

Handling metadata and annotations

AlpsNMR authors

2024-04-17

Abstract

This vignette shows some examples on how to explore sample metadata and add additional sample annotations, coming from one or more CSV or Excel files.

Package

AlpsNMR 4.5.0

Contents

1	Getting started	2
2	Exploring the sample metadata	2
3	Sample annotations	5
4	Further annotations	6
5	Summary	8
6	Session Information	8

1 Getting started

We start by loading AlpsNMR and some convenience libraries:

```
library(dplyr)
library(readxl)
library(AlpsNMR)
```

We also load the demo samples, see the introduction vignette for further details:

```
MeOH_plasma_extraction_dir <- system.file("dataset-demo", package = "AlpsNMR")
zip_files <- list.files(MeOH_plasma_extraction_dir, pattern = glob2rx("*.zip"), full.names = TRUE)
dataset <- nmr_read_samples(sample_names = zip_files)
dataset <- nmr_interpolate_1D(dataset, axis = NULL)
dataset
## An nmr_dataset_1D (3 samples)
```

```
plot(dataset, chemshift_range = c(3.4, 3.6))
```

2 Exploring the sample metadata

Most NMR formats include besides the actual NMR spectra, a lot of additional information describing the acquisition properties, instrument settings, and spectral processing information.

AlpsNMR parses all that information whenever possible, and stores it in the `nmr_dataset` object, so the user can inspect it. Since there may be a lot of information, the data is stored in several data frames.

The available data frames are:

```
nmr_meta_groups(dataset)
## [1] "info"      "orig"      "title"     "acqu"     "procs"     "levels"    "external"
```

We can further explore each of those groups.

For instance, for the `acqu` group we find 239 columns:

```
acqu_metadata <- nmr_meta_get(dataset, groups = "acqu")
acqu_metadata
## # A tibble: 3 x 239
##   NMRExperiment acqu_TITLE          acqu_JCAMPDX acqu_DATATYPE acqu_NPOINTS
##   <chr>          <chr>                <dbl> <chr>          <chr>
## 1 10            Parameter file, TopS~ 5 Parameter Val~ "13\t$$ modi~
## 2 20            Parameter file, TopS~ 5 Parameter Val~ "15\t$$ modi~
## 3 30            Parameter file, TopS~ 5 Parameter Val~ "13\t$$ modi~
## # i 234 more variables: acqu_ORIGIN <chr>, acqu_OWNER <chr>,
## #   acqu_Stamp <list>, acqu_ACQT0 <dbl>, acqu_AMP <list>,
## #   acqu_AMP_COIL <list>, acqu_ANAVPT <dbl>, acqu_AQSEQ <dbl>,
## #   acqu_AQ_mod <dbl>, acqu_AUNM <chr>, acqu_AUTOPOS <chr>, acqu_BF1 <dbl>,
## #   acqu_BF2 <dbl>, acqu_BF3 <dbl>, acqu_BF4 <dbl>, acqu_BF5 <dbl>,
## #   acqu_BF6 <dbl>, acqu_BF7 <dbl>, acqu_BF8 <dbl>, acqu_BWFAC <list>,
## #   acqu_BYTORDA <dbl>, acqu_CAGPARS <list>, acqu_CHEMSTR <chr>, ...
```

