

# Package ‘proBatch’

June 14, 2018

**Type** Package

**Title** Tools for Batch Effects Diagnostics and Correction

**Version** 1.1.0

**Author** Jelena Čuklina <chuklina.jelena@gmail.com>

**Maintainer** The package maintainer <chuklina.jelena@gmail.com>

**Description** The proBatch package contains functions for diagnosing and removing batch effects and other unwanted variation in high-throughput experiment, primarily designed for DIA proteomics data.

The diagnostic part of the package can be broadly divided in (1) Genome-wide and (2) Gene-specific functions, explained in corresponding vignettes. Since the diagnostic part for batch effects does require batch effect removal, here we provide a few convenience wrappers for common batch-effect removal approaches, namely, ComBat (Johnson et al. 2007 Biostatistics) and mean/median centering. However, proteomics data may require more complicated technical artifact correction approaches like non-linear fitting, which is also found in “normalization” section of this package.

The approaches are described in (Čuklina et al. 2019, MCP)

**License** What license is it under?

**URL** <https://github.com/jelenachuklina/proBatch>

**BugReports** <https://github.com/jelenachuklina/proBatch/issues>

**Depends** R (>= 3.4.0),  
dplyr (>= 0.7.4),  
ggplot2

**Encoding** UTF-8

**LazyData** true

**Imports** Biobase,  
corrplot,  
ggfortify,  
data.table,  
lazyeval,  
lubridate,  
magrittr,  
pheatmap,  
preprocessCore,  
pvca,  
RColorBrewer,

readr,  
 reshape2,  
 rlang,  
 scales,  
 sva,  
 tibble,  
 tidyverse ( $\geq 1.2.1$ ),  
 wesanderson,  
 WGCNA  
**Suggests** knitr,  
 rmarkdown,  
 SWATH2stats  
**VignetteBuilder** knitr  
**RoxygenNote** 6.0.1

## R topics documented:

boxplot_all_steps . . . . .	3
clean_requants . . . . .	4
color_list_to_df . . . . .	5
convert_to_matrix . . . . .	5
create_peptide_annotation . . . . .	6
dates_to_posix . . . . .	6
date_to_sample_order . . . . .	7
define_batches_by_MS_pauses . . . . .	8
define_batches_by_MS_pauses_within_instrument . . . . .	8
example_peptide_annotation . . . . .	9
example_proteome . . . . .	10
example_sample_annotation1 . . . . .	10
generate_colors_for_numeric . . . . .	11
get_peptide_corr_df . . . . .	12
get_sample_corr_distrib . . . . .	12
join_data_matrices . . . . .	13
matrix_to_long . . . . .	14
merge_rare_levels . . . . .	15
normalize . . . . .	15
plot_clustering . . . . .	16
plot_corr_matrix . . . . .	17
plot_heatmap . . . . .	18
plot_iRTs . . . . .	19
plot_pca . . . . .	20
plot_peptides_of_one_protein . . . . .	21
plot_peptide_correlation_distr_one_protein . . . . .	22
plot_peptide_level . . . . .	22
plot_protein_corrplot . . . . .	23
plot_prot_corr_distribution . . . . .	24
plot_pvca . . . . .	25
plot_samples_corrplot . . . . .	26
plot_sample_corr_distribution . . . . .	27
plot_sample_means_or_boxplots . . . . .	28
plot_spike_ins . . . . .	29

boxplot\_all\_steps

3

plot\_within\_prot\_corr\_distribution . . . . .

30

plot\_with\_fitting\_curve . . . . .

31

proBatch . . . . .

32

quantile\_normalize . . . . .

33

remove\_peptides\_with\_missing\_batch . . . . .

33

sample\_annotation\_to\_colors . . . . .

34

sample\_random\_peptides . . . . .

35

summarize\_peptides . . . . .

35

Index

37

boxplot_all_steps	<i>Plot boxplots to compare various data normalization steps/approaches</i> <i>WARNING: extremely slow for big dataframes</i>
-------------------	--

Description

Plot boxplots to compare various data normalization steps/approaches WARNING: extremely slow for big dataframes

Usage

boxplot\_all\_steps(list\_of\_dfs, sample\_annotation, batch\_col, step = NULL)

Arguments

- list\_of\_dfs

list of data frames of format, specified in ‘plot\_boxplot’
- sample\_annotation

data matrix with 1) ‘sample\_id\_col’ (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
- batch\_col

column in ‘sample\_annotation’ that should be used for batch comparison
- step

normalization step (e.g. ‘Raw’ or ‘Quantile\_normalized’ or ‘qNorm\_ComBat’). Useful if consecutive steps are compared in plots. Note that in plots these are usually ordered alphabetically, so it’s worth naming with numbers, e.g. ‘1\_raw’, ‘2\_quantile’

Value

ggplot object

See Also

[plot\\_boxplot](#)

---

clean_requants	<i>Remove requanted and sparse features</i>
----------------	---

---

## Description

Cleans dataset `df_long` ([proBatch](#)) by removing requanted features and features not meeting user defined sparsness criterias.

## Usage

```
clean_requants(df_long, sample_annotation, batch_col = "MS_batch.final",
               feature_id_col = "peptide_group_label", missing_frac_batch = 0.3,
               missing_frac_total = 0.3)
```

## Arguments

<code>df_long</code>	data frame where each row is a single feature in a single sample. It minimally has a <code>sample_id_col</code> , a <code>feature_id_col</code> and a <code>measure_col</code> , but usually also an <code>m_score</code> (in OpenSWATH output result file)
<code>sample_annotation</code>	data matrix with: <ol style="list-style-type: none"> <li>1. <code>sample_id_col</code> (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
<code>batch_col</code>	column in <code>sample_annotation</code> that should be used for batch comparison
<code>feature_id_col</code>	name of the column with feature/gene/peptide/protein ID used in the long format representation <code>df_long</code> . In the wide formatted representation <code>data_matrix</code> this corresponds to the row names.
<code>missing_frac_batch</code>	maximally tolerated fraction of missing values for a feature in a batch
<code>missing_frac_total</code>	maximally tolerated fraction of globally missing values for a feature

## Value

`df_long` ([proBatch](#)) like data frame filtered as follow:

- remove requant values
- remove features not meeting batch or global sparsness thresholds

## See Also

Other dataset cleaning functions: [remove\\_peptides\\_with\\_missing\\_batch](#), [summarize\\_peptides](#)

---

color_list_to_df	<i>Color list to data frame</i>
------------------	---------------------------------

---

**Description**

Turn color list to df (some plotting functions require the latter)

**Usage**

```
color_list_to_df(color_list, sample_annotation)
```

**Arguments**

color_list	list of colors
sample_annotation	factor-based configuration of the sample annotation

**Value**

a data frame representation of the input color list

---

convert_to_matrix	<i>Long to wide conversion</i>
-------------------	--------------------------------

