

Protein significance analysis of mass spectrometry-based proteomics experiments with R and MSstats (v3.7.3)

Meena Choi & Tsung-Heng Tsai

January 16, 2017

Contents

1. Statistical relative protein quantification: SRM, DDA and DIA experiments	1
1.1 Applicability	2
1.2 Statistical functionalities	2
1.3 Interoperability with existing computational tools	2
1.4 Availability	2
1.5 Overview of the functionalities	4
1.6 Troubleshooting	5
2. Allowable data formats	5
2.1 SRM with stable isotope labeled reference peptides	5
2.2 Label-free DDA	8
2.2 Label-free DIA	9
3. Prerequisites and setting for MSstats analysis	9
4. DDA analysis with MSstats	10
4.1 General workflow for DDA	10
4.2 Suggested workflow with Skyline output for DDA	30
4.3 Suggested workflow with MaxQuant output for DDA	33
4.4 Suggested workflow with Progenesis output for DDA	35
4.5 Suggested workflow with Proteome Discoverer output for DDA	40
5. DIA analysis with MSstats	44
5.1 Suggested workflow with Skyline output for DIA	44
5.2 Suggested workflow with Spectronaut output for DIA	47
5.3 Suggested workflow with OpenSWATH output for SWATH	50
Reference	54

1. Statistical relative protein quantification: SRM, DDA and DIA experiments

MSstats is an open-source R-based package for statistical relative quantification of peptides and proteins in mass spectrometry-based proteomic experiments. This document describes **MSstats**, the most recent version of the package, and its use through the command line.

1.1 Applicability

MSstats version 3.0 and above is applicable to multiple types of sample preparation, including label-free workflows, workflows that use stable isotope labeled reference proteins and peptides, and workflows that use fractionation. It is applicable to targeted Selected Reaction Monitoring (SRM), Data-Dependent Acquisition (DDA or shotgun), and Data-Independent Acquisition (DIA or SWATH-MS). It is applicable to experiments that make arbitrary complex comparisons of experimental conditions or times.

MSstats is currently not applicable to experiments that compare multiple metabolically labeled endogenous samples within a same run. It is not applicable to experiments with iTRAQ labeling. These experiments will be supported in the future.

1.2 Statistical functionalities

MSstats version 3.0 and above performs three analysis steps. The first step, *data processing, visualization, and run-level summarization*, transforms, normalizes and summarizes the intensities of the peaks per MS run and per protein, and generates workflow-specific and customizable numeric summaries for data visualization and quality control.

The second step, *statistical modeling and inference*, automatically detects the experimental design (e.g. group comparison, paired design or time course, presence of labeled reference peptides or proteins) from the data. It then reflects the experimental design and the type of spectral acquisition strategy, and fits an appropriate linear mixed model by means of `lm` and `lmer` functionalities in R. The model is used to detect differentially abundant proteins or peptides, or to summarize the protein or peptide abundance in a single biological replicate or condition (that can be used, e.g. as input to clustering or classification).

The third step, *statistical experimental design*, views the dataset being analyzed as a pilot study of a future experiment, utilizes the variance components of the current dataset, and calculates the minimal number of replicates necessary in the future experiment to achieve a pre-specified statistical power.

1.3 Interoperability with existing computational tools

MSstats takes as input data in a tabular .csv format, which can be generated by any spectral processing tool such as Skyline (MacLean et al. 2010), MaxQuant (Jürgen Cox and Mann 2008), Progenesis QI(Nonlinear dynamics/Waters), Proteome Discoverer (Thermo Scientific) MultiQuant(Applied Biosystems), OpenMS (Sturm et al. 2008), SuperHirn (Mueller et al. 2007), OpenSWATH (Röst et al. 2014) or Spectronaut(Biognosys). The functions to convert the required format from several processing tools are available from **MSstats** v3.6. Details are in the section below.

For statistics experts, **MSstats** 3.0 and above satisfies the interoperability requirements of Bioconductor, and takes as input data in the `MSnSet` format (Gatto and Lilley 2012). The command line-based workflow is partitioned into a series of independent steps, that facilitate the development and testing of alternative statistical approaches. It complies with the maintenance and documentation requirements of Bioconductor.

