hands[idx]),
               times = 5)

## Unit: seconds
##           expr      min       lq     mean   median      uq      max
##  mz(ms_mem[idx]) 51.514084 55.542641 59.32276 55.586931 60.34142 73.62871
##  mz(ms_dsk_hands[idx]) 3.892099 3.893295 13.41512 3.976054 26.77394 28.54019
## neval
##      5
##      5
```

Thus, for this combined subsetting and data access operation the on-disk backend performed better than the in-memory `MSnExp`, while even requiring much less memory.

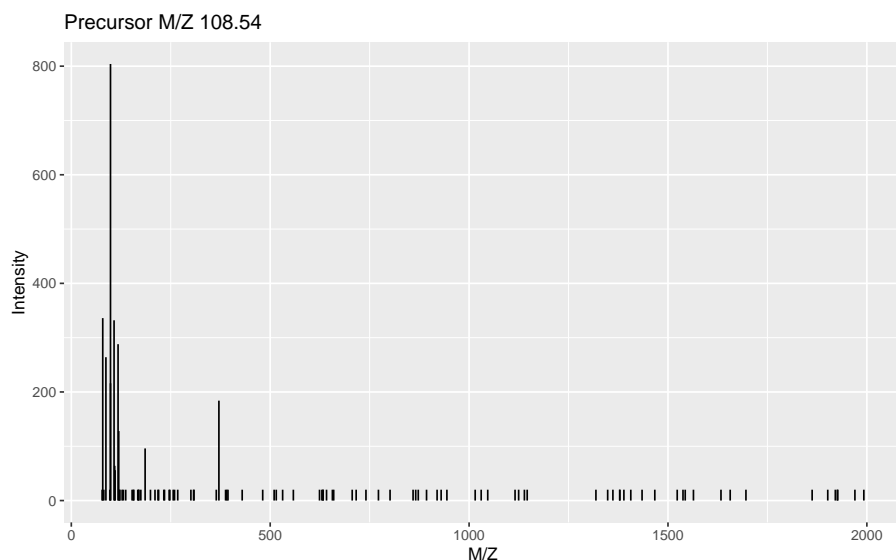
Next we extract all MS2 spectra with a retention time between 50 and 60 seconds and a precursor m/z of 108.5362 (± 5 ppm). This subsetting operation is performed on the on-disk `MSnExp` object representing the full experiment with the 1182 data files/samples. To assess the performance of the following operations we use `system.time` calls that record elapsed time in seconds.

```
system.time(
  ms_sub <- ms_dsk %>%
    filterMsLevel(2L) %>%
    filterRt(c(50, 60)) %>%
    filterPrecursorMz(mz = 108.5362, ppm = 5)
)["elapsed"]
## elapsed
## 6.698
```

Supplementary data for: 'MSnbase', efficient and elegant R-based processing and visualisation of raw mass spectrometry data

In total `length(ms_sub)` spectra were selected from in total 928 data files/samples. The plot below shows the data for the first spectrum.