---

**Description**

Convert from a long data frame representation to a wide matrix representation

**Usage**

```
convert_to_matrix(df_long, feature_id_col = "peptide_group_label",
  measure_col = "Intensity", sample_id_col = "FullRunName")
```

**Arguments**

df_long	data frame where each row is a single feature in a single sample. It minimally has a sample_id_col, a feature_id_col and a measure_col, but usually also an m_score (in OpenSWATH output result file)
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
measure_col	if df_long is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)

**Value**

data\_matrix ([proBatch](#)) like matrix (features in rows, samples in columns)

See Also

Other matrix manipulation functions: [join\\_data\\_matrices](#), [matrix\\_to\\_long](#)

---

create_peptide_annotation	<i>Create light-weight peptide annotation data frame for selection of illustrative proteins</i>
---------------------------	---

---

Description

Create light-weight peptide annotation data frame for selection of illustrative proteins

Usage

```
create_peptide_annotation(df_long, peptide_col = "peptide_group_label",
  protein_col = c("Uniprot_ID", "Gene"))
```

Arguments

- peptide\_col      column containing peptide ID
- protein\_col      one or more columns contatining protein ID

See Also

[plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_corrplot\\_protein](#), [plot\\_within\\_prot\\_distribution](#)

Examples

```
\donotrun{
  peptide_annotation = create_peptide_annotation(example_proteome)
  peptide_summary =
}
```

---

dates_to_posix	<i>Convert data/time to POSIXct</i>
----------------	-------------------------------------

---

Description

convert date/time column of sample\_annotation to POSIX format required to keep number-like behaviour

Usage

```
dates_to_posix(sample_annotation, time_column = c("RunDate", "RunTime"),
  new_time_column = NULL, dateTimeFormat = c("%b_%d", "%H:%M:%S"))
```

**Arguments**

sample\_annotation  
data matrix with:

1. sample\_id\_col (this can be repeated as row names)
2. biological covariates
3. technical covariates (batches etc)

time\_column name of the column(s) where run date & time are specified. These will be used to determine the run order

new\_time\_column name of the new column to which date&time will be converted to

dateTimeFormat POSIX format of the date and time. See [as.POSIXct](#) from base R for details

**Value**

sample annotation file with column names as 'new\_time\_column' with POSIX-formatted date

---

date\_to\_sample\_order    *Convert date/time to POSIXct and rank samples by it*

---

**Description**

Converts date/time columns fo sample\_annotation to POSIXct format and calculates sample run rank in order column

**Usage**

```
date_to_sample_order(sample_annotation, time_column = c("RunDate", "RunTime"),
  new_time_column = "DateTime", dateTimeFormat = c("%b_%d",
    "%H:%M:%S"), order_col = "order", instrument_col = "instrument")
```

**Arguments**

sample\_annotation  
data matrix with:

1. sample\_id\_col (this can be repeated as row names)
2. biological covariates
3. technical covariates (batches etc)

time\_column name of the column(s) where run date & time are specified. These will be used to determine the run order

new\_time\_column name of the new column to which date&time will be converted to

dateTimeFormat POSIX format of the date and time. See [as.POSIXct](#) from base R for details

**Value**

sample annotation file with column names as 'new\_time\_column' with POSIX-formatted date & order\_col used in some diagnostic plots (e.g. [plot\\_iRTs](#), [plot\\_sample\\_mean](#))

---

```
define_batches_by_MS_pauses
```

*Batch by date/time and instrument*

---

### Description

Identify long stretches of time between samples on a per instrument basis and split them into batches.

### Usage

```
define_batches_by_MS_pauses(sample_annotation, threshold,
  runtime_col = "RunDateTime", minimal_batch_size = 5,
  instrument_col = "instr", batch_name = "MS_batch")
```

### Arguments

sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>1. sample_id_col (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
threshold	time difference that would mean there was an interruption
runtime_col	POSIX or numeric-like column corresponding to the sample MS profile acquisition timepoint
minimal_batch_size	minimal number of samples in a batch
instrument_col	column specifying MS instrument used to acquired data (to account for the presence of multiple instruments in sample_annotation)
batch_name	string with a self-explanatory name for the batch (e.g. 'MS_batch' for MS-proteomics) to which batch number will be added

### Value

sample\_annotation data matrix with an additional column to indicate sample batching by MS run time and instrument.

---

```
define_batches_by_MS_pauses_within_instrument
```

*Batch by date/time*

---

### Description

Identify long stretches of time between samples and split them into batches. Most users are going to want to call define\_batches\_by\_MS\_pauses, rather than this function



**Usage**

```
define_batches_by_MS_pauses_within_instrument(date_vector, threshold,
  minimal_batch_size = 5, batch_name = "MS_batch")
```

**Arguments**

date_vector	POSIX or numeric-like vector corresponding to the sample MS profile acquisition timepoint
threshold	time difference that would mean there was an interruption
minimal_batch_size	minimal number of samples in a batch
batch_name	string with a self-explanatory name for the batch (e.g. 'MS_batch' for MS-proteomics) to which batch number will be added

**Value**

vector of batches for each sample

---

example_peptide_annotation
<i>Peptide annotation data</i>

---

**Description**

This is data from Evan's aging study annotated with gene names

**Usage**

```
example_peptide_annotation
```

**Format**

A data frame with 200625 rows and 8 variables:

**Peptide** peptide group label ID, identical to 'peptide\_group\_label' in 'example\_proteome'

**Gene** HUGO gene ID

**ProteinName** protein group name as specified in 'example\_proteome' ... other parameters determined by 'summarize\_peptides' function

---

example_proteome	<i>Example protein data</i>
------------------	-----------------------------

---

### Description

This is data from Evan's aging study with all iRT, spike-in peptides, few random peptides and QTL proteins for biological signal improvement demonstration

### Usage

example\_proteome

### Format

A data frame with 200625 rows and 8 variables:

**peptide\_group\_label** peptide ID, which is regular feature level. This column is mostly used as 'feature\_id\_col'

**RT** retention time. Relevant to identify retention time related bias

**Intensity** peptide group intensity in given sample. Used in function as 'measure\_col'

**ProteinName** Protein group ID, specified as N/UniProtID1|UniProtID2|..., where N is number of protein peptide group maps to. If 1/UniProtID, then this is proteotypic peptide

**assay\_rt** retention time as in DIA library

**m\_score** peptide group identification FDR as determined by pyProphet

**FullRunName** name of the file, in most functions used for 'sample\_id\_col' '#' ...

### Source

PRIDE ID will be added in future

---

example_sample_annotation1	<i>Sample annotation data version 1</i>
----------------------------	---

---

### Description

This is data from Evan's aging study with mock instruments to show how instrument-specific functionality works

### Usage

example\_sample\_annotation1

**Format**

A data frame with 375 rows and 18 variables:

**FullRunName** name of the file, in most functions used for 'sample\_id\_col'

**MS\_batch.final** mass-spectrometry batch: 7-level factor of manually annotated batches

**EarTag** mouse ID, i.e. ID of the biological object

**Strain** mouse strain ID - biological covariate #1

**Diet** diet - either 'HFD' = 'High Fat Diet' or 'CD' = 'Chow Diet'. 'Mix' stands for mixture of several samples

**Sex** mice sex - 3-level biological covariate. Possible values - "

**Age\_Days** mice age at sampling - numeric biological covariate

**RunDate** mass-spectrometry running date. In combination with 'RunTime' used for running order determination

**RunTime** mass-spectrometry running time. In combination with 'RunDate' used for running order determination

**SacrificeDate** date of mouse sacrifice - technical covariate

**ProteinPrepDate** date of protein preparation - technical covariate ...

---

```
generate_colors_for_numeric
```

*Generates color list*

---

**Description**

Generates a list of colors for a vector of numeric, POSIXct (i.e. the (signed) number of seconds since the beginning of 1970), or factors

**Usage**

```
generate_colors_for_numeric(num_col, palette_type = "brewer", i = 1,
  granularity = 10)
```

**Arguments**

num_col	a vector of type numeric or factor to generate colors for
palette_type	'brewer' or 'viridis'
i	if palette_type is 'brewer' the palette argument to brewer_pal. If palette_type is 'viridis' the option argument to viridis_pal()
granularity	the breaks to use when generating colors for num_col

**Value**

list, containing the following items:

1. 'color\_vector' - string-like vector of colors
2. 'new\_annotation' - factor representation of numeric vector (factor with number of levels equal to "granularity")

---

get_peptide_corr_df	<i>Transform square correlation matrix into long data frame of correlations</i>
---------------------	---

---

### Description

Transform square correlation matrix into long data frame of correlations

### Usage

```
get_peptide_corr_df(peptide_cor, peptide_annotation,
  protein_col = "ProteinName", feature_id_col = "peptide_group_label")
```

### Arguments

feature\_id\_col

---

get_sample_corr_distrib	<i>Calculates correlation distribution for all pairs of the replicated samples</i>
-------------------------	--

---

### Description

Calculates correlation distribution for all pairs of the replicated samples

### Usage

```
get_sample_corr_distrib(cor_proteome, sample_annotation,
  sample_id_col = "sample_id", biospecimen_id_col = "EarTag",
  batch_col = "batch")
```

### Arguments

cor_proteome	sample correlation matrix (square)
sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>1. sample_id_col (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
batch_col	column in 'sample_annotation' that should be used for batch comparison

**Value**

dataframe with the following columns, that are suggested to use for plotting in [plot\\_sample\\_corr\\_distribution](#) as plot\_param:

1. replicate
2. batch\_the\_same
3. batch\_replicate
4. batches

other columns are:

1. sample\_id\_1 & sample\_id\_2, both generated from sample\_id\_col variable
2. correlation - correlation of two corresponding samples
3. batch\_1 & batch\_2 or analogous, created the same as sample\_id\_1

---

join_data_matrices	<i>Join data matrices</i>
--------------------	---------------------------

---

**Description**

Joins 2 or more data matrices

**Usage**

```
join_data_matrices(matrices_list, step, sample_annotation,
                   measure_col = "Intensity")
```

**Arguments**

matrices_list	list of matrices in data_matrix ( <a href="#">proBatch</a> ) format to be joined
step	normalization step (e.g. 'Raw' or 'Quantile_normalized' or 'qNorm_ComBat'). Useful if consecutive steps are compared in plots. Note that in plots these are usually ordered alphabetically, so it's worth naming with numbers, e.g. '1_raw', '2_quantile'
sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>1. sample_id_col (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
measure_col	if df_long is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency

**Value**

df\_long ([proBatch](#)) like data frame with a row, having entries for:

1. feature\_id\_col (e.g. peptide name)
2. sample\_id\_col (e.g. filename)
3. measure\_col (e.g. intensity/expression)
4. step (e.g. 'raw', 'quantile\_norm')

**See Also**

Other matrix manipulation functions: [convert\\_to\\_matrix](#), [matrix\\_to\\_long](#)

---

matrix_to_long	<i>Wide to long conversion</i>
----------------	--------------------------------

---

**Description**

Convert from wide matrix to a long data frame representation

**Usage**

```
matrix_to_long(data_matrix, sample_annotation,
               feature_id_col = "peptide_group_label", measure_col = "Intensity",
               sample_id_col = "FullRunName", step = NA)
```

**Arguments**

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. Usually the log transformed version of the original data
sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>1. sample_id_col (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
measure_col	if df_long is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
step	normalization step (e.g. 'Raw' or 'Quantile_normalized' or 'qNorm_ComBat'). Useful if consecutive steps are compared in plots. Note that in plots these are usually ordered alphabetically, so it's worth naming with numbers, e.g. '1_raw', '2_quantile'

**Value**

df\_long ([proBatch](#)) like data frame

**See Also**

Other matrix manipulation functions: [convert\\_to\\_matrix](#), [join\\_data\\_matrices](#)

---

merge_rare_levels	<i>Replaces rare levels with other</i>
-------------------	--

---

### Description

Replaces levels with a maximal occurrence of 1 with other

### Usage

```
merge_rare_levels(column)
```

---

normalize	<i>Data normalization and batch adjustment methods</i>
-----------	--

---

### Description

Data normalization and batch adjustment methods

Median normalization of the data (per batch median)

Median normalization of the data (global)

normalize with the custom (continuous) fit

Standardized input-output ComBat normalization ComBat allows users to adjust for batch effects in datasets where the batch covariate is known, using methodology described in Johnson et al. 2007. It uses either parametric or non-parametric empirical Bayes frameworks for adjusting data for batch effects. Users are returned an expression matrix that has been corrected for batch effects. The input data are assumed to be cleaned and normalized before batch effect removal.