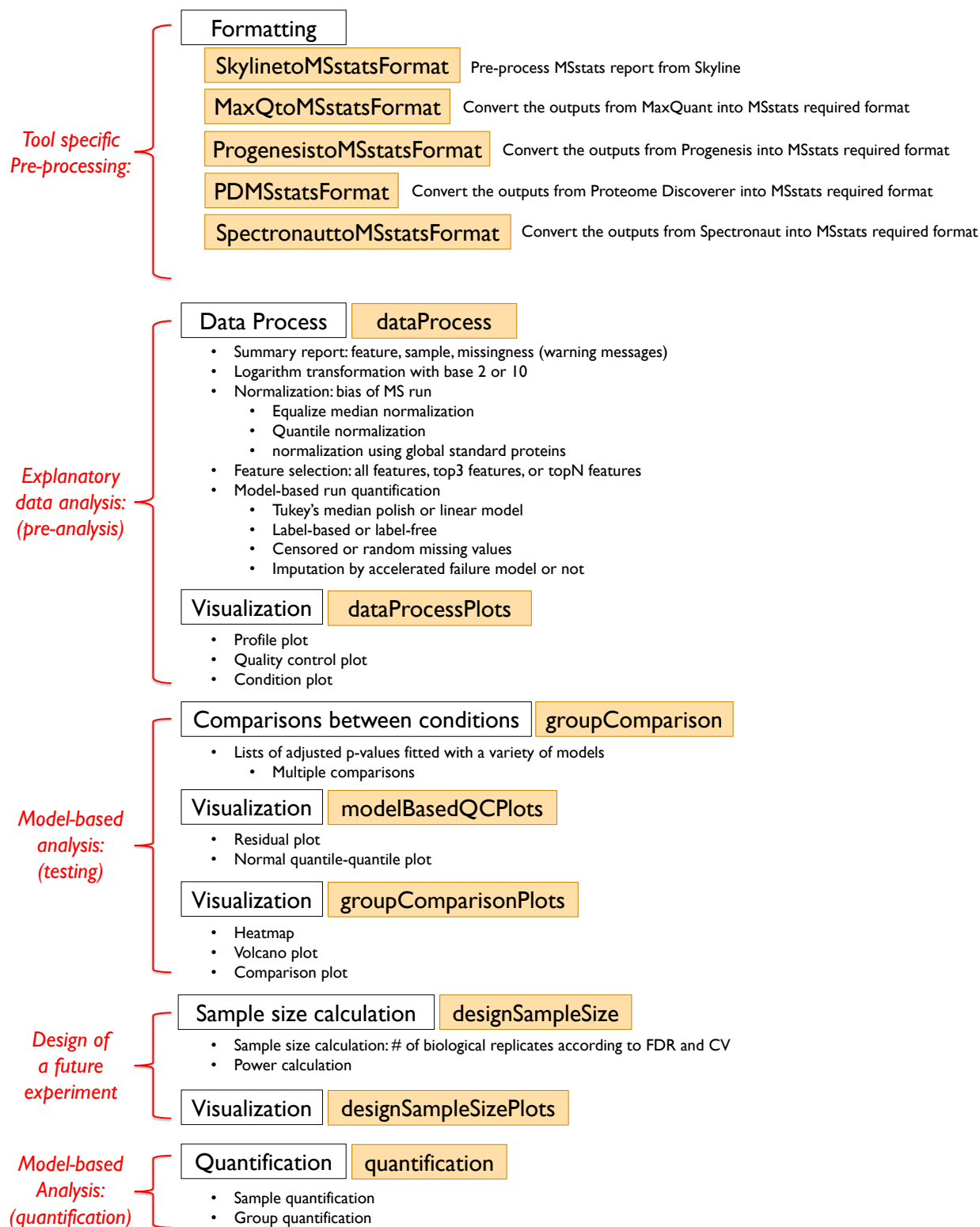
Finally, **MSstats** 3.0 and above is available as an external tool within Skyline. The external tool support within Skyline manages **MSstats** installation, point-and-click execution, parameter collection in Windows forms and output display. Skyline manages the annotations of the experimental design, and the processing of raw data. It outputs a custom report, that is fed as a single stream input into **MSstats**. This design buffers proteomics users from the details of the R implementation, while enabling rigorous statistical modeling.

1.4 Availability

MSstats is available under the Artistic-2.0 license at msstats.org. **MSstats** as an external tool for Skyline is available at <http://proteome.gs.washington.edu/software/Skyline/tools.html>. **MSstats** is now also available

in Bioconductor. The most recent version of the package is available at msstats.org or MSstats GitHub. We suggest to use that if possible. The versioning of the main package is updated several times a year, to synchronise with the Bioconductor release.

1.5 Overview of the functionalities



1.6 Troubleshooting

To help troubleshoot potential problems with installation or functionalities of **MSstats**, a progress report is generated in a separate log file *msstats.log*. The file includes information on the R session (R version, loaded software libraries), options selected by the user, checks of successful completion of intermediate analysis steps, and warning messages. If the analysis produces an error, the file contains suggestions for possible reasons for the errors. If a file with this name already exists in working directory, a suffix with a number will be appended to the file name. In this way a record of all the analyses is kept. Please see the file [KnownIssues-Skyline-MSstatsV3.6.pdf](#) on the “Installation” section of “MSstats” page in [msstats.org](#) for a list of known issues and possible solutions for installation problem of MSstats external tool in Skyline

2. Allowable data formats

2.1 SRM with stable isotope labeled reference peptides

2.1.1 10-column format

MSstats performs statistical analysis steps, that follow peak identification and quantitation. Therefore, input to **MSstats** is the output of other software tools (such as Skyline or MultiQuant) that read raw spectral files and identify and quantify spectral peaks. The preferred structure of data for use in **MSstats** is a .csv file in a “long” format with 10 columns representing the following variables: **ProteinName**, **PeptideSequence**, **PrecursorCharge**, **FragmentIon**, **ProductCharge**, **IsotopeLabelType**, **Condition**, **BioReplicate**, **Run**, **Intensity**. The variable names are fixed, but are case-insensitive.