Handling metadata and annotations

Here follows a long list of all the columns available:

```
colnames(acqus_metadata)
## [1] "NMRExperiment"      "acqus_TITLE"      "acqus_JCAMPDX"
## [4] "acqus_DATATYPE"     "acqus_NPOINTS"   "acqus_ORIGIN"
## [7] "acqus_OWNER"        "acqus_Stamp"     "acqus_ACQT0"
## [10] "acqus_AMP"          "acqus_AMPACOIL"  "acqus_ANAVPT"
## [13] "acqus_AQSEQ"        "acqus_AQ_mod"    "acqus_AUNM"
## [16] "acqus_AUTOPOS"     "acqus_BF1"       "acqus_BF2"
## [19] "acqus_BF3"          "acqus_BF4"       "acqus_BF5"
## [22] "acqus_BF6"          "acqus_BF7"       "acqus_BF8"
## [25] "acqus_BWFAC"        "acqus_BYTORDA"   "acqus_CAGPARS"
## [28] "acqus_CHEMSTR"      "acqus_CNST"      "acqus_CPDPRG"
## [31] "acqus_D"            "acqus_DATE"      "acqus_DE"
## [34] "acqus_DECBNUC"     "acqus_DECIM"     "acqus_DECNUC"
## [37] "acqus_DECSTAT"     "acqus_DIGMOD"    "acqus_DIGTYP"
## [40] "acqus_DQDMODE"     "acqus_DR"        "acqus_DS"
## [43] "acqus_DSPFIRM"     "acqus_DSPFVS"    "acqus_DTYPA"
## [46] "acqus_EXP"          "acqus_FCUCHAN"   "acqus_FL1"
## [49] "acqus_FL2"          "acqus_FL3"       "acqus_FL4"
## [52] "acqus_FN_INDIRECT" "acqus_FOV"       "acqus_FQ1LIST"
## [55] "acqus_FQ2LIST"     "acqus_FQ3LIST"   "acqus_FQ4LIST"
## [58] "acqus_FQ5LIST"     "acqus_FQ6LIST"   "acqus_FQ7LIST"
## [61] "acqus_FQ8LIST"     "acqus_FRQL03"    "acqus_FRQL03N"
## [64] "acqus_FS"           "acqus_FTLPGN"    "acqus_FW"
## [67] "acqus_FnILOOP"     "acqus_FnMODE"    "acqus_FnTYPE"
## [70] "acqus_GPNAM"        "acqus_GPX"        "acqus_GPY"
## [73] "acqus_GPZ"          "acqus_GRDPROG"   "acqus_GRPDLY"
## [76] "acqus_HDDUTY"      "acqus_HDRATE"    "acqus_HGAIN"
## [79] "acqus_HL1"          "acqus_HL2"       "acqus_HL3"
## [82] "acqus_HL4"          "acqus_HOLDER"    "acqus_HPMOD"
## [85] "acqus_HPPRGN"      "acqus_IN"         "acqus_INF"
## [88] "acqus_INP"          "acqus_INSTRUM"   "acqus_INTEGFAC"
## [91] "acqus_L"            "acqus_LFILTER"   "acqus_LGAIN"
## [94] "acqus_LINPSTP"     "acqus_LOCKED"    "acqus_LOCKFLD"
## [97] "acqus_LOCKGN"      "acqus_LOCKPOW"   "acqus_LOCKPPM"
## [100] "acqus_LOCNUC"      "acqus_LOCPHAS"   "acqus_LOCSHFT"
## [103] "acqus_LOCSW"       "acqus_LTIME"     "acqus_MASR"
## [106] "acqus_MASRLST"    "acqus_MULEXPNO"  "acqus_NBL"
## [109] "acqus_NC"          "acqus_NLOGCH"    "acqus_NOVFLW"
## [112] "acqus_NS"          "acqus_NUC1"      "acqus_NUC2"
## [115] "acqus_NUC3"        "acqus_NUC4"      "acqus_NUC5"
## [118] "acqus_NUC6"        "acqus_NUC7"      "acqus_NUC8"
## [121] "acqus_NUCLEUS"    "acqus_NUSLIST"   "acqus_NusAMOUNT"
## [124] "acqus_NusFPNZ"     "acqus_NusJSP"    "acqus_NusSEED"
## [127] "acqus_NusSPTYPE"   "acqus_NusT2"     "acqus_NusTD"
## [130] "acqus_01"          "acqus_02"        "acqus_03"
## [133] "acqus_04"          "acqus_05"        "acqus_06"
## [136] "acqus_07"          "acqus_08"        "acqus_OVERFLW"
## [139] "acqus_P"           "acqus_PACOIL"    "acqus_PAPS"
## [142] "acqus_PARMODE"     "acqus_PCPD"      "acqus_PEXSEL"
## [145] "acqus_PHCOR"       "acqus_PHLIST"    "acqus_PHP"
```