### Usage

```
normalize_medians_batch(data_long, sample_annotation = NULL,
  sample_id_col = "FullRunName", batch_col = "MS_batch.final",
  feature_id_col = "peptide_group_label", measure_col = "Intensity")
```

```
normalize_medians_global(data_long, sample_id_col = "FullRunName",
  measure_col = "Intensity")
```

```
normalize_custom_fit(data_matrix, sample_annotation,
  batch_col = "MS_batch.final", feature_id_col = "peptide_group_label",
  sample_id_col = "FullRunName", measure_col = "Intensity",
  sample_order_col = "order", fit_func = fit_nonlinear, ...)
```

```
correct_with_ComBat(data_matrix, sample_annotation,
  batch_col = "MS_batch.final", par.prior = TRUE)
```

**Arguments**

sample_annotation	data matrix with 1) 'sample_id_col' (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
batch_col	column in 'sample_annotation' that should be used for batch comparison
measure_col	if 'df_long' is among the parameters, it is the column with expression/abundance/intensity, otherwise, it is used internally for consistency
data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. Usually the log transformed version of the original data
sample_order_col	column, determining the order of sample MS run, used as covariate to fit the non-linear fit
fit_func	function to fit the (non)-linear trend
...	other parameters, usually those of the 'fit_func'
par.prior	
return_long	whether the result should be the "long" data frame (as 'df_long') or "wide" (as 'data_matrix')

**Value**

'data\_matrix'-size data matrix with batch-effect corrected by 'ComBat'

**See Also**

[fit\\_nonlinear](#)

---

plot\_clustering

*cluster the data matrix to visually inspect which confounder dominates*

---

**Description**

cluster the data matrix to visually inspect which confounder dominates

**Usage**

```
plot_clustering(data_matrix, color_df, distance = "euclidean",
  agglomeration = "complete", label_samples = T, label_font = 0.2,
  plot_title = NULL, ...)
```



**Arguments**

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. in most function, it is assumed that this is the log transformed version of the original data
color_df	data frame of colors, as created by 'sample_annotation_to_colors'
distance	distance metric used for clustering
agglomeration	agglomeration methods as used by 'hclust'
label_samples	if TRUE sample IDs (column names of data_matrix) will be printed
label_font	size of the font. Is active if label_samples is TRUE, ignored otherwise
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
...	other parameters of 'plotDendroAndColors' from 'WGCNA' package

**See Also**

[hclust](#), [sample\\_annotation\\_to\\_colors](#), [plotDendroAndColors](#)

---

plot_corr_matrix	<i>Visualise correlation matrix</i>
------------------	-------------------------------------

---

**Description**

Plot correlation of selected samples or peptides

**Usage**

```
plot_corr_matrix(corr_matrix, flavor = "corrplot", filename = NULL,
  width = NA, height = NA, unit = c("cm", "in", "mm"), plot_title = "",
  ...)
```

**Arguments**

corr_matrix	square correlation matrix
flavor	either corrplot from 'corrplot' package or heatmap, as in 'pheatmap'
filename	path where the results are saved. If null the object is returned to the active window; otherwise, the object is save into the file. Currently only pdf and png format is supported
width	option determining the output image width
height	option determining the output image width
unit	units: 'cm', 'in' or 'mm'
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
...	parameters for the <a href="#">corrplot.mixed</a> or <a href="#">pheatmap</a> visualisation, for details see examples and help to corresponding functions

**See Also**

[pheatmap](#), [corrplot.mixed](#)

---

plot_heatmap	<i>Plot the heatmap of samples</i>
--------------	------------------------------------

---

## Description

Plot the heatmap of samples

## Usage

```
plot_heatmap(data_matrix, sample_annotation = NULL, fill_the_missing = T,
  cluster_rows = T, cluster_cols = F, annotation_color_list = NA,
  heatmap_color = colorRampPalette(rev(RColorBrewer::brewer.pal(n = 7, name =
    "RdYlBu")))(100), color_for_missing = "black", filename = NA,
  plot_title = NA, ...)
```

## Arguments

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. in most function, it is assumed that this is the log transformed version of the original data
sample_annotation	data matrix with <ol style="list-style-type: none"> <li>1. 'sample_id_col' (this can be repeated as row names)</li> <li>2. biological and</li> <li>3. technical covariates (batches etc)</li> </ol> ; each column of sample annotation will get it's own row. If 'cluster_cols = T' this will indicate, whether sample proximity is driven by one of biological or technical factors
fill_the_missing	boolean value determining if missing values should be substituted with -1 (and colored with black)
cluster_rows	boolean value determining if rows should be clustered
cluster_cols	boolean value determining if columns should be clustered
annotation_color_list	list specifying colors for columns (samples). Best created by 'sample_annotation_to_colors'
heatmap_color	vector of colors used in heatmap (typicall a gradient)
color_for_missing	special color to make missing values. Usually black or white, depending on 'heatmap_color'
filename	filepath where to save the image
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
...	other parameters of link[pheatmap]{pheatmap}

## Value

object returned by link[pheatmap]{pheatmap}

**See Also**

[sample\\_annotation\\_to\\_colors](#), [pheatmap](#)

---

plot_iRTs	<i>Plot iRT measurements</i>
-----------	------------------------------

---

**Description**

Creates a iRT faceted ggplot2 plot of the value in `measure_col` vs `order_col` using [plot\\_peptide\\_level](#). Additionally, the resulting plot can also be faceted by batch.

**Usage**

```
plot_iRTs(df_long, sample_annotation, order_col = NULL, irt_pattern = "iRT",
          batch_col = "MS_batch.final", sample_id_col = "FullRunName",
          feature_id_col = "peptide_group_label", measure_col = "Intensity",
          plot_title = "iRT peptide profile", ...)
```

**Arguments**

<code>df_long</code>	data frame where each row is a single feature in a single sample. It minimally has a <code>sample_id_col</code> , a <code>feature_id_col</code> and a <code>measure_col</code> , but usually also an <code>m_score</code> (in OpenSWATH output result file)
<code>sample_annotation</code>	data matrix with: <ol style="list-style-type: none"> <li>1. <code>sample_id_col</code> (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
<code>order_col</code>	column in <code>sample_annotation</code> that determines sample order. It is used for certain diagnostics and normalisations.
<code>irt_pattern</code>	substring used to identify irts proteins in the column 'ProteinName'
<code>batch_col</code>	column in <code>sample_annotation</code> that should be used for batch comparison
<code>sample_id_col</code>	name of the column in <code>sample_annotation</code> file, where the filenames (colnames of the data matrix are found)
<code>feature_id_col</code>	name of the column with feature/gene/peptide/protein ID used in the long format representation <code>df_long</code> . In the wide formatted representation <code>data_matrix</code> this corresponds to the row names.
<code>measure_col</code>	if <code>df_long</code> is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
<code>plot_title</code>	the string indicating the source of the peptides
<code>...</code>	additional arguments to <a href="#">plot_peptide_level</a> function

**Value**

ggplot2 type plot of `measure_col` vs `order_col`, faceted by `irt_pattern` containing proteins and (optionally) by `batch_col`

**See Also**

Other feature-level diagnostic functions: [plot\\_peptide\\_level](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_spike\\_ins](#), [plot\\_with\\_fitting\\_curve](#)

---

plot_pca	<i>plot PCA plot</i>
----------	----------------------

---

**Description**

plot PCA plot

**Usage**

```
plot_pca(data_matrix, sample_annotation,
         feature_id_col = "peptide_group_label", color_by = "MS_batch",
         PC_to_plot = c(1, 2), fill_the_missing = 0, colors_for_factor = NULL,
         theme = "classic", plot_title = NULL)
```

**Arguments**

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. in most function, it is assumed that this is the log transformed version of the original data
sample_annotation	data matrix with 1) 'sample_id_col' (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
color_by	column name (as in 'sample_annotation') to color by
PC_to_plot	principal component numbers for x and y axis
fill_the_missing	boolean value determining if missing values should be substituted with -1 (and colored with black)
colors_for_factor	named vector of colors for the 'color_by' variable
theme	ggplot theme, by default 'classic'. Can be easily overridden (see examples)
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))

**Value**

ggplot scatterplot colored by factor levels of column specified in 'factor\_to\_color'

**See Also**

[autoplot.pca\\_common](#), [ggplot](#)

---

plot\_peptides\_of\_one\_protein

*Plot peptides of one protein*


---

## Description

Creates a spike-in faceted ggplot2 plot of the value in `measure_col` vs `order_col` using [plot\\_peptide\\_level](#). Additionally, the resulting plot can also be faceted by `batch`.

