

# Package ‘proBatch’

September 18, 2019

**Type** Package

**Title** Tools for Diagnostics and Corrections of Batch Effects in Proteomics

**Version** 1.0.0

**Author** Jelena Cuklina <chuklina.jelena@gmail.com>, Chloe H. Lee <chloe.h.lee94@gmail.com>, Patrick Pedrioli <pedrioli@gmail.com>

**Maintainer** Chloe H. Lee <chloe.h.lee94@gmail.com>

**Description** The proBatch package facilitates batch effects analysis and correction in high-throughput experiments. It was developed primarily for mass-spectrometry proteomics (DIA/SWATH), but could also be applicable to most omic data with minor adaptations. The package contains functions for diagnostics (proteome/genome-wide and feature-level), correction (normalization and batch effects correction) and quality control. Non-linear fitting based approaches were also included to deal with complex, mass spectrometry-specific signal drifts.

**biocViews** BatchEffect, Normalization, Preprocessing, Software, MassSpectrometry, Proteomics, QualityControl, GeneExpression

**License** GPL-3

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

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

**Depends** R (>= 3.6)

**Encoding** UTF-8

**LazyData** true

**Imports** Biobase, corrplot, dplyr, data.table, ggfortify, ggplot2, grDevices, lazyeval, lubridate, magrittr, pheatmap, preprocessCore, purrr, pvca, RColorBrewer, readr, reshape2, rlang, scales, viridis, stats, sva, tidyr, tibble, utils, wesanderson, WGCNA

**Suggests** knitr, rmarkdown, devtools, roxygen2, testthat

**VignetteBuilder** knitr

**RoxygenNote** 6.1.1

**git\_url** <https://git.bioconductor.org/packages/proBatch>

**git\_branch** RELEASE\_3\_9

**git\_last\_commit** 6640beb

**git\_last\_commit\_date** 2019-05-02

**Date/Publication** 2019-09-17

## R topics documented:

adjust_batch_trend . . . . .	3
center_peptide_batch_medians . . . . .	4
correct_batch_effects . . . . .	5
correct_with_ComBat . . . . .	6
create_peptide_annotation . . . . .	7
dates_to_posix . . . . .	8
date_to_sample_order . . . . .	8
example_peptide_annotation . . . . .	9
example_proteome . . . . .	10
example_proteome_matrix . . . . .	10
example_sample_annotation . . . . .	11
log_transform . . . . .	12
long_to_matrix . . . . .	12
matrix_to_long . . . . .	13
normalize . . . . .	14
normalize_data . . . . .	14
normalize_sample_medians . . . . .	15
plot_heatmap . . . . .	15
plot_hierarchical_clustering . . . . .	17
plot_iRT . . . . .	18
plot_PCA . . . . .	19
plot_peptides_of_one_protein . . . . .	20
plot_peptide_corr_distribution . . . . .	22
plot_protein_corrplot . . . . .	23
plot_PVCA . . . . .	24
plot_sample_corr_distribution . . . . .	25
plot_sample_corr_heatmap . . . . .	26
plot_sample_mean_or_boxplot . . . . .	27
plot_single_feature . . . . .	29
plot_spike_in . . . . .	30
plot_with_fitting_curve . . . . .	32
proBatch . . . . .	33
quantile_normalize . . . . .	34
sample_annotation_to_colors . . . . .	35
sample_color_scheme . . . . .	36

---

adjust\_batch\_trend      *adjust batch signal trend with the custom (continuous) fit*

---

### Description

adjust batch signal trend with the custom (continuous) fit

### Usage

```
adjust_batch_trend(data_matrix, sample_annotation,
  batch_col = "MS_batch", feature_id_col = "peptide_group_label",
  sample_id_col = "FullRunName", measure_col = "Intensity",
  sample_order_col = "order", fit_func = fit_nonlinear,
  abs_threshold = 5, pct_threshold = 0.2, ...)
```

### 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 frame with sample ID, technical (e.g. MS batches) and biological (e.g. Diet) covariates
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.
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
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_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
abs_threshold	the absolute threshold to filter data for curve fitting
pct_threshold	the percentage threshold to filter data for curve fitting
...	other parameters, usually those of the fit_func

### Value

list of two items: 1) data\_matrix, adjusted with continuous fit; 2) fit\_df, used to examine the fitting curves

### See Also

[fit\\_nonlinear](#)

