# Package 'MetaboCoreUtils'

March 9, 2025

Title Core Utils for Metabolomics Data

Version 1.15.0

**Description** MetaboCoreUtils defines metabolomics-related core functionality provided as low-level functions to allow a data structure-independent usage across various R packages. This includes functions to calculate between ion (adduct) and compound mass-to-charge ratios and masses or functions to work with chemical formulas. The package provides also a set of adduct definitions and information on some commercially available internal standard mixes commonly used in MS experiments.

**Depends** R (>= 4.0)

Imports utils, MsCoreUtils, BiocParallel, methods, stats

Suggests BiocStyle, testthat, knitr, rmarkdown, robustbase

License Artistic-2.0

LazyData no

VignetteBuilder knitr

BugReports https://github.com/RforMassSpectrometry/MetaboCoreUtils/issues

URL https://github.com/RforMassSpectrometry/MetaboCoreUtils

biocViews Infrastructure, Metabolomics, MassSpectrometry

**Roxygen** list(markdown=TRUE)

RoxygenNote 7.3.1

git\_url https://git.bioconductor.org/packages/MetaboCoreUtils

git\_branch devel

git\_last\_commit 8267815

git\_last\_commit\_date 2024-10-29

Repository Bioconductor 3.21

Date/Publication 2025-03-09

Author Johannes Rainer [aut, cre] (ORCID:

<https://orcid.org/0000-0002-6977-7147>),

Michael Witting [aut] (ORCID: <https://orcid.org/0000-0002-1462-4426>),

30

```
Andrea Vicini [aut],
Liesa Salzer [ctb] (ORCID: <https://orcid.org/0000-0003-0761-0656>),
Sebastian Gibb [aut] (ORCID: <https://orcid.org/0000-0001-7406-4443>),
Michael Stravs [ctb] (ORCID: <https://orcid.org/0000-0002-1426-8572>),
Roger Gine [aut] (ORCID: <https://orcid.org/0000-0003-0288-9619>),
Philippine Louail [aut] (ORCID:
<https://orcid.org/0009-0007-5429-6846>)
```

Maintainer Johannes Rainer <Johannes.Rainer@eurac.edu>

# Contents

addElements	 2
adductFormula	 3
adductNames	 4
calculateKendrickMass	 5
calculateMass	 6
containsElements	 7
convertMtime	 8
correctRindex	 9
countElements	 10
fit_lm	 10
formula2mz	 14
indexRtime	 15
internalStandardMixNames	 16
internalStandards	 16
isotopicSubstitutionMatrix	 17
isotopologues	 19
mass2mz	 21
mclosest	 22
multiplyElements	 23
mz2mass	 24
pasteElements	 25
quality_assessment	 26
standardizeFormula	 28
subtractElements	 29

# Index

addElements

Combine chemical formulae

## Description

addElements Add one chemical formula to another.

# adductFormula

## Usage

addElements(x, y)

#### Arguments

х	character strings with chemical formula
У	character strings with chemical formula that should be added from $\boldsymbol{x}$

# Value

character Resulting formula

# Author(s)

Michael Witting and Sebastian Gibb

## Examples

addElements("C6H12O6", "Na")

addElements("C6H12O6", c("Na", "H2O"))

adductFormula Calculate a table of adduct (ionic) formulas

## Description

adductFormula calculates the chemical formulas for the specified adducts of provided chemical formulas.

# Usage

```
adductFormula(formulas, adduct = "[M+H]+", standardize = TRUE)
```

## Arguments

formulas	character with molecular formulas for which adduct formulas should be cal- culated.
adduct	character or data.frame of valid adduct. to be used. Custom adduct definitions can be provided via a data.frame but its format must follow adducts()
standardize	logical(1) whether to standardize the molecular formulas to the Hill notation system before calculating their mass.

# Value

character matrix with *formula* rows and *adducts* columns containing all ion formulas. In case an ion can't be generated (eg. [M-NH3+H]+ in a molecule that doesn't have nitrogen), a NA is returned instead.

## Author(s)

Roger Gine

#### See Also

adductNames() for a list of all available predefined adducts and adducts() for the adduct data.frame definition style.

## Examples

adductNames

#### Retrieve names of supported adducts

## Description

adductNames returns all supported adduct definitions that can be used by mass2mz() and mz2mass(). adducts returns a data.frame with the adduct definitions.

#### Usage

```
adductNames(polarity = c("positive", "negative"))
```

adducts(polarity = c("positive", "negative"))

# Arguments

polarity character(1) defining the ion mode, either "positive" or "negative".

# Value

for adductNames: character vector with all valid adduct names for the selected ion mode. For adducts: data.frame with the adduct definitions.

## Author(s)

Michael Witting, Johannes Rainer

#### Examples

```
## retrieve names of adduct names in positive ion mode
adductNames(polarity = "positive")
## retrieve names of adduct names in negative ion mode
adductNames(polarity = "negative")
```

calculateKendrickMass Kendrick mass defects

#### Description

Kendrick mass defect analysis is a way to analyze high-resolution MS data in order to identify homologous series. The Kendrick mass (KM) is calculated by choosing a specific molecular fragment (e.g. CH2) and settings its mass to an integer mass. In case of CH2 the mass of 14.01565 would be set to 14.The Kendrick mass defect (KMD) is defined as the difference between the KM and the nominial (integer) KM. All molecules of homologoues series, e.g. only differing in the number of CH2, will have an identical KMD. In an additional step the KMD can be referenced to the mass defect of specific lipid backbone and by this normalize values to the referenced KMD (RKMD). This leads to values of 0 for saturated species or -1, -2, -3, etc for unsaturated species.