## Usage

```
plot_peptides_of_one_protein(proteinName, protein_col = "ProteinName",
  df_long, sample_annotation, peptide_annotation = NULL,
  order_col = "order", sample_id_col = "FullRunName",
  batch_col = "MS_batch", measure_col = "Intensity",
  feature_id_col = "peptide_group_label",
  plot_title = sprintf("Peptides of %s protein", proteinName), ...)
```

## Arguments

<code>proteinName</code>	name of the protein as defined in <code>ProteinName</code>
<code>protein_col</code>	column where protein names are specified
<code>df_long</code>	data frame where each row is a single feature in a single sample. It minimally has a <code>sample_id_col</code> , a <code>feature_id_col</code> and a <code>measure_col</code> , but usually also an <code>m_score</code> (in OpenSWATH output result file)
<code>sample_annotation</code>	data matrix with: <ol style="list-style-type: none"> <li>1. <code>sample_id_col</code> (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
<code>order_col</code>	column in <code>sample_annotation</code> that determines sample order. It is used for certain diagnostics and normalisations.
<code>sample_id_col</code>	name of the column in <code>sample_annotation</code> file, where the filenames (colnames of the data matrix are found)
<code>batch_col</code>	column in <code>sample_annotation</code> that should be used for batch comparison
<code>measure_col</code>	if <code>df_long</code> is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
<code>feature_id_col</code>	name of the column with feature/gene/peptide/protein ID used in the long format representation <code>df_long</code> . In the wide formatted representation <code>data_matrix</code> this corresponds to the row names.
<code>plot_title</code>	the string indicating the source of the peptides
<code>...</code>	additional arguments to <a href="#">plot_peptide_level</a> function

## Value

ggplot2 type plot of `measure_col` vs `order_col`, faceted by `spike_ins` containing proteins and (optionally) by `batch_col`

**See Also**

Other feature-level diagnostic functions: [plot\\_iRTs](#), [plot\\_peptide\\_level](#), [plot\\_spike\\_ins](#), [plot\\_with\\_fitting\\_curve](#)

---

plot\_peptide\_correlation\_distr\_one\_protein

*Plot distribution curves for each step to see if they shift by data normalization*

---

**Description**

Plot distribution curves for each step to see if they shift by data normalization

**Usage**

```
plot_peptide_correlation_distr_one_protein(data_matrix_list, protein_name,
  peptide_annotation, protein_col = "ProteinName",
  feature_id_col = "peptide_group_label",
  plot_title = sprintf("Distribution of peptide correlation at different correction steps,\n
  protein_name), theme = "classic")
```

**Arguments**

protein\_name     name of the protein, as specified in protein\_col of peptide\_annotation  
 theme

---

plot\_peptide\_level     *Plot peptide measurements*

---

**Description**

Creates a peptide faceted ggplot2 plot of the value in measure\_col vs order\_col. Additionally, the resulting plot can also be faceted by batch.

**Usage**

```
plot_peptide_level(pep_name, df_long, sample_annotation, order_col = NULL,
  sample_id_col = "FullRunName", batch_col = "MS_batch.final",
  measure_col = "Intensity", feature_id_col = "peptide_group_label",
  geom = c("point", "line"), color_by_batch = F, facet_by_batch = F,
  plot_title = NULL, requant = NULL, theme = "classic")
```

**Arguments**

pep_name	name of the peptide for diagnostic profiling
df_long	data frame where each row is a single feature in a single sample. It minimally has a sample_id_col, a feature_id_col and a measure_col, but usually also an m_score (in OpenSWATH output result file)
sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>1. sample_id_col (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
order_col	column in sample_annotation that determines sample order. It is used for certain diagnostics and normalisations.
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
batch_col	column in sample_annotation that should be used for batch comparison
measure_col	if df_long is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
geom	whether to show the feature as points and/or connect by lines
color_by_batch	(logical) whether to color points by batch
facet_by_batch	(logical) whether to plot each batch in its own facet
plot_title	the string indicating the source of the peptides
requant	if data frame: requant values; if logical: whether to indicate requant values (requires 'requant' or 'm_score' column in df_long)
theme	plot theme (default is 'classical'; other options not implemented)

**Value**

ggplot2 type plot of measure\_col vs order\_col, faceted by pep\_name and (optionally) by batch\_col

**See Also**

Other feature-level diagnostic functions: [plot\\_iRTs](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_spike\\_ins](#), [plot\\_with\\_fitting\\_curve](#)

---

plot\_protein\_corrplot *Peptide correlation matrix (heatmap)*

---

**Description**

Plots correlation plot of peptides from a single protein

**Usage**

```
plot_protein_corrplot(data_matrix, protein_name, peptide_annotation,
  protein_col = "ProteinName", peptide_col_name = "peptide_group_label",
  flavor = "corrplot", filename = NULL, width = NA, height = NA,
  unit = c("cm", "in", "mm"), plot_title = "peptide correlation matrix",
  ...)
```

**Arguments**

<code>data_matrix</code>	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. Usually the log transformed version of the original data
<code>protein_name</code>	the name of the protein
<code>peptide_annotation</code>	df with peptides and their corresponding proteins
<code>protein_col</code>	the column name in peptide_annotation with protein names
<code>peptide_col_name</code>	the column name in peptide_annotation with peptide names
<code>flavor</code>	either corrplot from 'corrplot' package or heatmap, as in 'pheatmap'
<code>filename</code>	path where the results are saved. If null the object is returned to the active window; otherwise, the object is save into the file. Currently only pdf and png format is supported
<code>width</code>	option determining the output image width
<code>height</code>	option determining the output image width
<code>unit</code>	units: 'cm', 'in' or 'mm'
<code>plot_title</code>	The title of the plot
<code>...</code>	parameters for the corrplot visualisation