**Examples**

```
trend_corrected_matrix <- adjust_batch_trend(example_proteome_matrix,
example_sample_annotation, span = 0.7,
abs_threshold = 5, pct_threshold = 0.20)
```

---

center\_peptide\_batch\_medians

*Median centering of the peptides (per batch median)*

---

**Description**

Median centering of the peptides (per batch median)

**Usage**

```
center_peptide_batch_medians(df_long, sample_annotation = NULL,
sample_id_col = "FullRunName", batch_col = "MS_batch",
feature_id_col = "peptide_group_label", measure_col = "Intensity")
```

**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 frame with sample ID, technical (e.g. MS batches) and biological (e.g. Diet) covariates
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
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

**Value**

df\_long-size long format data with batch-effect corrected with per-feature batch median centering in Intensity\_normalized column

**Examples**

```
median_centered_proteome <- center_peptide_batch_medians(
example_proteome, example_sample_annotation)
```

---

correct\_batch\_effects *Batch correction method allows correction of continuous sigal drift within batch and discrete difference across batches.*

---

### Description

Batch correction method allows correction of continuous sigal drift within batch and discrete difference across batches.

### Usage

```
correct_batch_effects(data_matrix, sample_annotation,
  fitFunc = "loess_regression", discreteFunc = c("MedianCentering",
  "ComBat"), batch_col = "MS_batch",
  feature_id_col = "peptide_group_label",
  sample_id_col = "FullRunName", measure_col = "Intensity",
  sample_order_col = "order", abs_threshold = 5, pct_threshold = 0.2,
  ...)
```

### 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 frame with sample ID, technical (e.g. MS batches) and biological (e.g. Diet) covariates
fitFunc	function to use for the fit (currently only loess_regression available)
discreteFunc	function to use for discrete batch correction (MedianCentering or ComBat)
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.
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
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_order_col	column, determining the order of sample MS run, used as covariate to fit the non-linear fit
abs_threshold	the absolute threshold to filter data for curve fitting
pct_threshold	the percentage threshold to filter data for curve fitting
...	other parameters, usually of normalize_custom_fit, and fit_func

### Value

data\_matrix-size data matrix with batch-effect corrected by fit and discrete functions

**Examples**

```
batch_corrected_matrix <- correct_batch_effects(
  example_proteome_matrix, example_sample_annotation,
  discreteFunc = 'MedianCentering',
  batch_col = 'MS_batch',
  span = 0.7,
  abs_threshold = 5, pct_threshold = 0.20)
```

---

correct\_with\_ComBat     *Adjusts for discrete batch effects using ComBat*

---

**Description**

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**

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

**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 frame with sample ID, technical (e.g. MS batches) and biological (e.g. Diet) covariates
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
par.prior	whether parametrical or non-parametrical prior should be used

**Value**

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

**Examples**

```
combat_corrected_matrix <- correct_with_ComBat(
  example_proteome_matrix, example_sample_annotation)
```

---

`create_peptide_annotation`

*Prepare peptide annotation from long format data frame Create light-weight peptide annotation data frame for selection of illustrative proteins*

---

## Description

Prepare peptide annotation from long format data frame

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

## Usage

```
create_peptide_annotation(df_long,  
  feature_id_col = "peptide_group_label",  
  annotation_col = c("ProteinName"))
```

## 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.

`annotation_col` one or more columns containing protein ID

## Value

data frame containing peptide annotations

## See Also

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

## Examples

```
generated_peptide_annotation <- create_peptide_annotation(  
  example_proteome, feature_id_col = "peptide_group_label",  
  annotation_col = c("ProteinName" ))
```

---

dates\_to\_posix      *Convert data/time to POSIXct*

---

### 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 = "DateTime", 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 a new column new\_time\_column with POSIX-formatted date

### Examples

```
date_to_posix <- dates_to_posix(example_sample_annotation,
  time_column = c('RunDate', 'RunTime'),
  new_time_column = 'DateTime',
  dateTimeFormat = c("%b_%d", "%H:%M:%S"))
```

---

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"),
  new_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

new\_order\_col name of column with generated the order of sample run based on time columns

instrument\_col column, denoting different instrument used for measurements

**Value**

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

**Examples**

```
sample_annotation_wOrder <- date_to_sample_order(
  example_sample_annotation,
  time_column = c('RunDate', 'RunTime'),
  new_time_column = 'new_DateTime',
  dateTimeFormat = c("%b_%d", "%H:%M:%S"),
  new_order_col = 'new_order',
  instrument_col = NULL)
```