Available functoins are:

- calculateKm: calculates the Kendrick mass from an exact mass for a specific molecular fragment, e.g. "CH2".
- calculateKmd: calculates the Kendrick mass defect from an exact mass for a specific molecular fragment, e.g. "CH2".
- calculateRkmd: calculates the referenced Kendrick mass defect from an exact mass for a specific molecular fragment, e.g. "CH2", and a reference KMD.
- isRkmd: Checks if a calculated RKMD falls within a specific error range around an negative integer corresponding the number of double bonds, in case of CH2 as fragment.

#### Usage

```
calculateKm(x, fragment = 14/14.01565)
calculateKmd(x, fragment = 14/14.01565)
calculateRkmd(x, fragment = 14/14.01565, rkmd = 0.749206)
isRkmd(x, rkmdTolerance = 0.1)
```

#### Arguments

х	numeric with exact masses or calculated RKMDs in case of isRkmd.
fragment	numeric(1) or character(1) corresponding factor or molecular formula of
	molecular fragment, e.g. 14 / 14.01565 or "CH2" for CH2.

calculateMass

rkmd	numeric(1) KMD used for referencing of KMDs.
rkmdTolerance	numeric(1) Tolerance to check if RKMD fall around a negative integer corresponding to the number of double bonds

# Value

numeric or boolean. All functions, except isRkmd return a numeric with same length as the input corresponding to the KM, KMD or RMKD. isRkmd returns a logical with TRUE or FALSE indicating if the RKMD falls within a specific range around a negative integer corresponding to the number of double bonds.

#### Author(s)

Michael Witting

## Examples

calculateKm(760.5851)

calculateKmd(760.5851)

calculateRkmd(760.5851, rkmd = 0.749206)

isRkmd(calculateRkmd(760.5851, rkmd = 0.749206))

calculateMass Calculate exact mass

#### Description

calculateMass calculates the exact mass from a formula. Isotopes are also supported. For isotopes, the isotope type needs to be specified as an element's prefix, e.g. "[13C]" for carbon 13 or "[2H]" for deuterium. A formula with 2 carbon 13 isotopes and 3 carbons would thus contain e.g. "[13C2]C3".

#### Usage

```
calculateMass(x)
```

## Arguments

х

character representing chemical formula(s) or a list of numeric with element counts such as returned by countElements(). Isotopes and deuterated elements are supported (see examples below).

# Value

numeric Resulting exact mass.

## containsElements

## Author(s)

Michael Witting

## Examples

```
calculateMass("C6H1206")
calculateMass("NH3")
calculateMass(c("C6H1206", "NH3"))
## Calculate masses for formulas containing isotope information.
calculateMass(c("C6H1206", "[13C3]C3H1206"))
## Calculate mass for a chemical with 5 deuterium.
```

calculateMass("C11[2H5]H7N2O2")

containsElements Check if one formula is contained in another

# Description

containsElements checks if one sum formula is contained in another.

## Usage

```
containsElements(x, y)
```

#### Arguments

х	character strings with a chemical formula
У	character strings with a chemical formula that shall be contained in $\boldsymbol{x}$

# Value

logical TRUE if y is contained in x

#### Author(s)

Michael Witting and Sebastian Gibb

# Examples

containsElements("C6H1206", "H20")
containsElements("C6H1206", "NH3")

```
convertMtime
```

# Description

convertMtime performs effective mobility scale transformation of CE(-MS) data, which is used to overcome variations of the migration times, caused by differences in the Electroosmotic Flow (EOF) between different runs. In order to monitor the EOF and perform the transformation, neutral or charged EOF markers are spiked into the sample before analysis. The information of the EOF markers (migration time and effective mobility) will be then used to perform the effective mobility transformation of the migration time scale.

## Usage

```
convertMtime(
  x = numeric(),
  rtime = numeric(),
  mobility = numeric(),
  tR = 0,
  U = numeric(),
  L = numeric()
)
```

## Arguments

х	numeric vector with migration times in minutes.
rtime	numeric vector that holds the migration times (in minutes) of either one or two EOF markers in the same run of which the migration time is going to be transformed.
mobility	numeric vector containing the respective effective mobility (in in mm <sup>2</sup> / (kV * min)) of the EOF markers. If two markers are used, one is expected to be the neutral marker, i.e. having a mobility of 0.
tR	numeric a single value that defines the time (in minutes) of the electrical field ramp. The default is 0.
U	numeric a single value that defines the voltage (in $kV$ ) applied. Note that for reversed polarity CE mode a negative value is needed. Is only used if the transformation is performed based on a single marker.
L	numeric a single value that defines the total length (in mm) of the capillary that was used for CE(-MS) analysis. Is only used if the transformation is performed based on a single marker.

#### Value

numeric vector of same length as x with effective mobility values.