**Examples**

```
\donotrunc{plot_corr_plot(q_norm_proteome, protein_name = 'Hao',
  peptide_annotation = peptide_annotation, prot.column = 'Gene',
  title = 'Hao protein peptides after quantile norm',
  number.cex=0.75, tl.cex = .75
  mar=c(0,0,1,0))}
\donotrunc{lower = "ellipse", upper = "number",
  tl.col = "black", diag = '1', tl.pos = "lt", number.cex=0.75, tl.cex = .75}
```

---

plot\_prot\_corr\_distribution

*Plot distribution of peptide correlations within one protein and between proteins*

---

**Description**

Plot distribution of peptide correlations within one protein and between proteins



**Usage**

```
plot_prot_corr_distribution(data_matrix, peptide_annotation,
  protein_col = "ProteinName", feature_id_col = "peptide_group_label",
  plot_title = "Distribution of peptide correlation", theme = "classic")
```

**Arguments**

theme

---

plot_pvca	<i>Plot variance distribution by variable</i>
-----------	---

---

**Description**

Plot variance distribution by variable

**Usage**

```
plot_pvca(data_matrix, sample_annotation, sample_id_col = "FullRunName",
  feature_id_col = "peptide_group_label",
  technical_covariates = c("MS_batch", "instrument"),
  biological_covariates = c("cell_line", "drug_dose"), fill_the_missing = 0,
  threshold_pca = 0.6, threshold_var = 0.01, colors_forBars = NULL,
  theme = "classic", plot_title = NULL)
```

**Arguments**

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. in most function, it is assumed that this is the log transformed version of the original data
sample_annotation	data matrix with 1) 'sample_id_col' (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
technical_covariates	vector 'sample_annotation' column names that are technical covariates
biological_covariates	vector 'sample_annotation' column names, that are biologically meaningful covariates
threshold_pca	the percentile value of the minimum amount of the variabilities that the selected principal components need to explain
threshold_var	the percentile value of weight each of the covariates needs to explain (the rest will be lumped together)

colors_forBars	four-item color vector, specifying colors for the following categories: c('residual', 'biological', 'biol:techn', 'technical')
theme	ggplot theme, by default 'classic'. Can be easily overridden (see examples)
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.

**Value**

list of two items: plot = gg, df = pvca\_res

**See Also**

[sample\\_annotation\\_to\\_colors](#), [ggplot](#)

---

plot\_samples\_corrplot *Sample correlation matrix (heatmap)*

---

**Description**

Plot correlation of selected samples

**Usage**

```
plot_samples_corrplot(data_matrix, samples_to_plot = NULL,
  flavor = "corrplot", filename = NULL, width = NA, height = NA,
  unit = c("cm", "in", "mm"), plot_title = "Correlation matrix of samples",
  ...)
```

**Arguments**

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. Usually the log transformed version of the original data
samples_to_plot	string vector of samples in data_matrix to be used in the plot
filename	path where the results are saved. If null the object is returned to the active window; otherwise, the object is save into the file. Currently only pdf and png format is supported
width	option determining the output image width
height	option determining the output image width
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
...	parameters for the <a href="#">corrplot.mixed</a> or <a href="#">pheatmap</a> visualisation, for details see examples and help to corresponding functions

**See Also**

[pheatmap](#), [corrplot.mixed](#)

**Examples**

```
### Example 1: Plot heatmap of pre-specified samples
#'\dontrun{specified_samples = sample_annotation %>%
  filter(RunID %in% paste('Run', 110:115, sep = '')) %>%
  pull(FullRunName)
plot_samples_corr_heatmap(data_matrix, sample_to_plot = specified_samples,
  flavor = 'pheatmap', cluster_rows = F, cluster_cols = F)

}

### Example 2: Plot corrplot of pre-specified samples
#'\dontrun{specified_samples = sample_annotation %>%
  filter(RunID %in% paste('Run', 110:115, sep = '')) %>%
  pull(FullRunName)
plot_samples_corr_heatmap(data_matrix, sample_to_plot = specified_samples,
  flavor = 'corrplot', lower = "ellipse", upper = "number",
  tl.col = "black", diag = '1', tl.pos = "lt", number.cex=0.75, tl.cex = .75)

}
```

---

plot\_sample\_corr\_distribution

*Create violin plot of correlation distribution*

---

**Description**

Useful to visualize within batch vs within replicate vs non-related sample correlation

**Usage**

```
plot_sample_corr_distribution(data_matrix, sample_annotation,
  repeated_samples = NULL, sample_id_col = "sample_id",
  batch_col = "batch", biospecimen_id_col = "EarTag",
  plot_title = "Correlation distribution", plot_param = "batch_replicate")
```

**Arguments**

repeated_samples	if 'NULL', only repeated sample correlation is plotted
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
batch_col	column in sample_annotation that should be used for batch comparison
plot_title	Title of the plot (usually, processing step + representation level (fragments, transcriptions, proteins))
plot_param	columns, defined in correlation_df, which is output of get_sample_corr_distrib, specifically, # <ol style="list-style-type: none"> <li>1. replicate</li> <li>2. batch_the_same</li> </ol>

```

3. batch_replicate
4. batches
;

```

**Value**

ggplot type object with violin plot for each plot\_param

**See Also**

[get\\_sample\\_corr\\_distrib](#), [ggplot](#)

---

plot\_sample\_means\_or\_boxplots

*Plot per-sample average or boxplot (distribution) vs order (if the real running order available)*

---

**Description**

Plot per-sample average or boxplot (distribution) vs order (if the real running order available)

**Usage**

```

plot_sample_mean(data_matrix, sample_annotation = NULL,
  sample_id_col = "FullRunName", order_col = "order", batch_col = NULL,
  facet_col = "instrument", color_by_batch = F, color_scheme = "brewer",
  theme = "classic", plot_title = NULL, order_per_facet = F)

```

```

plot_boxplot(df_long, sample_annotation = NULL,
  sample_id_col = "FullRunName", measure_col = "Intensity",
  order_col = "order", batch_col = "MS_batch.final",
  facet_col = "instrument", color_by_batch = T, color_scheme = "brewer",
  theme = "classic", plot_title = NULL, order_per_facet = F)

```

**Arguments**

data_matrix	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. in most function, it is assumed that this is the log transformed version of the original data
sample_annotation	data matrix with 1) 'sample_id_col' (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
order_col	column where running order is specified.
batch_col	column in 'sample_annotation' that should be used for batch comparison
facet_col	recommended if more than one batch covariate is present. Faceting is most suited to examine instruments separately
color_by_batch	should the each batch be represented with its own color?

color_scheme	named vector, names corresponding to unique batch values as specified in 'sample_annotation'
theme	ggplot theme, by default 'classic'. Can be easily overridden (see examples)
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
order_per_facet	if order is defined ignoring facets (usually instrument), re-define order per-batch
df_long	data frame where each row is a single feature in a single sample, thus it has minimally, 'sample_id_col', 'feature_id_col' and 'measure_col', but usually also 'm_score' (in OpenSWATH output result file)
measure_col	if 'df_long' is among the parameters, it is the column with expression/abundance/intensity, otherwise, it is used internally for consistency

### Details

functions for quick visual assessment of trends associated, overall or specific covariate-associated (see 'batch\_col' and 'facet\_col')

### Value

ggplot2 class object. Thus, all aesthetics can be overridden

### See Also

[ggplot](#)

---

plot_spike_ins	<i>Plot spike-in measurements</i>
----------------	-----------------------------------

---

### Description

Creates a spike-in faceted ggplot2 plot of the value in measure\_col vs order\_col using [plot\\_peptide\\_level](#). Additionally, the resulting plot can also be faceted by batch.