---

example\_peptide\_annotation

*Peptide annotation data*

---

**Description**

This is data from Aging study annotated with gene names

**Usage**

```
example_peptide_annotation
```

**Format**

A data frame with 535 rows and 10 variables:

**peptide\_group\_label** 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

---

example_proteome	<i>Example protein data in long format</i>
------------------	--

---

**Description**

This is data from 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 124655 rows and 5 variables:

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

**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

**Gene** conventional gene name of corresponding ProteinName

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

**Source**

PRIDE ID will be added in future

---

example_proteome_matrix	<i>Example protein data in matrix</i>
-------------------------	---------------------------------------

---

**Description**

This is measurement data from Aging study with columns representing samples and rows representing peptides

**Usage**

example\_proteome\_matrix

**Format**

A matrix with 534 rows and 233 columns:

**Source**

PRIDE ID will be added in future

---

example\_sample\_annotation

*Sample annotation data version 1*

---

**Description**

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

**Usage**

example\_sample\_annotation

**Format**

A data frame with 233 rows and 11 variables:

**FullRunName** name of the file, in most functions used for `sample_id_col`

**MS\_batch** 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 - "

**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

**DateTime** numeric date and time generated by `date_to_sample_order`

**order** order of samples generated by sorting `DateTime` in `date_to_sample_order`

**digestion\_batch** peptide digestion batch: 5-level factor of manually annotated batches ...

---

log_transform	<i>Log transformation of the data, ensuring that the row and column names are retained</i>
---------------	--

---

**Description**

Log transformation of the data, ensuring that the row and column names are retained

**Usage**

```
log_transform(data_matrix, log_base = 2)
```

**Arguments**

data_matrix	raw data matrix (features in rows and samples in columns)
log_base	base of the logarithm for transformation

**Value**

data\_matrix-size matrix, with columns log2 transformed

**Examples**

```
log_transformed_matrix <- log_transform(example_proteome_matrix)
```

---

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

---

**Description**

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

**Usage**

```
long_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: [matrix\\_to\\_long](#)

**Examples**

```
proteome_matrix <- long_to_matrix(example_proteome)
```

---

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 = NULL,
  feature_id_col = "peptide_group_label", measure_col = "Intensity",
  sample_id_col = "FullRunName", step = NULL)
```

**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: [long\\_to\\_matrix](#)

**Examples**

```
proteome_long <- matrix_to_long(example_proteome_matrix,
example_sample_annotation)
```

---

normalize

*Data normalization methods*

---

**Description**

Data normalization methods

**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_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
measure_col	if 'df_long' is among the parameters, it is the column with expression/abundance/intensity, otherwise, it is used internally for consistency

---

normalize\_data

*Normalization brings the samples to the same scale*

---

**Description**

Normalization brings the samples to the same scale

**Usage**

```
normalize_data(data_matrix, normalizeFunc = c("quantile",
"medianCentering"), log_base = NULL)
```

**Arguments**

data_matrix	raw data matrix (features in rows and samples in columns)
normalizeFunc	global batch normalization method ('quantile' or 'MedianCentering')
log_base	whether to log transform data matrix before normalization ('NULL', '2' or '10')

**Value**

data\_matrix-size matrix, with columns normalized

**Examples**

```
quantile_normalized_matrix <- normalize_data(example_proteome_matrix,  
normalizeFunc = "quantile", log_base = 2)
```

---

normalize\_sample\_medians

*Normalization by centering sample medians to global median of the data*

---

**Description**

Normalization by centering sample medians to global median of the data

**Usage**

```
normalize_sample_medians(df_long, sample_id_col = "FullRunName",  
measure_col = "Intensity")
```

**Arguments**

df_long	log transformed long format data matrix (see 'df_long')
sample_id_col	name of the column in sample_annotation file, where the filenames (colnames of the data matrix are found)
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'-size matrix, with intensity scaled to global median