## correctRindex

## Author(s)

Liesa Salzer

#### Examples

```
rtime <- c(10,20,30,40,50,60,70,80,90,100)
marker_rt <- c(20,80)
mobility <- c(0, 2000)
convertMtime(rtime, marker_rt, mobility)</pre>
```

correctRindex 2-point correction of RIs

# Description

correctRindex performs correction of retention indices (RIs) based on reference substances. Even after conversion of RTs to RIs slight deviations might exist. These deviations can be further normalized, if they are linear, by using two metabolites for which the RIs are known (e.g. internal standards).

## Usage

correctRindex(x, y)

#### Arguments

х	numeric vector with retention indices, calculated by indexRtime
У	data.frame containing two columns. The first is expected to contain the measured RIs of the reference substances and the second the reference RIs.

#### Value

numeric vector of same length than x with corrected retention indices. Values are floating point decimals. If integer values shall be used conversion has to be performed manually.

# Author(s)

Michael Witting

```
ref <- data.frame(rindex = c(110, 210),
refindex = c(100, 200))
rindex <- c(110, 210)
correctRindex(rindex, ref)
```

countElements

#### Description

countElements parses strings representing a chemical formula into a named vector of element counts.

#### Usage

countElements(x)

## Arguments

```
Х
```

character() representing a chemical formula.

# Value

list of integer with the element counts (names being elements).

#### Author(s)

Michael Witting and Sebastian Gibb

## See Also

pasteElements()

#### Examples

countElements(c("C6H1206", "C11H12N202"))

fit\_lm

Linear model-based normalization of abundance matrices

## Description

The fit\_lm and adjust\_lm functions facilitate linear model-based normalization of abundance matrices. The expected noise in a numeric data matrix can be modeled with a linear regression model using the fit\_lm function and the data can subsequently be adjusted using the adjust\_lm function (i.e., the modeled noise will be removed from the abundance values). A typical use case would be to remove injection index dependent signal drifts in a LC-MS derived metabolomics data: a linear model of the form  $y \sim injection_index$  could be used to model the measured abundances of each feature (each row in a data matrix) as a function of the injection index in which a specific sample was measured during the LC-MS measurement run. The fitted linear regression models can

subsequently be used to adjust the abundance matrix by removing these dependencies from the data. This allows to perform signal adjustments as described in (Wehrens et al. 2016).

The two functions are described in more details below:

fit\_lm allows to fit a linear regression model (defined with parameter formula) to each row of the numeric data matrix submitted with parameter y. Additional covariates of the linear model defined in formula are expected to be provided as columns in a data.frame supplied *via* the data parameter.

The linear model is expected to be defined by a formula starting with  $y \sim .$  To model for example an injection index dependency of values in y a formula  $y \sim injection_index$  could be used, with values for the injection index being provided as a column "injection\_index" in the data data frame. fit\_lm would thus fit this model to each row of y.

Linear models can be fitted either with the standard least squares of lm() by setting method = "lm" (the default), or with the more robust methods from the *robustbase* package with method = "lmrob".

adjust\_lm can be used to adjust abundances in a data matrix y based on linear regression models provided with parameter lm. Parameter lm is expected to be a list of length equal to the number of rows of y, each element being a linear model (i.e., a results from lm or lmrob). Covariates for the model need to be provided as columns in a data.frame provided with parameter data. The number of rows of that data.frame need to match the number of columns of y. The function returns the input matrix y with values in rows adjusted with the linear models provided by lm. No adjustment is performed for rows for which the respective element in lm is NA. See examples below for details or the vignette for more examples, descriptions and information.

## Usage

```
fit_lm(
  formula,
  data,
  y,
  method = c("lm", "lmrob"),
  control = NULL,
  minVals = ceiling(nrow(data) * 0.75),
  model = TRUE,
  ...,
 BPPARAM = SerialParam()
)
```

adjust\_lm(y = matrix(), data = data.frame(), lm = list(), ...)

#### Arguments

formula	formula defining the model that should be fitted to the data. See also $lm()$ for more information. Formulas should begin with y ~ as values in rows of y will be defined as y. See description of the fit_lm function for more information.
data	data.frame containing the covariates for the linear model defined by formula (for fit_lm) or used in lm (for adjust_lm). The number of rows has to match the number of columns of y.

У	for fit_lm: matrix of abundances on which the linear model pefined with formula should be fitted. For adjust_lm: matrix of abundances that should be adjusted using the models provided with parameter lm.
method	character(1) defining the linear regression function that should be used for model fitting. Can be either method = "lm" (the default) for standard least squares model fitting or method = "lmrob" for a robust alternative defined in the <i>robustbase</i> package.
control	a list speficying control parameters for lmrob. Only used if method = "lmrob". See help of lmrob.control in the robustbase package for details. By default control = NULL the <i>KS2014</i> settings are used and scale-finding iterations are increased to 10000.
minVals	numeric(1) defining the minimum number of non-missing values (per fea- ture/row) required to perform the model fitting. For rows in y for which fewer non-NA values are available no model will be fitted and a NA will be reported instead.
model	logical(1) whether the model frame are included in the returned linear models. Passed to the lm or lmrob functions.
	for fit_lm: additional parameters to be passed to the downstream calls to lm or lmrob. For adjust_lm: ignored.
BPPARAM	parallel processing setup. See <pre>bpparam()</pre> for more information. Parallel pro- cessing can improve performance especially for method = "lmrob".
lm	list of linear models (as returned by lm or lmrob) such as returned by the fit_lm function. The length of the list is expected to match the number of rows of y, i.e., each element should be a linear model to adjust the specific row, or NA to skip adjustment for that particular row in y.