### Usage

```
plot_spike_ins(df_long, sample_annotation, order_col = "order",
  spike_ins = "BOVIN", sample_id_col = "FullRunName",
  batch_col = "MS_batch", measure_col = "Intensity",
  feature_id_col = "peptide_group_label",
  plot_title = "Spike-in BOVINE protein peptides", ...)
```

### Arguments

df_long	data frame where each row is a single feature in a single sample. It minimally has a sample_id_col, a feature_id_col and a measure_col, but usually also an m_score (in OpenSWATH output result file)
sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>sample_id_col (this can be repeated as row names)</li> </ol>

	2. biological covariates
	3. technical covariates (batches etc)
order_col	column in sample_annotation that determines sample order. It is used for certain diagnostics and normalisations.
spike_ins	substring used to identify spike-in proteins in the column 'ProteinName'
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
batch_col	column in sample_annotation that should be used for batch comparison
measure_col	if df_long is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
plot_title	the string indicating the source of the peptides
...	additional arguments to <a href="#">plot_peptide_level</a> function

**Value**

ggplot2 type plot of measure\_col vs order\_col, faceted by spike\_ins containing proteins and (optionally) by batch\_col

**See Also**

Other feature-level diagnostic functions: [plot\\_iRTs](#), [plot\\_peptide\\_level](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_with\\_fitting\\_curve](#)

---

plot\_within\_prot\_corr\_distribution

*Plot distribution of median correlations of peptides within same protein*

---

**Description**

Plot distribution of median correlations of peptides within same protein

**Usage**

```
plot_within_prot_corr_distribution(data_matrix_list, peptide_annotation,
  protein_col = "ProteinName", feature_id_col = "peptide_group_label",
  plot_title = "Distribution of peptide correlation", theme = "classic")
```

**Arguments**

theme

---

plot\_with\_fitting\_curve

*Plot peptide measurements across multi-step analysis*


---

## Description

Plot Intensity of a few representative peptides for each step of the analysis including the fitting curve

## Usage

```
plot_with_fitting_curve(pep_name, data_df_all_steps, sample_annotation, fit_df,
  fit_value_var = "fit", fit_step = "3_loess_fit", order_col = NULL,
  sample_id_col = "FullRunName", batch_col = "MS_batch",
  measure_col = "Intensity", feature_id_col = "peptide_group_label",
  geom = c("point", "line"), color_by_batch = F, facet_by_batch = F,
  plot_title = NULL, requant = NULL, theme = "classic")
```

## Arguments

pep_name	name of the peptide for diagnostic profiling
data_df_all_steps	data frame, similar to df_long <a href="#">proBatch</a> , where each row is a single feature in a single sample, at a certain step of the analysis (minimally raw and after linear normalization) thus it has minimally the following columns: sample_id_col, feature_id_col, measure_col, and fit_step, but usually also m_score
sample_annotation	data matrix with: <ol style="list-style-type: none"> <li>1. sample_id_col (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
order_col	column in sample_annotation that determines sample order. It is used for certain diagnostics and normalisations.
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
batch_col	column in sample_annotation that should be used for batch comparison
measure_col	if df_long is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.
geom	for the intensity measure_col profile:
color_by_batch	(logical) whether to color points by batch
facet_by_batch	(logical) whether to plot each batch in its own facet
plot_title	the string indicating the source of the peptides
requant	if data frame: requant values; if logical: whether to indicate requant values (requires 'requant' or 'm_score' column in df_long)
theme	plot theme (default is 'classical'; other options not implemented)

**Value**

ggplot-class plot with minimally two facets (before and after non-linear fit) with `measure_col` (Intensity) vs `order_col` (injection order) for selected peptides (specified in `pep_name`)

**See Also**

Other feature-level diagnostic functions: [plot\\_iRTs](#), [plot\\_peptide\\_level](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_spike\\_ins](#)

---

proBatch	<i>proBatch: A package for diagnostics and correction of batch effects, primarily in proteomics</i>
----------	---

---

**Description**

It addresses the following needs:

- prepare the original data (e.g. OpenSWATH output matrix and sample annotation file) for analysis. However, you might need to use ‘SWATH2stats’ additionally
- Diagnose batch effects, sample-wide and feature-level
- Correct for batch effects (normalize the data). Other useful package for this purpose is ‘Normalizer’.

**Arguments**

<code>df_long</code>	data frame where each row is a single feature in a single sample. It minimally has a <code>sample_id_col</code> , a <code>feature_id_col</code> and a <code>measure_col</code> , but usually also an <code>m_score</code> (in OpenSWATH output result file)
<code>data_matrix</code>	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. Usually the log transformed version of the original data
<code>sample_annotation</code>	data matrix with: <ol style="list-style-type: none"> <li>1. <code>sample_id_col</code> (this can be repeated as row names)</li> <li>2. biological covariates</li> <li>3. technical covariates (batches etc)</li> </ol>
<code>sample_id_col</code>	name of the column in <code>sample_annotation</code> file, where the filenames (colnames of the data matrix are found)
<code>batch_col</code>	column in <code>sample_annotation</code> that should be used for batch comparison
<code>order_col</code>	column in <code>sample_annotation</code> that determines sample order. It is used for certain diagnostics and normalisations.
<code>measure_col</code>	if <code>df_long</code> is among the parameters, it is the column with expression/abundance/intensity; otherwise, it is used internally for consistency
<code>feature_id_col</code>	name of the column with feature/gene/peptide/protein ID used in the long format representation <code>df_long</code> . In the wide formatted representation <code>data_matrix</code> this corresponds to the row names.
<code>plot_title</code>	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
<code>theme</code>	ggplot theme, by default ‘classic’. Can be easily overridden (see examples)



**Details**

To learn more about proBatch, start with the vignettes: `'browseVignettes(package = "proBatch")'`

---

quantile_normalize	<i>Quantile normalization of the data, ensuring that the row and column names are retained</i>
--------------------	--

---

**Description**

Quantile normalization of the data, ensuring that the row and column names are retained

**Usage**

```
quantile_normalize(data_matrix)
```

**Arguments**

data\_matrix      log transformed data matrix (features in rows and samples in columns)

**Value**

'data\_matrix'-size matrix, with columns quantile-normalized

---

remove_peptides_with_missing_batch	<i>Remove features missing in at least one batch</i>
------------------------------------	--

---

**Description**

Cleans dataset df\_long ([proBatch](#)) by removing all features that are not present in every batch

**Usage**

```
remove_peptides_with_missing_batch(df_long, batch_col = "MS_batch.final",
  feature_id_col = "peptide_group_label")
```

**Arguments**

df_long	data frame where each row is a single feature in a single sample. It minimally has a sample_id_col, a feature_id_col and a measure_col, but usually also an m_score (in OpenSWATH output result file)
batch_col	column in sample_annotation that should be used for batch comparison
feature_id_col	name of the column with feature/gene/peptide/protein ID used in the long format representation df_long. In the wide formatted representation data_matrix this corresponds to the row names.