- (a) **ProteinName**: This column needs information about Protein id. Statistical analysis will be done separately for each unique label in this column. For peptide-level modeling and analysis, use peptide id in this column.
- (b)-(e) **PeptideSequence**, **PrecursorCharge**, **FragmentIon**, **ProductCharge**: The combination of these 4 columns defines a *feature* of a protein (in SRM experiments, it is a transition that is identified and quantified across runs). If the information for one or several of these columns is not available, please do not discard these columns but use a single fixed value across the entire dataset. For example, if the original raw data does not contain the information of **ProductCharge**, assign the value 0 to the entries in the column **ProductCharge** for the entire dataset. If the peptide sequences should be distinguished based on post-translational modifications, this column can be renamed to **PeptideModifiedSequence**. For example, this allows us to use the **PeptideModifiedSequence** column from the Skyline report.
- (f) **IsotopeLabelType**: This column indicates whether this measurement is based on the endogenous peptides (use “L”) or labeled reference peptides (use “H”).
- (g) **Condition**: For group comparison experiments, this column indicates groups of interest (such as “Disease” or “Control”). For time-course experiments, this column indicates time points (such as “T1”, “T2”, etc). If the experimental design contains both distinct groups of subjects and multiple time points per subject, this column should indicate a combination of these values (such as “Disease_T1”, “Disease_T2”, “Control_T1”, “Control_T2”, etc.).
- (h) **BioReplicate**: This column should contain a unique identifier for each biological replicate in the experiment. For example, in a clinical proteomic investigation this should be a unique patient id. Patients from distinct groups should have distinct ids. **MSstats** does not require the presence of technical replicates in the experiment. If the technical replicates are present, all samples or runs from a same biological replicate should have a same id. **MSstats** automatically detects the presence of technical replicates and accounts for them in the model-based analysis.

- (i) **Run:** This column contains the identifier of a mass spectrometry run. Each mass spectrometry run should have a unique identifier, regardless of the origin of the biological sample. In SRM experiments, if all the transitions of a biological or a technical replicate are split into multiple “methods” due to the technical limitations, each method should have a separate identifier. When processed by Skyline, distinct values of runs correspond to distinct input file names. It is possible to use the actual input file names as values in the column **Run**.
- (j) **Intensity:** This column should contain the quantified signal of a feature in a run without any transformation (in particular, no logarithm transform). The signals can be quantified as the peak height or the peak of area under curve. Any other quantitative representation of abundance can also be used.

An example of an acceptable input dataset is shown below. This example dataset is from an SRM experiment with stable isotope labeled reference peptides. The dataset is stored in a .csv file in a “long” format. Each row corresponds to a single intensity. More details on assigning the values of **Condition**, **BioReplicate** and **Run**, depending on the structure of the experimental design, are given below.

	A	B	C	D	E	F	G	H	I	J
1	ProteinName	PeptideSequence	PrecursorCharge	FragmentIon	ProductCharge	IsotopeLabelType	Condition	BioReplicate	Run	Intensity
2	ACEA	EILGHEIFFDWELP	3 y3		0 H		1	ReplA	1	66472.3847
3	ACEA	EILGHEIFFDWELP	3 y3		0 L		1	ReplA	1	5764.16228
4	ACEA	EILGHEIFFDWELP	3 y4		0 H		1	ReplA	1	101005.166
5	ACEA	EILGHEIFFDWELP	3 y4		0 L		1	ReplA	1	61.65238
6	ACEA	EILGHEIFFDWELP	3 y5		0 H		1	ReplA	1	90055.4993
7	ACEA	EILGHEIFFDWELP	3 y5		0 L		1	ReplA	1	472.691803
8	ACEA	TDSEAAATLISSTID'	2 y10		0 H		1	ReplA	1	43506.5425
9	ACEA	TDSEAAATLISSTID'	2 y10		0 L		1	ReplA	1	217.203553
10	ACEA	TDSEAAATLISSTID'	2 y7		0 H		1	ReplA	1	68023.0377
11	ACEA	TDSEAAATLISSTID'	2 y7		0 L		1	ReplA	1	725.284308
12	ACEA	TDSEAAATLISSTID'	2 y8		0 H		1	ReplA	1	68276.0489
13	ACEA	TDSEAAATLISSTID'	2 y8		0 L		1	ReplA	1	243.658527

2.1.2 Assigning the values of Condition, BioReplicate and Run

The values of **Condition**, **BioReplicate**, **Run** depend on the design of the specific experiment.