# Value

For `fit\_lm`: a `list` with linear models (either of type \*lm\* or \*lmrob\*) or length equal to the number of rows of `y`. `NA` is reported for rows with too few non-missing data points (depending on parameter `minValues`). For `adjust\_lm`: a numeric matrix (same dimensions as input matrix `y`) with the values adjusted with the provided linear models.

## Author(s)

Johannes Rainer

# References

Wehrens R, Hageman JA, van Eeuwijk F, Kooke R, Flood PJ, Wijnker E, Keurentjes JJ, Lommen A, van Eekelen HD, Hall RD Mumm R and de Vos RC. Improved batch correction in untargeted MS-based metabolomics. *Metabolomics* 2016; 12:88.

## fit\_lm

```
## See also the vignette for more details and examples.
## Load a test matrix with abundances of features from a LC-MS experiment.
vals <- read.table(system.file("txt", "feature_values.txt",</pre>
                                package = "MetaboCoreUtils"), sep = "\t")
vals <- as.matrix(vals)</pre>
## Define a data.frame with the covariates to be used to model the noise
sdata <- data.frame(injection_index = seq_len(ncol(vals)))</pre>
## Fit a linear model describing the feature abundances as a
## function of the index in which samples were injected during the LC-MS
## run. We're fitting the model to log2 transformed data.
## Note that such a model should **only** be fitted if the samples
## were randomized, i.e. the injection index is independent of any
## experimental covariate. Alternatively, the injection order dependent
## signal drift could be estimated using QC samples (if they were
## repeatedly injected) - see vignette for more details.
ii_lm <- fit_lm(y ~ injection_index, data = sdata, y = log2(vals))</pre>
## The result is a list of linear models
ii_lm[[1]]
## Plotting the data for one feature:
plot(x = sdata$injection_index, y = log2(vals[2, ]),
   ylab = expression(log[2]~abundance), xlab = "injection index")
grid()
## plot also the fitted model
abline(ii_lm[[2]], lty = 2)
## For this feature (row) a decreasing signal intensity with injection
## index was observed (and modeled).
## For another feature an increasing intensity can be observed.
plot(x = sdata$injection_index, y = log2(vals[3, ]),
    ylab = expression(log[2]~abundance), xlab = "injection index")
grid()
## plot also the fitted model
abline(ii_lm[3]), lty = 2)
## This trend can be removed from the data using the `adjust_lm` function
## by providing the linear models descring the drift. Note that, because
## we're adjusting log2 transformed data, the resulting abundances are
## also in log2 scale.
vals_adj <- adjust_lm(log2(vals), data = sdata, lm = ii_lm)</pre>
## Plotting the data before (open circles) and after adjustment (filled
## points)
plot(x = sdata$injection_index, y = log2(vals[2, ]),
    ylab = expression(log[2]~abundance), xlab = "injection index")
points(x = sdata$injection_index, y = vals_adj[2, ], pch = 16)
```

```
grid()
## Adding the line fitted through the raw data
abline(ii_lm[[2]], lty = 2)
## Adding a line fitted through the adjusted data
abline(lm(vals_adj[2, ] ~ sdata$injection_index), lty = 1)
## After adjustment there is no more dependency on injection index.
```

formula2mz

Calculate mass-to-charge ratio from a formula

## Description

formula2mz calculates the m/z values from a list of molecular formulas and adduct definitions.

Custom adduct definitions can be passed to the adduct parameter in form of a data.frame. This data.frame is expected to have columns "mass\_add" and "mass\_multi" defining the *additive* and *multiplicative* part of the calculation. See adducts() for examples.

#### Usage

```
formula2mz(formula, adduct = "[M+H]+", standardize = TRUE)
```

# Arguments

formula	character with one or more valid molecular formulas for which their adduct m/z shall be calculated.
adduct	either a character specifying the name(s) of the adduct(s) for which the m/z should be calculated or a data.frame with the adduct definition. See adductNames() for supported adduct names and the description for more information on the expected format if a data.frame is provided.
standardize	logical whether to standardize the molecular formulas to the Hill notation system before calculating their mass.

## Value

Numeric matrix with same number of rows than elements in formula and number of columns being equal to the length of adduct (adduct names are used as column names). Each column thus represents the m/z of formula for each defined adduct.

# Author(s)

Roger Gine

14

# indexRtime

## Examples

```
## Calculate m/z values of adducts of a list of formulas
formulas <- c("C6H1206", "C9H11NO3", "C16H13C1N20")
ads <- c("[M+H]+", "[M+Na]+", "[2M+H]+", "[M]+")
formula2mz(formulas, ads)
formula2mz(formulas, adductNames()) #All available adducts
## Use custom-defined adducts as input
custom_ads <- data.frame(mass_add = c(1, 2, 3), mass_multi = c(1, 2, 0.5))
formula2mz(formulas, custom_ads)
## Use standardize = FALSE to keep formula unaltered
formula2mz("H12C606")
formula2mz("H12C606", standardize = FALSE)
```

```
indexRtime
```

Convert retention times to retention indices

## Description

indexRtime uses a list of known substances to convert retention times (RTs) to retention indices (RIs). By this retention information is normalized for differences in experimental settings, such as gradient delay volume, dead volume or flow rate. By default linear interpolation is performed, other ways of calculation can supplied as function.