**Details**

useful for some downstream functions as ComBat normalization, that would not work otherwise

**Value**

df\_long ([proBatch](#)) like data frame freed of features that were not detected in each batch

**See Also**

Other dataset cleaning functions: [clean\\_requants](#), [summarize\\_peptides](#)

---

sample\_annotation\_to\_colors

*Generate colors for sample annotation*

---

**Description**

Convert the sample annotation data frame to list of colors the list is named as columns included to use in potting functions

**Usage**

```
sample_annotation_to_colors(sample_annotation, columns_for_plotting = NULL,
                             sample_id_col = "FullRunName", factor_columns = c("subtype",
                                         "caseControl"), not_factor_columns = c("RunDate", "ProteinPrepDate"),
                             rare_categories_to_other = T, numerics_to_log = F,
                             numeric_palette_type = "brewer", granularity = 10)
```

**Arguments**

sample\_annotation

data matrix with:

1. sample\_id\_col (this can be repeated as row names)
2. biological covariates
3. technical covariates (batches etc)

columns\_for\_plotting

only consider these columns from sample\_annotation

sample\_id\_col name of the column in sample\_annotation file, where the filenames (colnames of the data matrix are found)

factor\_columns columns of sample\_annotation to be treated as factors. Note that factor and character columns are treated as factors by default.

not\_factor\_columns

don't treat these columns as factors. This can be used to override the default behaviour of considering factors and character columns as factors.

rare\_categories\_to\_other

if True rare categories will be merged as 'other'

numerics\_to\_log

NOT IMPLEMENTED!

numeric\_palette\_type

palette to be used for numeric values coloring

granularity

number of colors to map to the number vector (equally spaced between minimum and maximum)

**Value**

list of colors

---

sample\_random\_peptides

*sample random peptides for diagnostics*

---

**Description**

sample random peptides for diagnostics

**Usage**

```
sample_random_peptides(proteome, seed = 1, pep_per_group = 3,
  groups_RT = 10, groups_intensity = 5)
```

**Arguments**

proteome	required columns:
	1. m_score
	2. Intensity
	3. peptide_group_label
	4. RT
pep_per_group	number of peptides to sample per group
groups_intensity	

---

summarize_peptides	<i>Summarize run features</i>
--------------------	-------------------------------

---

**Description**

Summarizes various peptide properties on a per sample basis. By default will summarize RT, Intensity and m\_score. If your feature does not have some of these set them to NULL when calling.

**Usage**

```
summarize_peptides(df_long, sample_id_col = "FullRunName",
  feature_id_col = "peptide_group_label", RT = "RT",
  measure_col = "Intensity", m_score = "m_score")
```

**Details**

summarize peptides by sample (ranking) and on the contrary, across peptide-wise across samples

**Value**

a data frame summarizing features in a dataset on a per sample basis. The following columns are returned: 'RT\_mean', 'Int\_mean', 'numb\_requants', 'median\_m\_score', 'mean\_m\_score', 'median\_good\_m\_score' (median of 'm\_score' excluding requants)

**See Also**

Other dataset cleaning functions: [clean\\_requants](#), [remove\\_peptides\\_with\\_missing\\_batch](#)

# Index

## \*Topic **datasets**

example\_peptide\_annotation, 9  
example\_proteome, 10  
example\_sample\_annotation1, 10

as.POSIXct, 7

autoplot.pca\_common, 20

boxplot\_all\_steps, 3

clean\_requants, 4, 34, 36

color\_list\_to\_df, 5

convert\_to\_matrix, 5, 14

correct\_with\_ComBat (normalize), 15

corrplot.mixed, 17, 26, 27

create\_peptide\_annotation, 6

date\_to\_sample\_order, 7

dates\_to\_posix, 6

define\_batches\_by\_MS\_pauses, 8

define\_batches\_by\_MS\_pauses\_within\_instrument,  
8

example\_peptide\_annotation, 9

example\_proteome, 10

example\_sample\_annotation1, 10

fit\_nonlinear, 16

generate\_colors\_for\_numeric, 11

get\_peptide\_corr\_df, 12

get\_sample\_corr\_distrib, 12, 28

ggplot, 20, 26, 28, 29

hclust, 17

join\_data\_matrices, 6, 13, 14

matrix\_to\_long, 6, 14, 14

merge\_rare\_levels, 15

normalize, 15

normalize\_custom\_fit (normalize), 15

normalize\_medians\_batch (normalize), 15

normalize\_medians\_global (normalize), 15

pheatmap, 17, 19, 26, 27

plot\_boxplot, 3

plot\_boxplot  
(plot\_sample\_means\_or\_boxplots),  
28

plot\_clustering, 16

plot\_corr\_matrix, 17

plot\_corrplot\_protein, 6

plot\_heatmap, 18

plot\_iRTs, 7, 19, 22, 23, 30, 32

plot\_pca, 20

plot\_peptide\_correlation\_distr\_one\_protein,  
22

plot\_peptide\_level, 19–22, 22, 29, 30, 32

plot\_peptides\_of\_one\_protein, 6, 20, 21,  
23, 30, 32

plot\_prot\_corr\_distribution, 24

plot\_protein\_corrplot, 23

plot\_pvca, 25

plot\_sample\_corr\_distribution, 13, 27

plot\_sample\_mean, 7

plot\_sample\_mean  
(plot\_sample\_means\_or\_boxplots),  
28

plot\_sample\_means\_or\_boxplots, 28

plot\_samples\_corrplot, 26

plot\_spike\_ins, 20, 22, 23, 29, 32

plot\_with\_fitting\_curve, 20, 22, 23, 30,  
31

plot\_within\_prot\_corr\_distribution, 30

plot\_within\_prot\_distribution, 6

plotDendroAndColors, 17

proBatch, 4, 5, 13, 14, 31, 32, 33, 34

proBatch-package (proBatch), 32

quantile\_normalize, 33

remove\_peptides\_with\_missing\_batch, 4,  
33, 36

sample\_annotation\_to\_colors, 17, 19, 26,  
34

sample\_random\_peptides, 35

summarize\_peptides, 4, 34, 35