1) Group comparison

In a group comparison design, the conditions (e.g., disease states) are profiled across **non-overlapping sets of biological replicates** (i.e., subjects). In this example there are 2 conditions, Disease and Control (in general the number of conditions can vary). There are 3 subjects (i.e., biological replicates) per condition (in general an equal number of replicates per condition is not required). Each subject has 2 technical replicate runs (in general technical replicates are not required, and their number per sample may vary). Overall, in this example there are $2 \times 3 \times 2 = 12$ mass spectrometry runs.

Table below shows the values of the columns **Condition**, **BioReplicate** and **Run** for this situation. It is important to note two things. First, the order of subjects and conditions in the experiment should be randomized, and run id does not need to represent the order of spectral acquisition. Second, the values of the columns are repeated for every quantified transition. For example, if in each run the experiment quantifies 50 endogenous transitions and 50 labeled reference counterparts, then the input file has $12 \times 50 \times 2 = 1200$ lines. When a feature intensity is missing in a run, the data structure should contain a separate row for each missing value. The rows should include all the information (from **ProteinName** to **Run**), and indicate missing intensities with NA.

Condition	BioReplicate	Run
Disease	Subject1	1
Disease	Subject1	2
Disease	Subject2	3

Condition	BioReplicate	Run
Disease	Subject2	4
Disease	Subject3	5
Disease	Subject3	6
Control	Subject4	7
Control	Subject4	8
Control	Subject5	9
Control	Subject5	10
Control	Subject6	11
Control	Subject6	12

2) Time course

The important feature of a time course experimental design is that **a same subject (i.e., biological replicate) is repetitively measured across multiple time points**. In this example there are 2 time points, Time1 and Time2 (in general the number of times can vary). There are 4 subjects (i.e., biological replicates) measured across times (in general an equal number of times per replicate is not required). There are no technical replicates (in general the number of technical replicates per sample may vary). Overall, in this example there are $2 \times 4 \times 1 = 8$ mass spectrometry runs.

Table below shows the values of the columns **Condition**, **BioReplicate** and **Run** for this situation. Comments on the order of the runs, on the number of lines in the input data structure, and on the handling of missing peak intensities are as in the group comparison design.

Condition	BioReplicate	Run
Time1	Subject1	1
Time2	Subject1	2
Time1	Subject2	3
Time2	Subject2	4
Time1	Subject3	5
Time2	Subject3	6
Time1	Subject4	7
Time2	Subject4	8

3) Paired design

Another frequently used experimental design is a *paired design*, where measurements from multiple conditions (such as healthy biopsy and disease biopsy) are taken from a same subject. The statistical model for this experimental design is the same as in the time course experiment, however the values in the columns of the input data may have a different appearance. In this example there are 2 subjects, PatientA and PatientB (in general the number of patients can vary). There are two conditions per subject, BiopsyHealthy and BiopsyTumor (in general the number of conditions per subject can exceed two). In this example there are 3 technical replicates of each type (in this example, the technical replicates are biopsies; in general these can also be replicate sample preparations or replicate mass spectrometry runs). Overall, in this example there are $2 \times 2 \times 3 = 12$ mass spectrometry runs.

Table below shows the values of the columns **Condition**, **BioReplicate** and **Run** for this situation. Comments on the order of the runs, on the number of lines in the input data structure, and on the handling of missing peak intensities are as in the group comparison design.

Condition	BioReplicate	Run
BiopsyHealthy	PatientA	1

Condition	BioReplicate	Run
BiopsyHealthy	PatientA	2
BiopsyHealthy	PatientA	3
BiopsyTumor	PatientA	4
BiopsyTumor	PatientA	5
BiopsyTumor	PatientA	6
BiopsyHealthy	PatientB	7
BiopsyHealthy	PatientB	8
BiopsyHealthy	PatientB	9
BiopsyTumor	PatientB	10
BiopsyTumor	PatientB	11
BiopsyTumor	PatientB	12

2.1.3 MSnSet format

MSstats also allows data to be in the format of MSnSet, which is consistent with the requirements of Bioconductor. The MSnSet format has several components, of which the most commonly accessed are `assayData`, `phenoData`, and `featureData`. `assayData` is a matrix of intensities, where each row corresponds a transition, and the columns correspond to sample ids. `phenoData` contains columns that describe the biological samples, conditions in the experiment. `featureData` contains columns describing the peptide features, such as the name or id of the underlying protein and information of features.