Handling metadata and annotations

```
## [148] "acqu_ PH_ref"      "acqu_ PL"      "acqu_ PLSTEP"
## [151] "acqu_ PLSTRT"     "acqu_ PLW"     "acqu_ PLWMAX"
## [154] "acqu_ PQPHASE"    "acqu_ PQSCALE" "acqu_ PR"
## [157] "acqu_ PRECHAN"    "acqu_ PRGAIN"  "acqu_ PROBHD"
## [160] "acqu_ PULPROG"    "acqu_ PW"      "acqu_ PYNM"
## [163] "acqu_ ProjAngle"  "acqu_ QNP"     "acqu_ RD"
## [166] "acqu_ RECCHAN"    "acqu_ RECPH"   "acqu_ RECPRE"
## [169] "acqu_ RECPRFX"    "acqu_ RECSEL"  "acqu_ RG"
## [172] "acqu_ R0"         "acqu_ RSEL"    "acqu_ S"
## [175] "acqu_ SELREC"     "acqu_ SF01"    "acqu_ SF02"
## [178] "acqu_ SF03"       "acqu_ SF04"    "acqu_ SF05"
## [181] "acqu_ SF06"       "acqu_ SF07"    "acqu_ SF08"
## [184] "acqu_ SOLVENT"    "acqu_ SOLVOLD" "acqu_ SP"
## [187] "acqu_ SPECTR"     "acqu_ SPINCNT" "acqu_ SPNAM"
## [190] "acqu_ SPOAL"      "acqu_ SPOFFS"  "acqu_ SPPEX"
## [193] "acqu_ SPW"        "acqu_ SUBNAM"  "acqu_ SW"
## [196] "acqu_ SWIBOX"     "acqu_ SW_h"    "acqu_ SWfinal"
## [199] "acqu_ SigLockShift" "acqu_ TD"      "acqu_ TD0"
## [202] "acqu_ TD_INDIRECT" "acqu_ TDav"    "acqu_ TE"
## [205] "acqu_ TE1"        "acqu_ TE2"     "acqu_ TE3"
## [208] "acqu_ TE4"        "acqu_ TEG"     "acqu_ TE_MAGNET"
## [211] "acqu_ TE_PIDX"    "acqu_ TE_STAB" "acqu_ TL"
## [214] "acqu_ TOTROT"     "acqu_ TUBE_TYPE" "acqu_ USERA1"
## [217] "acqu_ USERA2"     "acqu_ USERA3"  "acqu_ USERA4"
## [220] "acqu_ USERA5"     "acqu_ V9"      "acqu_ VALIDCODE"
## [223] "acqu_ VALIST"     "acqu_ VCLIST"  "acqu_ VDLIST"
## [226] "acqu_ VPLIST"     "acqu_ VTLIST"  "acqu_ WBST"
## [229] "acqu_ WBSW"       "acqu_ XGAIN"   "acqu_ XL"
## [232] "acqu_ YL"         "acqu_ YMAX_a"  "acqu_ YMIN_a"
## [235] "acqu_ ZGOPTNS"    "acqu_ ZL1"     "acqu_ ZL2"
## [238] "acqu_ ZL3"        "acqu_ ZL4"
```

We can check for instance that the nuclei used on all samples is 1H:

```
acqu_metadata[, c("NMRExperiment", "acqu_NUC1")]
## # A tibble: 3 x 2
##   NMRExperiment acqu_NUC1
##   <chr>         <chr>
## 1 10             1H
## 2 20             1H
## 3 30             1H
```

Similarly, we can obtain the processing settings:

```
procs_metadata <- nmr_meta_get(dataset, groups = "procs")
procs_metadata
## # A tibble: 3 x 137
##   NMRExperiment procs_TITLE          procs_JCAMPDX procs_DATATYPE procs_NPOINTS
##   <chr>         <chr>                <dbl> <chr>         <chr>
## 1 10             Parameter file, TopS~      5 Parameter Val~ "6\t$$ modif~
## 2 20             Parameter file, TopS~      5 Parameter Val~ "11\t$$ modi~
## 3 30             Parameter file, TopS~      5 Parameter Val~ "6\t$$ modif~
```