```
system.time(  
  plot(ms_sub[[1]])  
)["elapsed"]
```



```
## elapsed  
## 0.354
```

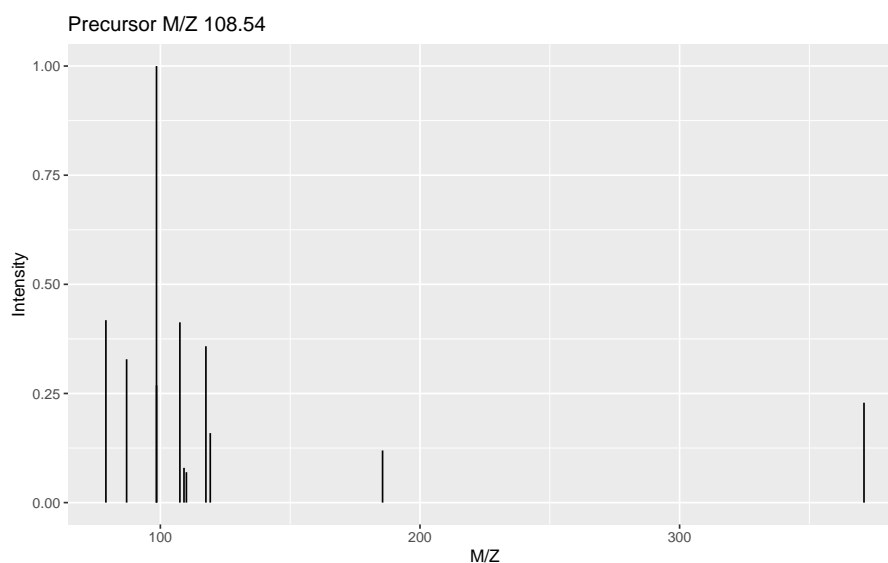
Since there seems to be quite some background noise in the MS2 spectrum we next remove peaks with an intensity below 50 by first replacing their intensities with 0 (with the `removePeaks` call) and subsequently removing all 0-intensity peaks from each spectrum with the `clean` call. In addition we *normalize* each spectrum by dividing the maximum intensity per spectrum from the spectrum's intensities.

```
system.time(  
  ms_sub <- ms_sub %>%  
    removePeaks(t = 50) %>%  
    clean(all = TRUE) %>%  
    normalize(method = "max")  
)["elapsed"]  
## elapsed  
## 0.043
```

The result on the first spectrum is shown below.

```
system.time(  
  plot(ms_sub[[1]])  
)["elapsed"]
```

Supplementary data for: 'MSnbase', efficient and elegant R-based processing and visualisation of raw mass spectrometry data



Supplementary data for: 'MSnbase', efficient and elegant R-based processing and visualisation of raw mass spectrometry data

```
## [8] methods    base
##
## other attached packages:
## [1] microbenchmark_1.4-7 BiocParallel_1.22.0 pryr_0.1.4
## [4] magrittr_1.5         MSnbase_2.14.2      ProtGenerics_1.20.0
## [7] S4Vectors_0.26.1     mzR_2.22.0         Rcpp_1.0.4.6
## [10] Biobase_2.48.0       BiocGenerics_0.34.0 BiocStyle_2.16.0
## [13] rmarkdown_2.3
##
## loaded via a namespace (and not attached):
## [1] tinytex_0.24          tidyselect_1.1.0      xfun_0.15
## [4] purrr_0.3.4           lattice_0.20-41       colorspace_1.4-1
## [7] vctrs_0.3.1           generics_0.0.2        htmltools_0.5.0
## [10] yaml_2.2.1            vsn_3.56.0            XML_3.99-0.3
## [13] rlang_0.4.6           pillar_1.4.4          glue_1.4.1
## [16] affy_1.66.0           foreach_1.5.0         affyio_1.58.0
## [19] lifecycle_0.2.0       plyr_1.8.6            mzID_1.26.0
## [22] stringr_1.4.0         zlibbioc_1.34.0       munsell_0.5.0
## [25] pcaMethods_1.80.0     gtable_0.3.0          codetools_0.2-16
## [28] evaluate_0.14         labeling_0.3           knitr_1.29
## [31] IRanges_2.22.2        doParallel_1.0.15     preprocessCore_1.50.0
## [34] scales_1.1.1          BiocManager_1.30.10   limma_3.44.3
## [37] farver_2.0.3          impute_1.62.0         ggplot2_3.3.2
## [40] digest_0.6.25         stringi_1.4.6         bookdown_0.20
## [43] dplyr_1.0.0           ncd4_1.17             grid_4.0.0
## [46] tools_4.0.0           tibble_3.0.1          crayon_1.3.4
## [49] pkgconfig_2.0.3       ellipsis_0.3.1        MASS_7.3-51.6
## [52] iterators_1.0.12      R6_2.4.1              MALDIquant_1.19.3
## [55] compiler_4.0.0
```