## Usage

```
indexRtime(x, y, FUN = rtiLinear, ...)
```

# Arguments

х	numeric vector with retention times
У	data.framedata.frame containing two columns, where the first holds the reten- tion times of the indexing substances and the second the actual index value
FUN	function function defining how the conversion is performed, default is linear interpolation
	additional parameter used by FUN

# Value

numeric vector of same length as x with retention indices. Values floating point decimals. If integer values shall be used conversion has to be performed manually

## Author(s)

Michael Witting

## Examples

```
rti <- data.frame(rtime = c(1,2,3),
rindex = c(100,200,300))
rtime <- c(1.5, 2.5)
indexRtime(rtime, rti)
```

internalStandardMixNames

Get names of internal standard mixes provided by the package

#### Description

internalStandardMixNames returns available names of internal standard mixes provided by the MetaboCoreUtils package.

## Usage

```
internalStandardMixNames()
```

# Value

character names of available IS mixes

#### Author(s)

Michael Witting

#### Examples

internalStandardMixNames()

internalStandards Get definitions for internal standards

#### Description

internalStandards returns a table with metabolite standards available in commercial internal standard mixes. The returned data frame contains the following columns:

- "name": the name of the standard
- "formula\_salt": chemical formula of the salt that was used to produce the standard mix
- "formula\_metabolite": chemical formula of the metabolite in free form
- "smiles\_salt": SMILES of the salt that was used to produced the standard mix
- "smiles\_metabolite": SMILES of the metabolite in free form

16

#### isotopicSubstitutionMatrix

- "mol\_weight\_salt": molecular (average) weight of the salt (can be used for calculation of molar concentration, etc.)
- "exact\_mass\_metabolite": exact mass of free metabolites
- "conc": concentration of the metabolite in ug/mL (of salt form)
- "mix": name of internal standard mix

## Usage

```
internalStandards(mix = "QReSS")
```

## Arguments

mix

character(1) Name of the internal standard mix that shall be returned. One of internalStandardMixNames().

## Value

data.frame data on internal standards

#### Author(s)

Michael Witting

#### See Also

internalStandardMixNames() for provided internal standard mixes.

## Examples

```
internalStandards(mix = "QReSS")
internalStandards(mix = "UltimateSplashOne")
```

isotopicSubstitutionMatrix

Definitions of isotopic substitutions

## Description

In order to identify potential isotopologues based on only m/z and intensity values with the isotopologues() function, sets of pre-calculated parameters are required. This function returns such parameter sets estimated on different sources/databases. The nomenclature used to describe isotopes follows the following convention: the number of neutrons is provided in [ as a prefix to the element and the number of atoms of the element as suffix. [13]C2[37]Cl3 describes thus an isotopic substitution containing 2 [13]C isotopes and 3 [37]Cl isotopes.

Each row in the returned data.frame is associated with an isotopic substitution (which can involve isotopes of several elements or different isotopes of the same element). In general for each isotopic substitution multiple rows are present in the data.frame. Each row provides parameters to compute

bounds (for the ratio between the isotopologue peak and the monoisotopic one) on a certain mass range. The provided isotopic substitutions are in general the most frequently observed substitutions in the database (e.g. HMDB) on which they were defined. Parameters (columns) defined for each isotopic substitution are:

- "minmass": the minimal mass of a compound for which the isotopic substitution was found. Peaks with a mass lower than this will most likely not have the respective isotopic substitution.
- "maxmass": the maximal mass of a compound for which the isotopic substitution was found. Peaks with a mass higher than this will most likely not have the respective isotopic substitution.
- "md": the mass difference between the monoisotopic peak and a peak of an isotopologue characterized by the respective isotopic substitution.
- "leftend": left endpoint of the mass interval.
- "rightend": right endpoint of the mass interval.
- "LBint": intercept of the lower bound line on the mass interval whose endpoints are "leftend" and "rightend".
- "LBslope": slope of the lower bound line on the mass interval.
- "UBint": intercept of the upper bound line on the mass interval.
- "UBslope": slope of the upper bound line on the mass interval.

#### Usage

```
isotopicSubstitutionMatrix(source = c("HMDB_NEUTRAL"))
```

# Arguments

source character(1) defining the set of predefined parameters and isotopologue definitions to return.

# Value

data. frame with parameters to detect the defined isotopic substitutions

#### Available pre-calculated substitution matrices

• source = "HMDB": most common isotopic substitutions and parameters for these have been calculated for all compounds from the Human Metabolome Database (HMDB, July 2021). Note that the substitutions were calculated on the **neutral masses** (i.e. the chemical formulas of the compounds, not considering any adducts).