**Examples**

```
median_normalized_matrix <- normalize_sample_medians(example_proteome)
```

---

plot\_heatmap

*Plot the heatmap of samples*

---

**Description**

Plot the heatmap of samples

**Usage**

```
plot_heatmap(data_matrix, sample_annotation = NULL,
             sample_id_col = "FullRunName", sample_annotation_col = NULL,
             sample_annotation_row = NULL, fill_the_missing = TRUE,
             cluster_rows = TRUE, cluster_cols = FALSE,
             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

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

; 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

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

**sample\_annotation\_col** biological or technical factors to be annotated in heatmap columns

**sample\_annotation\_row** biological or technical factors to be annotated in heatmap rows

**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](#)

**Examples**

```
color_scheme <- sample_annotation_to_colors (example_sample_annotation,
factor_columns = c('MS_batch', 'EarTag', "Strain",
"Diet", "digestion_batch", "Sex"),
not_factor_columns = 'DateTime',
numeric_columns = c('order'))
```

```
heatmap_plot <- plot_heatmap(example_proteome_matrix,
example_sample_annotation,
sample_annotation_col = c("MS_batch", "digestion_batch", "Diet"),
cluster_cols = TRUE,
annotation_color_list = color_scheme$list_of_colors,
show_rownames = FALSE, show_colnames = FALSE)
```

---

plot\_hierarchical\_clustering

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

---

**Description**

cluster the data matrix to visually inspect which confounder dominates

**Usage**

```
plot_hierarchical_clustering(data_matrix, color_df,
distance = "euclidean", agglomeration = "complete",
label_samples = TRUE, 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

**Value**

No return

**See Also**

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

**Examples**

```
color_scheme <- sample_annotation_to_colors (example_sample_annotation,
factor_columns = c('MS_batch', 'EarTag', "Strain", "Diet", "digestion_batch", "Sex"),
not_factor_columns = 'DateTime',
numeric_columns = c('order'))

color_annotation <- color_scheme$color_df

hierarchical_clustering_plot <- plot_hierarchical_clustering(
example_proteome_matrix, color_annotation,
distance = "euclidean", agglomeration = 'complete',
label_samples = FALSE)
```

---

plot\_iRT

*Plot iRT measurements*

---

**Description**

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

**Usage**

```
plot_iRT(df_long, sample_annotation, peptide_annotation = NULL,
protein_col = "ProteinName", order_col = "order",
irt_pattern = "iRT", sample_id_col = "FullRunName",
batch_col = "MS_batch", measure_col = "Intensity",
feature_id_col = "peptide_group_label", color_by_batch = FALSE,
color_scheme = "brewer", facet_by_batch = FALSE,
color_by_col = NULL, color_by_value = NULL,
plot_title = "iRT peptide profile", ...)
```

**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:

1. `sample_id_col` (this can be repeated as row names)
2. biological covariates
3. technical covariates (batches etc)

`peptide_annotation` long format data with peptide ID and their corresponding protein annotations

`protein_col` column where protein names are specified

order_col	column in sample_annotation that determines sample order. It is used for certain diagnostics and normalisations.
irt_pattern	substring used to identify irts 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.
color_by_batch	(logical) whether to color points by batch
color_scheme	color scheme for ggplot representation
facet_by_batch	(logical) whether to plot each batch in its own facet
color_by_col	column to color by certain value denoted by color_by_value
color_by_value	value in color_by_col to color
plot_title	the string indicating the source of the peptides
...	additional arguments to <a href="#">plot_single_feature</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\\_peptides\\_of\\_one\\_protein](#), [plot\\_single\\_feature](#), [plot\\_spike\\_in](#), [plot\\_with\\_fitting\\_curve](#)

**Examples**

```
irt_plot <- plot_iRT(example_proteome,
  example_sample_annotation,
  protein_col = 'Gene', irt_pattern = "BOVINE_A1ag")
```

---

plot\_PCA

*plot PCA plot*


---

**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**

<code>data_matrix</code>	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
<code>sample_annotation</code>	data matrix with 1) <code>sample_id_col</code> (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
<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>color_by</code>	column name (as in <code>sample_annotation</code> ) to color by
<code>PC_to_plot</code>	principal component numbers for x and y axis
<code>fill_the_missing</code>	boolean value determining if missing values should be substituted with -1 (and colored with black)
<code>colors_for_factor</code>	named vector of colors for the <code>color_by</code> variable
<code>theme</code>	ggplot theme, by default <code>classic</code> . Can be easily overridden (see examples)
<code>plot_title</code>	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](#)

**Examples**

```
pca_plot <- plot_PCA(example_proteome_matrix, example_sample_annotation,
  color_by = 'MS_batch', plot_title = "PCA colored by MS batch")
```

---

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\\_single\\_feature](#). Additionally, the resulting plot can also be faceted by batch.