If the data are stored in the format `expressionSet`, information for group labels is required. If more than one variable is listed in the argument `group`, then a concatenated variable is created based on all of the specified `group` variables. The remaining information (peptide feature ids, biological replicate ids, and abundance) can be extracted from the rows and columns of `featureData` and `phenoData`, or assigned by the users based on the experimental design.

2.2 Label-free DDA

For label-free DDA experiments the required input is the 10-column format, the same as described in section 2.1 for SRM experiments. In DDA experiments spectral features are defined as peptide ions, which are identified and quantified across runs. Since for label-free DDA experiments some of the columns `PeptideSequence`, `PrecursorCharge`, `FragmentIon`, and `ProductCharge` are not relevant, these columns will have a constant fixed value (such as NA) across the entire dataset. Furthermore, the column `IsotopeLabelType` will be set to “L” for the entire dataset.

ProteinName	PeptideSequence	PrecursorCharge	Fragmention	ProductCharge	IsotopeLabelType	Condition	BioReplicate	Run	Intensity
bovine	S.PVDIDTK	5	NA	NA	L	C1	1	1	2636791.5
bovine	S.PVDIDTK	5	NA	NA	L	C1	1	2	1992418.5
bovine	S.PVDIDTK	5	NA	NA	L	C1	1	3	1982146.38
bovine	S.PVDIDTK	5	NA	NA	L	C2	1	4	5019594
bovine	S.PVDIDTK	5	NA	NA	L	C2	1	5	4560467.5
bovine	S.PVDIDTK	5	NA	NA	L	C2	1	6	3627848.75
bovine	S.PVDIDTK	5	NA	NA	L	C5	1	13	145511.83
bovine	S.PVDIDTK	5	NA	NA	L	C5	1	14	291829.69
bovine	S.PVDIDTK	5	NA	NA	L	C6	1	16	786667.38
bovine	S.PVDIDTK	5	NA	NA	L	C6	1	17	705295.31
bovine	S.PVDIDTK	5	NA	NA	L	C6	1	18	453448.78
bovine	S.PVDIDTK	5	NA	NA	L	C3	1	7	NA

2.2 Label-free DIA

For label-free DIA experiments, the required input is the 10-column format, the same as described in section 2.1 for SRM experiments. The values of the required columns can be extracted from the output of signal processing software such as Skyline or OpenSWATH. By default, the combination of the values in the columns `PeptideSequence`, `PrecursorCharge`, `FragmentIon`, `ProductCharge` uniquely identifies each spectral feature (i.e., a fragment ion identified and quantified across multiple runs). If the signal processing software does not provide the information on some of these columns but provides a unique feature identifier, it is possible to use this unique identifier instead of one of these columns. Furthermore, the column `IsotopeLabelType` is set to “L” for the entire dataset.

An example dataset is shown below. In this example, feature id generated by OpenSWATH is used instead of `ProductCharge` to uniquely characterize each feature.

ProteinName	PeptideSequence	PrecursorCharge	FragmentIon	ProductCharge	IsotopeLabelType	Condition	BioReplicate	Run	Intensity
350748	TPPAAVLLK	2	y7	109401	L	2	1	3	257486
350748	TPPAAVLLK	2	y7	109401	L	2	2	4	141159
350748	TPPAAVLLK	2	y7	109401	L	1	1	1	452908
350748	TPPAAVLLK	2	y7	109401	L	1	2	2	348222
515084	NIC[160]VNAIAPGFIESDMTGVLPEK	3	y3	7717	L	2	1	3	12753
515084	NIC[160]VNAIAPGFIESDMTGVLPEK	3	y3	7717	L	2	2	4	12857
515084	NIC[160]VNAIAPGFIESDMTGVLPEK	3	y3	7717	L	1	1	1	89652
515084	NIC[160]VNAIAPGFIESDMTGVLPEK	3	y3	7717	L	1	2	2	76724
515084	MVNEAIESLGSIDVLVNNAGITNDK	3	y9	57971	L	2	1	3	2052
515084	MVNEAIESLGSIDVLVNNAGITNDK	3	y9	57971	L	2	2	4	1050
515084	MVNEAIESLGSIDVLVNNAGITNDK	3	y9	57971	L	1	1	1	10772
515084	MVNEAIESLGSIDVLVNNAGITNDK	3	y9	57971	L	1	2	2	10516

3. Prerequisites and setting for MSstats analysis

MSstats is an R-based package. It is assumed that you already have R installed. You can install MSstats from Bioconductor:

```
source("https://bioconductor.org/biocLite.R")
biocLite("MSstats")
```

Alternatively, you can download the package from the MSstats installation page and install it as follows:

```
install.packages(pkgs = 'MSstats_3.7.3.tar.gz', repos = NULL, type = 'source')
```