## Author(s)

Andrea Vicini

```
## Get the substitution matrix calculated on HMDB
isotopicSubstitutionMatrix("HMDB_NEUTRAL")
```

## Description

Given a spectrum (i.e. a peak matrix with m/z and intensity values) the function identifies groups of potential isotopologue peaks based on pre-defined mass differences and intensity (probability) ratios that need to be passed to the function with the substDefinition parameter. Each isotopic substitution in a compound determines a certain isotopologue and it is associated with a certain mass difference of that with respect to the monoisotopic isotopologue. Also each substitution in a compound is linked to a certain ratio between the intensities of the peaks of the corresponding isotopologue and the monoisotopic one. This ratio isn't the same for isotopologues corresponding to the same isotopic substitution but to different compounds. Through the substDefinition parameter we provide upper and lower values to compute bounds for each isotopic substitution dependent on the peak's mass.

#### Usage

```
isotopologues(
    x,
    substDefinition = isotopicSubstitutionMatrix(),
    tolerance = 0,
    ppm = 20,
    seedMz = numeric(),
    charge = 1,
    .check = TRUE
)
```

#### Arguments

X	matrix or data.frame with spectrum data. The first column is expected to con- tain $m/z$ and the second column intensity values. The $m/z$ values in that matrix are expected to be increasingly ordered and no NA values should be present.
substDefinition	
	<pre>matrix or data.frame with definition of isotopic substitutions (columns "name" and "md" are among the mandatory columns). The rows in this matrix have to be ordered by column md in increasing order. See isotopicSubstitutionMatrix() for more information on the format and content.</pre>
tolerance	<pre>numeric(1) representing the absolute tolerance for the relaxed matching of m/z values of peaks. See MsCoreUtils::closest() for details.</pre>
ppm	numeric(1) representing a relative, value-specific parts-per-million (PPM) tol- erance for the relaxed matching of m/z values of peaks. See MsCoreUtils::closest() for details.
seedMz	numeric vector of <b>ordered</b> m/z values. If provided, the function checks if there are peaks in x which m/z match them. If any, it looks for groups where the first peak is one of the matched ones.

isotopologues

charge	numeric(1) representing the expected charge of the ionized compounds.
. check	logical(1) to disable input argument check. Should only be set to FALSE if provided $m/z$ values are guaranteed to be increasingly ordered and don't contain NA values.

# Details

The function iterates over the peaks (rows) in x. For each peak (which is assumed to be the monoisotopic peak) it searches other peaks in x with a difference in mass matching (given ppm and tolerance) any of the pre-defined mass differences in substDefinitions (column "md"). The mass is obtained by multiplying the m/z of the peaks for the charge expected for the ionized compounds.

For matching peaks, the function next evaluates whether their intensity is within the expected (pre-defined) intensity range. Using "LBint", "LBslope", "UBint", "UBslope" of the previously matched isotopic substitution in substDefinition, the function estimates a (mass dependent) lower and upper intensity ratio limit based on the peak's mass.

When some peaks are grouped together their indexes are excluded from the set of indexes that are searched for further groups (i.e. peaks already assigned to an isotopologue group are not considered/tested again thus each peak can only be part of one isotopologue group).

# Value

list of integer vectors. Each integer vector contains the indixes of the rows in x with potential isotopologues of the same compound.

#### Author(s)

Andrea Vicini

```
## Read theoretical isotope pattern (high resolution) from example file
x <- read.table(system.file("exampleSpectra",
    "serine-alpha-lactose-caffeine.txt", package = "MetaboCoreUtils"),
    header = TRUE)
x <- x[order(x$mz), ]
plot(x$mz, x$intensity, type = "h")
isos <- isotopologues(x, ppm = 5)
isos
## highlight them in the plot
for (i in seq_along(isos)) {
    z <- isos[[i]]
    points(x$mz[z], x$intensity[z], col = i + 1)
}
```

mass2mz

## Description

mass2mz calculates the m/z value from a neutral mass and an adduct definition.

Custom adduct definitions can be passed to the adduct parameter in form of a data.frame. This data.frame is expected to have columns "mass\_add" and "mass\_multi" defining the *additive* and *multiplicative* part of the calculation. See adducts() for examples.

## Usage

mass2mz(x, adduct = "[M+H]+")

#### Arguments

х	numeric neutral mass for which the adduct m/z shall be calculated.
adduct	either a character specifying the name(s) of the adduct(s) for which the m/z
	should be calculated or a data.frame with the adduct definition. See adductNames()
	for supported adduct names and the description for more information on the ex-
	pected format if a data.frame is provided.

#### Value

numeric matrix with same number of rows than elements in x and number of columns being equal to the length of adduct (adduct names are used as column names). Each column thus represents the m/z of x for each defined adduct.

#### Author(s)

Michael Witting, Johannes Rainer

### See Also

mz2mass() for the reverse calculation, adductNames() for supported adduct definitions.

```
exact_mass <- c(100, 200, 250)
adduct <- "[M+H]+"
## Calculate m/z of [M+H]+ adduct from neutral mass
mass2mz(exact_mass, adduct)
exact_mass <- 100
adduct <- "[M+Na]+"
## Calculate m/z of [M+Na]+ adduct from neutral mass</pre>
```

mclosest

```
mass2mz(exact_mass, adduct)
## Calculate m/z of multiple adducts from neutral mass
mass2mz(exact_mass, adduct = adductNames())
## Provide a custom adduct definition.
adds <- data.frame(mass_add = c(1, 2, 3), mass_multi = c(1, 2, 0.5))
rownames(adds) <- c("a", "b", "c")
mass2mz(c(100, 200), adds)</pre>
```

mclosest

Extract closest values in a pairwise manner between two matrices

#### Description

The mclosest function calculates the closest rows between two matrices (or data frames) considering pairwise differences between values in columns of x and table. It returns the index of the closest row in table for each row in x.