**Usage**

```
plot_peptides_of_one_protein(protein_name, 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", color_by_batch = FALSE,
  color_scheme = "brewer", facet_by_batch = FALSE,
  color_by_col = NULL, color_by_value = NULL,
  plot_title = sprintf("Peptides of %s protein", protein_name), ...)
```

**Arguments**

protein_name	name of the protein as defined in ProteinName
protein_col	column where protein names are specified
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>
peptide_annotation	long format data with peptide ID and their corresponding protein annotations
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.
color_by_batch	(logical) whether to color points by batch
color_scheme	color scheme for ggplot representation
facet_by_batch	(logical) whether to plot each batch in its own facet
color_by_col	column to color by certain value denoted by color_by_value
color_by_value	value in color_by_col to color
plot_title	the string indicating the source of the peptides
...	additional arguments to <a href="#">plot_single_feature</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\\_iRT](#), [plot\\_single\\_feature](#), [plot\\_spike\\_in](#), [plot\\_with\\_fitting\\_curve](#)

**Examples**

```
peptides_of_one_protein_plot <- plot_peptides_of_one_protein (
  protein_name = "Hao",
  protein_col = "Gene", df_long = example_proteome,
  example_sample_annotation,
  order_col = 'order', sample_id_col = 'FullRunName',
  batch_col = 'MS_batch')
```

---

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

**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>peptide_annotation</code>	long format data with peptide ID and their corresponding protein annotations
<code>protein_col</code>	the column name in <code>peptide_annotation</code> with protein names
<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
<code>theme</code>	ggplot theme, by default <code>classic</code> . Can be easily overridden
<code>...</code>	parameters for the ggplot visualisation

**Value**

ggplot type object with violin plot for each `plot_param`

**Examples**

```
peptide_corr_distribution <- plot_peptide_corr_distribution(
  example_proteome_matrix,
  example_peptide_annotation, protein_col = 'Gene')
```

---

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", feature_id_col = "peptide_group_label",
  flavor = c("pheatmap", "corrplot"), filename = NULL, width = NA,
  height = NA, unit = c("cm", "in", "mm"),
  plot_title = sprintf("Peptide correlation matrix of %s protein",
  protein_name), ...)
```

**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
protein_name	the name of the protein
peptide_annotation	df with peptides and their corresponding proteins
protein_col	the column name in peptide_annotation with protein names
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.
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	The title of the plot
...	parameters for the corrplot visualisation

**Value**

corrplot or pheatmap object depending on flavor

**Examples**

```
protein_corrplot_plot <- plot_protein_corrplot(example_proteome_matrix, protein_name = 'Haao',
  peptide_annotation = example_peptide_annotation,
  protein_col = 'Gene', flavor = "pheatmap")
```

plot\_PVCA

*Plot variance distribution by variable***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

**fill\_the\_missing** numeric value that the missing values are substituted with

**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))

**Value**

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

**See Also**

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

**Examples**

```
matrix <- example_proteome_matrix[1:50, ]
pvca_plot <- plot_PVCA(matrix, example_sample_annotation,
  technical_covariates = c('MS_batch', 'digestion_batch'),
  biological_covariates = c("Diet", "Sex", "Strain"))
```

---

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 = "FullRunName",
  batch_col = "MS_batch", biospecimen_id_col = "EarTag",
  plot_title = "Correlation distribution",
  plot_param = "batch_replicate")
```

**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 1) sample_id_col (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
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

biospecimen_id_col	column in sample_annotation that captures the biological sample, that (possibly) was profiled several times as technical replicates. Tip: if such ID is absent, but can be defined from several columns, create new biospecimen_id column
plot_title	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
plot_param	columns, defined in correlation_df, which is output of get_sample_corr_distrib, specifically, #' <ul style="list-style-type: none"> <li>1. replicate</li> <li>2. batch_the_same</li> <li>3. batch_replicate</li> <li>4. batches</li> </ul> ;

**Value**

ggplot type object with violin plot for each plot\_param

**See Also**

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

**Examples**

```
sample_corr_distribution_plot <- plot_sample_corr_distribution(
  example_proteome_matrix,
  example_sample_annotation, batch_col = 'MS_batch',
  biospecimen_id_col = "EarTag",
  plot_param = 'batch_replicate')
```