Handling metadata and annotations

```
## # i 132 more variables: procs_ORIGIN <chr>, procs_OWNER <chr>,  
## #   procs_Stamp <list>, procs_ABSF1 <dbl>, procs_ABSF2 <dbl>, procs_ABSG <dbl>,  
## #   procs_ABSL <dbl>, procs_ALPHA <dbl>, procs_AQORDER <dbl>,  
## #   procs_ASSFAC <dbl>, procs_ASSFACI <dbl>, procs_ASSFACX <dbl>,  
## #   procs_ASSWID <dbl>, procs_AUNMP <chr>, procs_AXLEFT <dbl>,  
## #   procs_AXNAME <chr>, procs_AXNUC <chr>, procs_AXRIGHT <dbl>,  
## #   procs_AXTYPE <dbl>, procs_AXUNIT <chr>, procs_AZFE <dbl>, ...
```

3 Sample annotations

Besides the sample metadata, most studies usually have design variables or annotations, that describe the biological sample. These annotations do not come from the instrument itself, but rather usually are defined on an *external* CSV or Excel file.

AlpsNMR supports adding *external* annotations from data frames.

Let's load a table from an Excel file, that has some annotations for our demo dataset:

```
excel_file <- file.path(MeOH_plasma_extraction_dir, "dummy_metadata.xlsx")  
subject_timepoint <- read_excel(excel_file, sheet = 1)  
subject_timepoint  
## # A tibble: 3 x 3  
##   NMRExperiment SubjectID TimePoint  
##   <chr>         <chr>    <chr>  
## 1 10           Ana      baseline  
## 2 20           Ana      3 months  
## 3 30           Elia     baseline
```

Note how this table includes a first column named `NMRExperiment`. This column allows us to match the rows in the table with our samples.

We can embed these external annotations in our dataset:

```
dataset <- nmr_meta_add(dataset, metadata = subject_timepoint, by = "NMRExperiment")
```

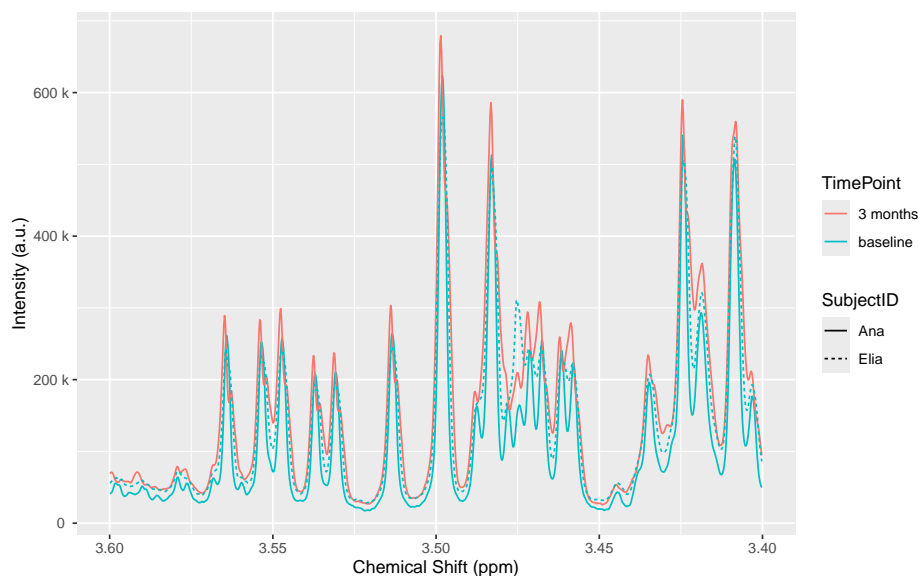
We can retrieve these *external* columns from the dataset:

```
nmr_meta_get(dataset, groups = "external")  
## # A tibble: 3 x 3  
##   NMRExperiment SubjectID TimePoint  
##   <chr>         <chr>    <chr>  
## 1 10           Ana      baseline  
## 2 20           Ana      3 months  
## 3 30           Elia     baseline
```

After adding the annotations to the dataset, we can use them in plots:

```
plot(dataset, color = "TimePoint", linetype = "SubjectID", chemshift_range = c(3.4, 3.6))  
## Warning: ! Passing aes_string arguments to plot(nmr_dataset, ...) is deprecated.  
## i Please pass aes arguments instead  
## This warning is displayed once every 8 hours.
```

Handling metadata and annotations



4 Further annotations

Sometimes due to the study design we have more than one table that we want to match with our data.