Once you have the package installed, load MSstats into an R session and verify that you have the correct version (3.6). Note that in order to use MSstats, the package needs to be loaded every time you restart R.

```
library('MSstats', warn.conflicts = F, quietly = T, verbose = F)
?MSstats
```

Finally, set the working directory to where you saved files. Note that you may have a different path on your computer from the example.

```
setwd('/Users/meenachoi/Dropbox/MSstats_GitHub_document/MSstats_v3.7.3')
```

You can check your working directory by:

```
getwd()
```

```
## [1] "/Users/meenachoi/Dropbox/MSstats_GitHub_document/MSstats_v3.7.3"
```

4. DDA analysis with MSstats

4.1 General workflow for DDA

This section describes a typical workflow for DDA analysis with `MSstats`. Controlled mixture DDA data will be used for demonstration. This dataset is available as an example data(`DDARawData`) in `MSstats`. Also the csv file for the same dataset, `RawData.DDA.csv`, is available in `MSstats` material GitHub in the folder named ‘example dataset/DDA_controlledMixture2009’. It is processed by Superhirn. (original reference link)

4.1.1 Preparing the data for MSstats input

The first step in using the `MSstats` is to format the data as described in Section 2. `DDARawData` is already formatted for `MSstats` input.

```
# Check the first 6 rows in DDARawData
head(DDARawData)
```

```
## ProteinName PeptideSequence PrecursorCharge FragmentIon ProductCharge
## 1 bovine S.PVDIDTK_5 5 NA NA
## 2 bovine S.PVDIDTK_5 5 NA NA
## 3 bovine S.PVDIDTK_5 5 NA NA
## 4 bovine S.PVDIDTK_5 5 NA NA
## 5 bovine S.PVDIDTK_5 5 NA NA
## 6 bovine S.PVDIDTK_5 5 NA NA
## IsotopeLabelType Condition BioReplicate Run Intensity
## 1 L C1 1 1 2636792
## 2 L C1 1 2 1992418
## 3 L C1 1 3 1982146
## 4 L C2 1 4 5019594
## 5 L C2 1 5 4560468
## 6 L C2 1 6 3627849
```

4.1.2 Processing the data

Normalizing and summarizing data with `dataProcess`

After reading the datasets, `MSstats` performs 1) logarithm transformation of `Intensity` column, 2) normalization, 3) feature selection, (all features vs subset of features), 4) imputation for censored missing value, which are below the cutoff and undetectable, 5) run-level summarization.

To get started with this function, visit the help section of `dataProcess` first:

```
?dataProcess
```

NOTE At the logarithm transformation step, zero value in `Intensity` is problematic. When `Intensity=0`, `Inf` is the output from logarithm transformed intensities. Also, logarithm transformed intensities, when `Intensity < 1`, are negative values and it can make overestimated between log fold change. Therefore, logarithm transformed intensities for original intensity between 0 and 1 will be replaced with zero value after normalization.

Default normalization and summarization options

`dataProcess` provides a variety of options in consideration of different experimental protocols. Default values for all options are our suggestion for general cases. However, the default options may not be appropriate for

all possible scenarios. It is important to understand their underlying assumption to avoid misuse. Below is the additional explanation for main options.

- **logTrans** : logarithm transformation with base 2(default) of **Intensity** column.
- **Normalization** :
 - ‘**equalizeMedians**’ : The default option for normalization is **equalizeMedians**, where all the intensities in a run are shifted by a constant, to equalize the median of intensities across runs for label-free experiment. This normalization method is appropriate when we can assume that the majority of proteins do not change across runs. *Be cautious when using the **equalizeMedians** option for a label-free DDA dataset with only a small number of proteins.* For label based experiment, **equalizeMedians** equalizes the median of reference intensities across runs and is generally proper even for a dataset with a small number of proteins.
 - ‘**globalStandards**’ : Instead, if you have a spiked in standard, you may set this to **globalStandards** and define the standard with **nameStandards** option.
 - ‘**quantile**’ : The distribution of all the intensities in each run will become the same across runs for label-free experiment. For label-based experiment, the distribution of all the reference intensities will be become the same across runs and all the endogenous intensities are shifted by a constant corresponding to reference intensities.
 - **FALSE** : No normalization is performed. If you had your own normalization before MSstats, you should use **Normalization=FALSE**.
 - NOTE : If there are multiple fractionations or injections for one sample, normalization is perform by each fractionation or different m/z range from multiple injections.
- **nameStandards** : Only for **Normalization='globalStandards'**, global standard peptide or Protein names, which you can assume that they have the same abundance across MS runs, should be assigned in the vector for this option.
- **featureSubset** :
 - ‘**all**’ : Use all features in the dataset.
 - ‘**top3**’ : Use top 3 features which have highest average of $\log_2(\text{intensity})$ across runs.
 - ‘**topN**’ : Use top N features which have highest average of $\log_2(\text{intensity})$ across runs. It needs the input for **n_top_featureoption** (ex. **n_top_feature=5** for top 5 features).
- **summaryMethod** : Method for run-level summarization.
 - ‘**TMP**’ : Default. Tukey’s median polish (**medpolish** function in stats). Robust parameter estimation method with median across rows and columns.
 - ‘**linear**’ : Linear model (**lm** function). Average-based summarization.
- **MBimpute** : whether model-based imputation will be performed or not. Only for **summaryMethod='TMP'**.
 - **TRUE** : Default. Censored missing values will be imputed by Accelerated Failure Time model. Censored missing values will be determined by other options, **censoredInt** and **maxQuantileforCensored**
 - **FALSE** : No model-based imputation.
- **maxQuantileforCensored** : Maximum quantile for deciding censored missing value. Default is 0.999. If you don’t want to apply the threshold of noise intensity in your data, you can use **maxQuantileforCensored=NULL**.
- **censoredInt** : The processing tools report missing values differently. This option is for distinguish which value should be considered as missing, and further whether it is censored or at random.
 - ‘**NA**’ : Default. It assumes that all NAs in **Intensity** column are censored.
 - ‘**0**’ : It assumes that all values between 0 and 1 in **Intensity** column are censored. If there are NAs in **Intensity** with this option, NAs will be considered as random missing.
 - **NULL** : It assumes that all missing values are randomly missing.