---

plot\_sample\_corr\_heatmap

*Sample correlation matrix (heatmap)*

---

**Description**

Plot correlation of selected samples

**Usage**

```
plot_sample_corr_heatmap(data_matrix, samples_to_plot = NULL,
  flavor = c("pheatmap", "corrplot"), filename = NULL, width = NA,
  height = NA, unit = c("cm", "in", "mm"),
  plot_title = sprintf("Correlation matrix of sample %s",
  samples_to_plot), ...)
```

**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>samples_to_plot</code>	string vector of samples in <code>data_matrix</code> to be used in the plot
<code>flavor</code>	either <code>corrplot</code> 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>	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
<code>...</code>	parameters for the <a href="#">corrplot.mixed</a> or <a href="#">pheatmap</a> visualisation, for details see examples and help to corresponding functions

**Value**

corrplot or pheatmap object depending on flavor

**See Also**

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

**Examples**

```
specified_samples = example_sample_annotation$FullRunName[
  which(example_sample_annotation$order %in% 110:115)]

sample_corr_heatmap <- plot_sample_corr_heatmap(example_proteome_matrix,
  samples_to_plot = specified_samples,
  flavor = 'pheatmap',
  cluster_rows= FALSE, cluster_cols=FALSE,
  annotation_names_col = TRUE, annotation_legend = FALSE,
  show_colnames = FALSE)
```

---

plot\_sample\_mean\_or\_boxplot

*Plot per-sample mean or boxplot (showing median and quantiles) vs order (if the real running order available)*

---

**Description**

Plot per-sample mean or boxplot (showing median and quantiles) 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 = "MS_batch", facet_col = NULL, color_by_batch = FALSE,
  color_scheme = "brewer", theme = "classic", plot_title = NULL,
  order_per_facet = FALSE, vline_color = "grey", ylimits = NULL)
```

```
plot_boxplot(df_long, sample_annotation = NULL,
  sample_id_col = "FullRunName", measure_col = "Intensity",
  order_col = "order", batch_col = "MS_batch", facet_col = NULL,
  color_by_batch = TRUE, color_scheme = "brewer", theme = "classic",
  plot_title = NULL, order_per_facet = FALSE)
```

**Arguments**

<code>data_matrix</code>	features (in rows) vs samples (in columns) matrix, with feature IDs in rownames and file/sample names as colnames. in most function,
<code>sample_annotation</code>	data matrix with 1) <code>sample_id_col</code> (this can be repeated as row names) 2) biological and 3) technical covariates (batches etc)
<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>order_col</code>	column where running order is specified.
<code>batch_col</code>	column in <code>sample_annotation</code> that should be used for batch comparison
<code>facet_col</code>	recommended if more than one batch covariate is present. Faceting is most suited to examine instruments separately
<code>color_by_batch</code>	should the each batch be represented with its own color?
<code>color_scheme</code>	named vector, names corresponding to unique batch values as specified in <code>sample_annotation</code>
<code>theme</code>	ggplot theme, by default <code>classic</code> . Can be easily overridden (see examples)
<code>plot_title</code>	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
<code>order_per_facet</code>	if order is defined ignoring facets (usually instrument), re-define order per-batch
<code>vline_color</code>	color of vertical lines, typically denoting different MS batches in ordered runs; should be <code>NULL</code> for experiments without intrinsic order
<code>ylimits</code>	range of y-axis to plot feature-level trends
<code>df_long</code>	data frame where each row is a single feature in a single sample, thus it has minimally, <code>sample_id_col</code> , <code>feature_id_col</code> and <code>measure_col</code> , but usually also <code>m_score</code> (in OpenSWATH output result file)
<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

**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](#)**Examples**

```
mean_plot <- plot_sample_mean(example_proteome_matrix, example_sample_annotation,
  order_col = 'order', batch_col = "MS_batch")
```