For instance, a collaborator just sent us this table:

```
additional_annotations <- data.frame(
  NMRExperiment = c("10", "20", "30"),
  SampleCollectionDay = c(1, 91, 3)
)
additional_annotations
##   NMRExperiment SampleCollectionDay
## 1             10                    1
## 2             20                   91
## 3             30                    3
```

Since we have the NMRExperiment column it is very easy to include it:

```
dataset <- nmr_meta_add(dataset, additional_annotations)
```

And the column has been added:

```
nmr_meta_get(dataset, groups = "external")
## # A tibble: 3 x 4
##   NMRExperiment SubjectID TimePoint SampleCollectionDay
##   <chr>         <chr>    <chr>          <dbl>
## 1 10            Ana      baseline        1
## 2 20            Ana      3 months        91
## 3 30            Elia     baseline        3
```

We received further information, but this time it is related to the SubjectID that we added before:

Handling metadata and annotations

```
subject_related_information <- data.frame(  
  SubjectID = c("Ana", "Elia"),  
  Age = c(33, 3),  
  Sex = c("female", "female")  
)  
subject_related_information  
##   SubjectID Age   Sex  
## 1     Ana  33 female  
## 2     Elia   3 female
```

Note how in this case we only have two rows, and we don't have the `NMRExperiment` column anymore.

We can specify the `by` argument in `nmr_meta_add()` to use another column for merging:

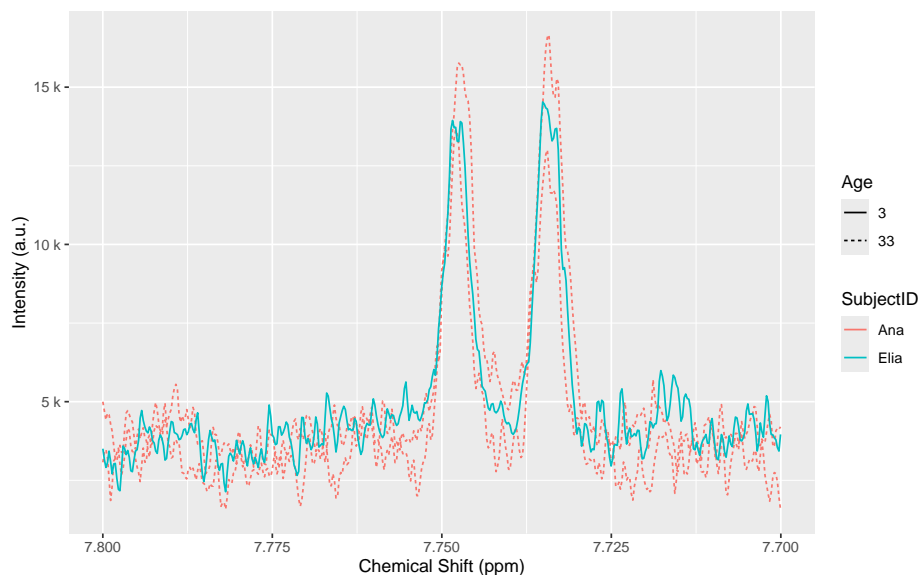
```
dataset <- nmr_meta_add(dataset, subject_related_information, by = "SubjectID")
```

And the `Sex` and `Age` columns will have been added:

```
nmr_meta_get(dataset, groups = "external")  
## # A tibble: 3 x 6  
##   NMRExperiment SubjectID TimePoint SampleCollectionDay   Age Sex  
##   <chr>         <chr>    <chr>          <dbl> <dbl> <chr>  
## 1 10           Ana      baseline         1    33 female  
## 2 20           Ana      3 months        91    33 female  
## 3 30           Elia     baseline         3     3 female
```

We can also use it in a plot:

```
plot(dataset, color = "SubjectID", linetype = "as.factor(Age)", chemshift_range = c(7.7, 7.8)) + ggplot2::lab
```



5 Summary

In this vignette we have seen how to explore the sample metadata, including acquisition and processing settings, and how to embed external annotations and use them in plots.

AlpsNMR is able to merge external annotations as long as there is a common annotation in the data that can be used as merging key.