## Usage

mclosest(x, table, ppm = 0, tolerance = Inf)

## Arguments

X	numeric matrix or data frame representing the query data. Each row in $x$ will be compared to every row in table. Both $x$ and table are expected to have the same number of columns, and the columns are expected to be in the same order.
table	numeric matrix or data frame containing the reference data to be matched with each row of x. Each row in table will be compared to every row in x. Both table and x are expected to have the same number of columns, and the columns are expected to be in the same order.
ppm	numeric representing a relative, value-specific parts-per-million (PPM) toler- ance that is added to tolerance (default is 0).
tolerance	numeric accepted tolerance. Defaults to tolerance = $Inf$ , thus for each row in x the closest row in table is reported, regardless of the magnitude of the (absolute) difference.

#### Details

If, for a row of x, two rows of table are closest only the index of first row will be returned.

For both the tolerance and ppm arguments, if their length is different to the number of columns of x and table, the input argument will be replicated to match it.

## Value

integer vector of indices indicating the closest row of table for each row of x. If no suitable match is found for a row in x based on the specified tolerance and ppm, the corresponding index is set to NA.

22

# multiplyElements

## Author(s)

Philippine Louail

#### Examples

```
x <- data.frame(a = 1:5, b = 3:7)
table <- data.frame(c = c(11, 23, 3, 5, 1), d = c(32:35, 45))
## Get for each row of `x` the index of the row in `table` with the smallest
## difference of values (per column)
mclosest(x, table)
## If the absolute difference is larger than `tolerance`, return `NA`. Note
## that the tolerance value of `25` is used for difference for each pairwise
## column in `x` and `table`.
mclosest(x, table, tolerance = 25)</pre>
```

multiplyElements *Multiply chemical formulas by a scalar* 

## Description

multiplyElements Multiply the number of atoms of each element by a constant, positive, integer

#### Usage

```
multiplyElements(x, k)
```

### Arguments

Х	character strings with chemical formula
k	numeric(1) positive integer by which each formula will be multiplied

# Value

character strings with the standardized chemical formula.

# Author(s)

Roger Gine

```
multiplyElements("H2O", 3)
```

```
multiplyElements(c("C6H1206", "Na", "CH40"), 2)
```

mz2mass

#### Description

mz2mass calculates the neutral mass from a given m/z value and adduct definition.

Custom adduct definitions can be passed to the adduct parameter in form of a data.frame. This data.frame is expected to have columns "mass\_add" and "mass\_multi" defining the *additive* and *multiplicative* part of the calculation. See adducts() for examples.

#### Usage

mz2mass(x, adduct = "[M+H]+")

## Arguments

х	numeric m/z value for which the neutral mass shall be calculated.
adduct	either a character specifying the name(s) of the adduct(s) for which the m/z
	should be calculated or a data.frame with the adduct definition. See adductNames()
	for supported adduct names and the description for more information on the ex-
	pected format if a data.frame is provided.

#### Value

numeric matrix with same number of rows than elements in x and number of columns being equal to the length of adduct (adduct names are used as column names. Each column thus represents the neutral mass of x for each defined adduct.

## Author(s)

Michael Witting, Johannes Rainer

#### See Also

mass2mz() for the reverse calculation, adductNames() for supported adduct definitions.

```
ion_mass <- c(100, 200, 300)
adduct <- "[M+H]+"
## Calculate m/z of [M+H]+ adduct from neutral mass
mz2mass(ion_mass, adduct)
ion_mass <- 100
adduct <- "[M+Na]+"
## Calculate m/z of [M+Na]+ adduct from neutral mass</pre>
```

# pasteElements

```
mz2mass(ion_mass, adduct)
## Provide a custom adduct definition.
adds <- data.frame(mass_add = c(1, 2, 3), mass_multi = c(1, 2, 0.5))
rownames(adds) <- c("a", "b", "c")
mz2mass(c(100, 200), adds)</pre>
```

pasteElements Create chemical formula from a named vector

# Description

pasteElements creates a chemical formula from element counts (such as returned by countElements()).

## Usage

```
pasteElements(x)
```

## Arguments

х

list/integer with element counts, names being individual elements.

## Value

character() with the chemical formulas.

# Author(s)

Michael Witting and Sebastian Gibb

# See Also

countElements()

```
elements <- c("C" = 6, "H" = 12, "O" = 6)
pasteElements(elements)</pre>
```

quality\_assessment

#### Description

The following functions allow to calculate basic quality assessment estimates typically employed in the analysis of metabolomics data. These functions are designed to be applied to entire rows of data, where each row corresponds to a feature. Subsequently, these estimates can serve as a foundation for feature filtering.

- rsd and rowRsd are convenience functions to calculate the relative standard deviation (i.e. coefficient of variation) of a numerical vector or for rows of a numerical matrix, respectively.
- rowDratio computes the D-ratio or *dispersion ratio*, defined as the standard deviation for QC (Quality Control) samples divided by the standard deviation for biological test samples, for each feature (row) in the matrix.
- percentMissing and rowPercentMissing determine the percentage of missing values in a vector or for each row of a matrix, respectively.
- rowBlank identifies rows (i.e., features) where the mean of test samples is lower than a specified multiple (defined by the threshold parameter) of the mean of blank samples. This can be used to flag features that result from contamination in the solvent of the samples.