```
boxplot <- plot_boxplot(example_proteome, example_sample_annotation,
  batch_col = "MS_batch")
```

---

plot\_single\_feature     *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_single_feature(pep_name, df_long, sample_annotation,
  order_col = "order", sample_id_col = "FullRunName",
  batch_col = "MS_batch", measure_col = "Intensity",
  feature_id_col = "peptide_group_label", geom = c("point", "line"),
  color_by_batch = FALSE, color_scheme = "brewer",
  facet_by_batch = FALSE, color_by_col = NULL, color_by_value = NULL,
  plot_title = NULL, vline_color = "red", 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
color_scheme	color scheme for ggplot representation
facet_by_batch	(logical) whether to plot each batch in its own facet
color_by_col	column to color by certain value denoted by color_by_value
color_by_value	value in color_by_col to color
plot_title	the string indicating the source of the peptides
vline_color	color of vertical lines, typically denoting different MS batches in ordered runs; should be NULL for experiments without intrinsic order
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\\_iRT](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_spike\\_in](#), [plot\\_with\\_fitting\\_curve](#)

**Examples**

```
single_feature_plot <- plot_single_feature(
  pep_name = "46213_NVGVSFYADKPEVTQEYK_2",
  df_long = example_proteome, example_sample_annotation,
  color_by_col = NULL)
```

---

plot_spike_in	<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\\_single\\_feature](#). Additionally, the resulting plot can also be faceted by batch.

**Usage**

```
plot_spike_in(df_long, sample_annotation, peptide_annotation = NULL,
  protein_col = "ProteinName", order_col = "order",
  spike_ins = "BOVIN", sample_id_col = "FullRunName",
  batch_col = "MS_batch", measure_col = "Intensity",
  feature_id_col = "peptide_group_label", color_by_batch = FALSE,
  color_scheme = "brewer", facet_by_batch = FALSE,
  color_by_col = NULL, color_by_value = NULL,
  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>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>
peptide_annotation	long format data with peptide ID and their corresponding protein annotations
protein_col	column where protein names are specified
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.
color_by_batch	(logical) whether to color points by batch
color_scheme	color scheme for ggplot representation
facet_by_batch	(logical) whether to plot each batch in its own facet
color_by_col	column to color by certain value denoted by color_by_value
color_by_value	value in color_by_col to color
plot_title	the string indicating the source of the peptides
...	additional arguments to <a href="#">plot_single_feature</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\\_iRT](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_single\\_feature](#), [plot\\_with\\_fitting\\_curve](#)

**Examples**

```
spike_in_plot <- plot_spike_in(example_proteome, example_sample_annotation,
  protein_col = 'Gene', spike_ins = "BOVINE_A1ag",
  plot_title = "Spike-in BOVINE protein peptides")
```

---

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, df_long, sample_annotation, fit_df,
  fit_value_var = "fit", order_col = "order",
  sample_id_col = "FullRunName", batch_col = "MS_batch",
  measure_col = "Intensity", feature_id_col = "peptide_group_label",
  geom = c("point", "line"), color_by_batch = FALSE,
  color_scheme = "brewer", facet_by_batch = FALSE,
  plot_title = sprintf("Fitting curve of %s peptide", pep_name),
  color_by_col = NULL, color_by_value = NULL, theme = "classic",
  vline_color = "grey", ...)
```

## 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>
fit_df	data frame typically output generated from nonlinear curve fitting by <code>normalize_custom_fit</code>
fit_value_var	column denoting intensity values, typically fitted to curve
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
color_scheme	color scheme for ggplot representation



facet_by_batch	(logical) whether to plot each batch in its own facet
plot_title	the string indicating the source of the peptides
color_by_col	column to color by certain value denoted by color_by_value
color_by_value	value in color_by_col to color
theme	plot theme (default is 'classical'; other options not implemented)
vline_color	color of vertical lines, typically denoting different MS batches in ordered runs; should be NULL for experiments without intrinsic order
...	additional arguments to <a href="#">plot_single_feature</a> function

### 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\\_iRT](#), [plot\\_peptides\\_of\\_one\\_protein](#), [plot\\_single\\_feature](#), [plot\\_spike\\_in](#)

### Examples

```
loess_fit_70 <- adjust_batch_trend(example_proteome_matrix,
example_sample_annotation, span = 0.7)
```

```
fitting_curve_plot <- plot_with_fitting_curve(
pep_name = "10231_QDVDVWLWQQEGSSK_2",
df_long = example_proteome, example_sample_annotation,
fit_df = loess_fit_70$fit_df, plot_title = "Curve fitting with 70% span")
```

---

proBatch

*proBatch: A package for diagnostics and correction of batch effects, primarily in proteomics*

---

### Description

The proBatch package contains functions for analyzing and correcting batch effects and other unwanted technical variation from high-throughput experiments. Although the package has primarily been developed for mass spectrometry proteomics (DIA/SWATH), it should also be applicable to most omic data with minor adaptations. It addresses the following needs:

- prepare the data for analysis
- Visualize batch effects in sample-wide and feature-level;
- Normalize and correct for batch effects.

**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)
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>
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
order_col	column in sample_annotation that determines sample order. It is used for certain diagnostics and normalisations.
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	Title of the plot (usually, processing step + representation level (fragments, transitions, proteins))
theme	ggplot theme, by default classic. Can be easily overridden

**Details**

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

**Section**

Common arguments to the functions.

---

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

**Examples**