- **cutoffCensored** : cutoff value for AFT model. It is only with `censoredInt='NA'` or `censoredInt='0'`. If you have `censoredInt=NULL`, it assumes that there is no censored missing and any imputation will not be performed.
 - **'minFeature'** : cutoff for AFT model will be the minimum value for each feature across runs.
 - **'minRun'** : cutoff for AFT model will be the minimum value for each run across features.
 - **'minFeatureNRun'** : cutoff for AFT model will be the smallest value between minimum value of corresponding feature and minimum value of corresponding run.

A typical label-free DDA dataset may have many missing values and noisy features with outliers. `MSstats` supports several ways to deal with this. The default option for summarization is `TMP` (robust parameter estimation method with median across rows and columns) after imputation by AFT (accelerated failure time model, `MBimpute=TRUE`) based on censored intensity for NA (`censoredInt="NA"`) with a cutoff as the minimum value for a feature (`cutoffCensored="minFeature"`).

This process handles missing values through imputation and reduces the influence of the outliers using the `TMP` estimation. Note, however, that those runs with no measurements at all will be removed and not be used for any calculation.

```
# default option
DDA2009.proposed <- dataProcess(raw = DDARawData,
                               normalization = 'equalizeMedians',
                               summaryMethod = 'TMP',
                               censoredInt = "NA", cutoffCensored = "minFeature",
                               MBimpute = TRUE,
                               maxQuantileforCensored=0.999)
```

```
## Log2 intensities under cutoff = 13.456 were considered as censored missing values.
```

```
## * Use all features that the dataset originally has.
```

```
##
```

```
## Summary of Features :
```

```
##          count
## # of Protein      6
## # of Peptides/Protein 11-32
## # of Transitions/Peptide 1-1
```

```
## Summary of Samples :
```

```
##          C1 C2 C3 C4 C5 C6
## # of MS runs      3 3 3 3 3 3
## # of Biological Replicates 1 1 1 1 1 1
## # of Technical Replicates 3 3 3 3 3 3
```

```
##
```

```
## Summary of Missingness :
```

```
## # transitions are completely missing in one condition: 90
```

```
## -> D.GPLTGTYR_23_23_NA_NA, F.HFHGSSDDQGSEHTVDR_402_402_NA_NA, G.PLTGTYR_8_8_NA_NA, H.SFNVEYDSSQ
```

```
##
```

```
## # run with 75% missing observations: 0
```

```
##
```

```
## == Start the summarization per subplot...
```

```
## Getting the summarization by Tukey's median polish per subplot for protein bovine ( 1 of 6 )
```

```
## Getting the summarization by Tukey's median polish per subplot for protein chicken ( 2 of 6 )
```

```
## Getting the summarization by Tukey's median polish per subplot for protein cyc_horse ( 3 of 6 )
## Getting the summarization by Tukey's median polish per subplot for protein myg_horse ( 4 of 6 )
## Getting the summarization by Tukey's median polish per subplot for protein rabbit ( 5 of 6 )
## Getting the summarization by Tukey's median polish per subplot for protein yeast ( 6 of 6 )
##
## == the summarization per subplot is done.
```