These functions are based on standard filtering methods described in the literature, and they are implemented to assist in preprocessing metabolomics data.

### Usage

rsd(x, na.rm = TRUE, mad = FALSE)
rowRsd(x, na.rm = TRUE, mad = FALSE)
rowDratio(x, y, na.rm = TRUE, mad = FALSE)
percentMissing(x)
rowPercentMissing(x)
rowBlank(x, y, threshold = 2, na.rm = TRUE)

#### Arguments

x	numeric For rsd, a numeric vector; for rowRsd, rowDratio, percentMissing and rowBlank, a numeric matrix representing the biological samples.
na.rm	logical(1) indicates whether missing values (NA) should be removed prior to the calculations.
mad	logical(1) indicates whether the <i>Median Absolute Deviation</i> (MAD) should be used instead of the standard deviation. This is suggested for non-gaussian distributed data.

У	numeric For rowDratio and rowBlank, a numeric matrix representing feature abundances in QC samples or blank samples, respectively.
threshold	numeric For rowBlank, indicates the minimum difference required between the mean of a feature in samples compared to the mean of the same feature in blanks for it to not be considered a possible contaminant. For example, the default threshold of 2 signifies that the mean of the features in samples has to be at least twice the mean in blanks for it not to be flagged as a possible contaminant.

# Value

See individual function description above for details.

# Note

For rsd and rowRsd the feature abundances are expected to be provided in natural scale and not e.g. log2 scale as it may lead to incorrect interpretations.

## Author(s)

Philippine Louail, Johannes Rainer

## References

Broadhurst D, Goodacre R, Reinke SN, Kuligowski J, Wilson ID, Lewis MR, Dunn WB. Guidelines and considerations for the use of system suitability and quality control samples in mass spectrometry assays applied in untargeted clinical metabolomic studies. Metabolomics. 2018;14(6):72. doi: 10.1007/s11306-018-1367-3. Epub 2018 May 18. PMID: 29805336; PMCID: PMC5960010.

```
## coefficient of variation
a <- c(4.3, 4.5, 3.6, 5.3)
rsd(a)
A <- rbind(a, a, a)
rowRsd(A)
## Dratio
x \le c(4.3, 4.5, 3.6, 5.3)
X <- rbind(a, a, a)
rowDratio(X, X)
#' ## Percent Missing
b <- c(1, NA, 3, 4, NA)
percentMissing(b)
B <- matrix(c(1, 2, 3, NA, 5, 6, 7, 8, 9), nrow = 3)
rowPercentMissing(B)
## Blank Rows
test_samples <- matrix(c(13, 21, 3, 4, 5, 6), nrow = 2)</pre>
```

```
blank_samples <- matrix(c(0, 1, 2, 3, 4, 5), nrow = 2)
rowBlank(test_samples, blank_samples)</pre>
```

standardizeFormula Standardize a chemical formula

# Description

standardizeFormula standardizes a supplied chemical formula according to the Hill notation system.

## Usage

```
standardizeFormula(x)
```

## Arguments

x character, strings with the chemical formula to standardize.

## Value

character strings with the standardized chemical formula.

# Author(s)

Michael Witting and Sebastian Gibb

# See Also

pasteElements() countElements()

# Examples

standardizeFormula("C606H12")

28

subtractElements subtract two chemical formula

# Description

subtractElements subtracts one chemical formula from another.

# Usage

```
subtractElements(x, y)
```

# Arguments

x	character strings with chemical formula
y	character strings with chemical formula that should be subtracted from >

# Value

character Resulting formula

# Author(s)

Michael Witting and Sebastian Gibb

```
subtractElements("C6H1206", "H20")
```

```
subtractElements("C6H12O6", "NH3")
```

# Index

```
addElements, 2
adductFormula, 3
adductNames, 4
adductNames(), 4, 14, 21, 24
adducts (adductNames), 4
adducts(), 3, 4, 14, 21, 24
adjust_lm (fit_lm), 10
```

bpparam(), *12* 

```
calculateKendrickMass, 5
calculateKm (calculateKendrickMass), 5
calculateKmd (calculateKendrickMass), 5
calculateMass, 6
calculateRkmd (calculateKendrickMass), 5
containsElements, 7
convertMtime, 8
correctRindex, 9
countElements, 10
countElements(), 6, 25, 28
```

```
fit_lm, 10
formula2mz, 14
```

```
indexRtime, 15
internalStandardMixNames, 16
internalStandardMixNames(), 17
internalStandards, 16
isotopicSubstitutionMatrix, 17
isotopologues, 19
isotopologues(), 17
isRkmd (calculateKendrickMass), 5
```

```
lm(), 11
```

```
mass2mz, 21
mass2mz(), 4, 24
mclosest, 22
MsCoreUtils::closest(), 19
multiplyElements, 23
```

mz2mass, 24
mz2mass(), 4, 21

pasteElements, 25
pasteElements(), 10, 28
percentMissing (quality\_assessment), 26

quality\_assessment, 26

standardizeFormula, 28
subtractElements, 29