```
quantile_normalized_matrix <- quantile_normalize(example_proteome_matrix)
```

---

```
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("MS_batch", "EarTag", "Strain", "Diet", "Sex"),
  not_factor_columns = "DateTime", numeric_columns = "order",
  rare_categories_to_other = TRUE, 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.

numeric\_columns

columns of sample\_annotation to be treated as continuous numeric values.

rare\_categories\_to\_other

if True rare categories will be merged as 'other'

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

**Examples**

```
color_scheme <- sample_annotation_to_colors (example_sample_annotation,  
factor_columns = c('MS_batch', 'EarTag', "Strain",  
"Diet", "digestion_batch", "Sex"),  
not_factor_columns = 'DateTime',  
numeric_columns = c('order'))
```

---

sample\_color\_scheme    *Sample color annotation*

---

**Description**

This is an color scheme generated from example sample annotation

**Usage**

```
sample_color_scheme
```

**Format**

A list of 3 components: list\_of\_colors, color\_df and sample\_annotation

**list\_of\_colors** a list of colors for 11 variables, including MS\_batch, EarTag, Strain, Diet, digestion\_batch, Sex, FullRunName, RunDate, RunTime, DateTime, order

**color\_df** a data frame with 233 samples and 11 variables describing a color for each component

**sample\_annotaion** a data frame containing 233 samples and 11 variables annotating samples to facilitate conversion to a color scheme

# Index

## \*Topic **datasets**

- example\_peptide\_annotation, 9
  - example\_proteome, 10
  - example\_proteome\_matrix, 10
  - example\_sample\_annotation, 11
  - sample\_color\_scheme, 36
- adjust\_batch\_trend, 3
- as.POSIXct, 8, 9
- autoplot.pca\_common, 20
- center\_peptide\_batch\_medians, 4
- correct\_batch\_effects, 5
- correct\_with\_ComBat, 6
- corrplot.mixed, 27
- create\_peptide\_annotation, 7
- date\_to\_sample\_order, 8
- dates\_to\_posix, 8
- example\_peptide\_annotation, 9
- example\_proteome, 10
- example\_proteome\_matrix, 10
- example\_sample\_annotation, 11
- fit\_nonlinear, 3
- get\_sample\_corr\_distrib, 26
- ggplot, 20, 25, 26, 29
- hclust, 18
- log\_transform, 12
- long\_to\_matrix, 12, 14
- matrix\_to\_long, 13, 13
- normalize, 14
- normalize\_data, 14
- normalize\_sample\_medians, 15
- pheatmap, 17, 27
- plot\_boxplot  
(plot\_sample\_mean\_or\_boxplot),  
27
- plot\_heatmap, 15
- plot\_hierarchical\_clustering, 17
- plot\_iRT, 9, 18, 22, 30, 31, 33
- plot\_PCA, 19
- plot\_peptide\_corr\_distribution, 22
- plot\_peptides\_of\_one\_protein, 7, 19, 20,  
30, 31, 33
- plot\_protein\_corrplot, 7, 23
- plot\_PVCA, 24
- plot\_sample\_corr\_distribution, 25
- plot\_sample\_corr\_heatmap, 26
- plot\_sample\_mean, 9
- plot\_sample\_mean  
(plot\_sample\_mean\_or\_boxplot),  
27
- plot\_sample\_mean\_or\_boxplot, 27
- plot\_single\_feature, 18–22, 29, 30, 31, 33
- plot\_spike\_in, 19, 22, 30, 30, 33
- plot\_with\_fitting\_curve, 19, 22, 30, 31,  
32
- plotDendroAndColors, 18
- proBatch, 13, 33
- proBatch-package (proBatch), 33
- quantile\_normalize, 34
- sample\_annotation\_to\_colors, 17, 18, 25,  
35
- sample\_color\_scheme, 36