

# Handling metadata and annotations

*AlpsNMR authors*

2024-10-29

**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 4.8.0

## 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_AMPCOIL <list>, acqu_ANAVPT <dbl>, acqu_AQSEQ <dbl>,
## #   acqu_AQ_mod <dbl>, acqu_AUNM <chr>, acqu_AUTOP0S <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_AMPCOIL"     "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_ZGPTNS"      "acqu_ZL1"      "acqu_ZL2"
## [238] "acqu_ZL3"         "acqu_ZL4"
```

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

```
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~
```

```
## # 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:

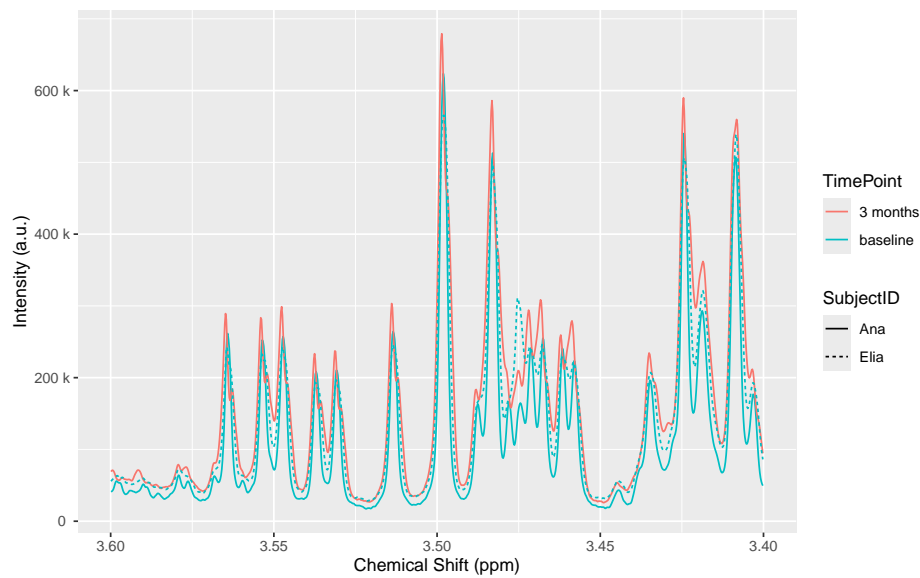
```
dataset <- nmr_meta_add(dataset, metadata = subject_timepoint, by = "NMRExperiment")
```

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.
```



## 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:

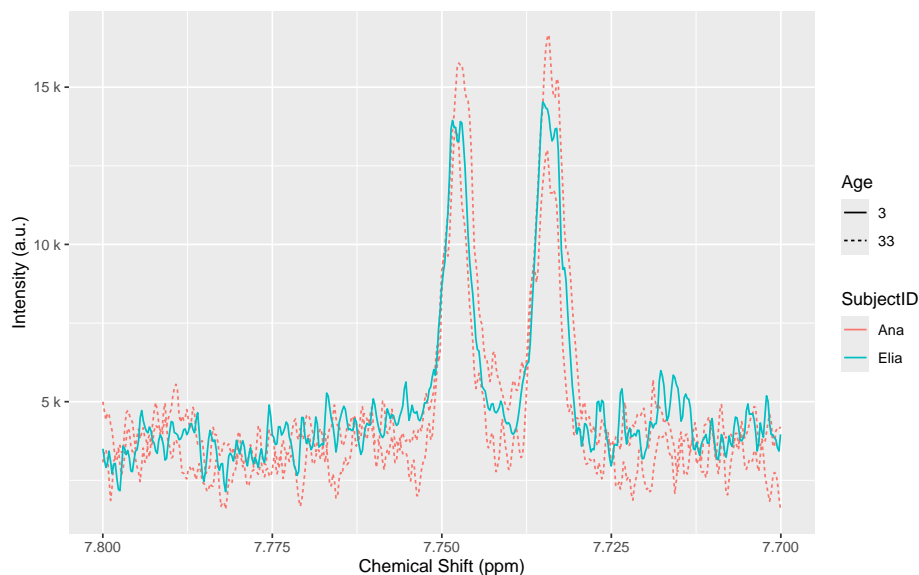
```
dataset <- nmr_meta_add(dataset, subject_related_information, by = "SubjectID")
```

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::labs
```