To import external data, you may want to use the following functions:

File type	Suggested function
CSV	<code>readr::read_csv()</code>
TSV	<code>readr::read_tsv()</code>
SPSS	<code>haven::read_spss()</code>
xls/xlsx	<code>readxl::read_excel()</code>

6 Session Information

```
sessionInfo()
## R version 4.4.0 beta (2024-04-15 r86425)
## Platform: x86_64-pc-linux-gnu
## Running under: Ubuntu 22.04.4 LTS
##
## Matrix products: default
## BLAS: /home/biocbuild/bbs-3.19-bioc/R/lib/libRblas.so
## LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.10.0
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
## [3] LC_TIME=en_GB             LC_COLLATE=C
## [5] LC_MONETARY=en_US.UTF-8   LC_MESSAGES=en_US.UTF-8
## [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
## [9] LC_ADDRESS=C              LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## time zone: America/New_York
## tzcode source: system (glibc)
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] AlpsNMR_4.5.0      future_1.33.2      BiocParallel_1.37.1
## [4] readxl_1.4.3      ggplot2_3.5.0      dplyr_1.1.4
## [7] BiocStyle_2.31.0
##
## loaded via a namespace (and not attached):
## [1] baseline_1.3-5      gridExtra_2.3      rlang_1.1.3
## [4] magrittr_2.0.3      MassSpecWavelet_1.69.0 matrixStats_1.3.0
```


Handling metadata and annotations

```
## [7] compiler_4.4.0      vctrs_0.6.5      reshape2_1.4.4
## [10] RcppZiggurat_0.1.6   quadprog_1.5-8   rvest_1.0.4
## [13] stringr_1.5.1        pkgconfig_2.0.3  crayon_1.5.2
## [16] fastmap_1.1.1        labeling_0.4.3   utf8_1.2.4
## [19] promises_1.3.0       rmarkdown_2.26   ps_1.7.6
## [22] itertools_0.1-3     tinytex_0.50     purrr_1.0.2
## [25] xfun_0.43            Rfast_2.1.0      randomForest_4.7-1.1
## [28] jsonlite_1.8.8       limSolve_1.5.7.1 later_1.3.2
## [31] parallel_4.4.0       cluster_2.1.6    R6_2.5.1
## [34] stringi_1.8.3        RColorBrewer_1.1-3 parallelly_1.37.1
## [37] cellranger_1.1.0     Rcpp_1.0.12      bookdown_0.39
## [40] iterators_1.0.14     knitr_1.46       snow_0.4-4
## [43] Matrix_1.7-0         igraph_2.0.3     tidysselect_1.2.1
## [46] yaml_2.3.8           websocket_1.4.1   codetools_0.2-20
## [49] processx_3.8.4       listenv_0.9.1    doRNG_1.8.6
## [52] lattice_0.22-6       tibble_3.2.1     plyr_1.8.9
## [55] withr_3.0.0          rARPACK_0.11-0   evaluate_0.23
## [58] signal_1.8-0         speaq_2.7.0      RcppParallel_5.1.7
## [61] xml2_1.3.6           lpSolve_5.6.20   pillar_1.9.0
## [64] BiocManager_1.30.22 rngtools_1.5.2   foreach_1.5.2
## [67] ellipse_0.5.0        pcaPP_2.0-4       generics_0.1.3
## [70] chromote_0.2.0       munsell_0.5.1    scales_1.3.0
## [73] globals_0.16.3      glue_1.7.0       tools_4.4.0
## [76] data.table_1.15.4    SparseM_1.81     RSpectra_0.16-1
## [79] fs_1.6.3             mvtnorm_1.2-4    cowplot_1.1.3
## [82] grid_4.4.0           impute_1.77.0    missForest_1.5
## [85] tidyr_1.3.1          colorspace_2.1-0 cli_3.6.2
## [88] fansi_1.0.6          mixOmics_6.27.0  corpcor_1.6.10
## [91] doSNOW_1.0.20        gtable_0.3.4     digest_0.6.35
## [94] progressr_0.14.0     ggrepel_0.9.5    farver_2.1.1
## [97] htmltools_0.5.8.1   lifecycle_1.0.4  httr_1.4.7
## [100] MASS_7.3-60.2
```