## 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 4.4.1 (2024-06-14)
## Platform: x86_64-apple-darwin20
## Running under: macOS Monterey 12.7.6
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.4-x86_64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.4-x86_64/Resources/lib/libRlapack.dylib; LAPACK versi
##
## locale:
## [1] C/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: America/New_York
## tzcode source: internal
##
## attached base packages:
## [1] stats graphics grDevices utils datasets methods base
##
## other attached packages:
## [1] AlpsNMR_4.8.0 BiocParallel_1.40.0 readxl_1.4.3
## [4] ggplot2_3.5.1 dplyr_1.1.4 BiocStyle_2.34.0
##
## loaded via a namespace (and not attached):
## [1] baseline_1.3-5 gridExtra_2.3 rlang_1.1.4
## [4] magrittr_2.0.3 MassSpecWavelet_1.72.0 matrixStats_1.4.1
## [7] compiler_4.4.1 vctrs_0.6.5 reshape2_1.4.4
## [10] RcppZigurat_0.1.6 quadprog_1.5-8 rvest_1.0.4
## [13] stringr_1.5.1 pkgconfig_2.0.3 crayon_1.5.3
## [16] fastmap_1.2.0 labeling_0.4.3 utf8_1.2.4
## [19] promises_1.3.0 rmarkdown_2.28 ps_1.8.1
## [22] itertools_0.1-3 tinytex_0.53 purrr_1.0.2
```



## Handling metadata and annotations

## [25] xfun_0.48	Rfast_2.1.0	randomForest_4.7-1.2
## [28] jsonlite_1.8.9	limSolve_1.5.7.1	later_1.3.2
## [31] parallel_4.4.1	cluster_2.1.6	R6_2.5.1
## [34] stringi_1.8.4	RColorBrewer_1.1-3	cellranger_1.1.0
## [37] Rcpp_1.0.13	bookdown_0.41	iterators_1.0.14
## [40] knitr_1.48	snow_0.4-4	Matrix_1.7-1
## [43] igraph_2.1.1	tidyselect_1.2.1	yaml_2.3.10
## [46] websocket_1.4.2	codetools_0.2-20	processx_3.8.4
## [49] doRNG_1.8.6	lattice_0.22-6	tibble_3.2.1
## [52] plyr_1.8.9	withr_3.0.2	rARPACK_0.11-0
## [55] evaluate_1.0.1	signal_1.8-1	speaq_2.7.0
## [58] RcppParallel_5.1.9	xml2_1.3.6	lpSolve_5.6.21
## [61] pillar_1.9.0	BiocManager_1.30.25	rngtools_1.5.2
## [64] foreach_1.5.2	ellipse_0.5.0	pcaPP_2.0-5
## [67] generics_0.1.3	chromote_0.3.1	munSELL_0.5.1
## [70] scales_1.3.0	glue_1.8.0	tools_4.4.1
## [73] data.table_1.16.2	SparseM_1.84-2	RSpectra_0.16-2
## [76] fs_1.6.4	mvtnorm_1.3-1	cowplot_1.1.3
## [79] grid_4.4.1	impute_1.80.0	missForest_1.5
## [82] tidyr_1.3.1	colorspace_2.1-1	cli_3.6.3
## [85] fansi_1.0.6	mixOmics_6.30.0	corpcor_1.6.10
## [88] doSNOW_1.0.20	gtable_0.3.6	digest_0.6.37
## [91] progressr_0.15.0	ggrepel_0.9.6	farver_2.1.2
## [94] htmltools_0.5.8.1	lifecycle_1.0.4	httr_1.4.7
## [97] MASS_7.3-61		