Output of dataProcess

Output of the `dataProcess` function contains the processed and run-level summarized data as well as relevant information for the summarization step.

```
# output of dataProcess includes several data types.
```

```
names(DDA2009.proposed)
```

```
## [1] "ProcessedData"      "RunlevelData"      "SummaryMethod"
```

```
## [4] "ModelQC"           "PredictBySurvival"
```

```
# the data after reformatting and normalization
```

```
head(DDA2009.proposed$ProcessedData)
```

```
##      PROTEIN                PEPTIDE TRANSITION
## 55   bovine                D.GPLTGTyr_23_23      NA_NA
## 937  bovine F.HFWGSSDDQGSEHTVDR_402_402      NA_NA
## 1628 bovine F.HWGSSDDQGSEHTVDR_229_229      NA_NA
## 19   bovine                G.PLTGTyr_8_8       NA_NA
## 1081 bovine                H.SFNVEYDDSQDK_465_465      NA_NA
## 469  bovine                K.AVVQDPALKPL_156_156    NA_NA
##
##                                FEATURE LABEL GROUP_ORIGINAL
## 55                                D.GPLTGTyr_23_23_NA_NA      L      C1
## 937 F.HFWGSSDDQGSEHTVDR_402_402_NA_NA      L      C1
## 1628 F.HWGSSDDQGSEHTVDR_229_229_NA_NA      L      C1
## 19                                G.PLTGTyr_8_8_NA_NA      L      C1
## 1081 H.SFNVEYDDSQDK_465_465_NA_NA      L      C1
## 469 K.AVVQDPALKPL_156_156_NA_NA      L      C1
##
## SUBJECT_ORIGINAL RUN GROUP SUBJECT SUBJECT_NESTED INTENSITY ABUNDANCE
## 55                1 1 1 1 1 1.1 757400.1 19.83052
## 937                1 1 1 1 1 1.1 2087125.8 21.29291
## 1628               1 1 1 1 1 1.1 1485145.8 20.80200
## 19                 1 1 1 1 1 1.1 4986404.0 22.54939
## 1081               1 1 1 1 1 1.1 2488141.2 21.54646
## 469                1 1 1 1 1 1.1 7519322.0 23.14200
##
## METHOD originalRUN censored
## 55      1          1 FALSE
## 937     1          1 FALSE
## 1628    1          1 FALSE
## 19      1          1 FALSE
## 1081    1          1 FALSE
## 469     1          1 FALSE
```

`DDA2009.TMP$ProcessedData` has the data after normalization and deciding the data-specific threshold for censored missing value. There are several new columns in the datasets. Also dataset is reformatted. **Intensity** column includes original intensities values in the input of `dataProcess`. **ABUNDANCE** column contains the *log2 transformed and normalized intensities and it will used for run-level summarization*. **censored** column has

the decision about censored missing or not, based on `censoredInt` and `maxQuantileforCensored` options. `ABUNDANCE` with `TRUE` value in `censored` column will be considered as censored missing and imputed with `MBimpute=TRUE` option. Censored missing will be distinguished in Profile plot from `dataProcessPlots`.

```
# run-level summarized data
head(DDA2009.proposed$RunlevelData)

##   RUN Protein LogIntensities NumMeasuredFeature MissingPercentage
## 1  1 bovine    21.28437          14          0.00000000
## 2  2 bovine    20.85653          14          0.00000000
## 3  3 bovine    20.67521          13          0.07142857
## 4  4 bovine    21.60443          13          0.07142857
## 5  5 bovine    21.82186          14          0.00000000
## 6  6 bovine    21.20445          13          0.07142857
##   more50missing NumImputedFeature originalRUN GROUP GROUP_ORIGINAL
## 1          FALSE              0           1     1           C1
## 2          FALSE              0           2     1           C1
## 3          FALSE              1           3     1           C1
## 4          FALSE              1           4     2           C2
## 5          FALSE              0           5     2           C2
## 6          FALSE              1           6     2           C2
##   SUBJECT_ORIGINAL SUBJECT_NESTED SUBJECT
## 1                1             1.1     1
## 2                1             1.1     1
## 3                1             1.1     1
## 4                1             2.1     1
## 5                1             2.1     1
## 6                1             2.1     1
```

`DDA2009.TMP$RunlevelData` includes run-level summarized data based on `DDA2009.TMP$ProcessedData`. `LogIntensities` is run-level summarized data and will be used for `groupComparison` function in next step. It will also be used for summarized profile plot (`summaryPlot=TRUE` for `dataProcessPlots` function with `type='ProfilePlot'`). `NumMeasuredFeature` shows how many features were used for summarization of the corresponding run and protein. `MissingPercentage` means the percentage of random and censored missing in the corresponding run and protein out of the total number of feature in the corresponding protein. `more50missing` means whether `MissingPercentage` is greater than 50% or not. `NumImputedFeature` shows how many features were imputed in the corresponding run and protein.

```
# here 'TMP' : It shows which summaryMethod is used for run-level summarization.
head(DDA2009.proposed$SummaryMethod)
```

```
## [1] "TMP"
```

4.1.3 Visualization of processed data

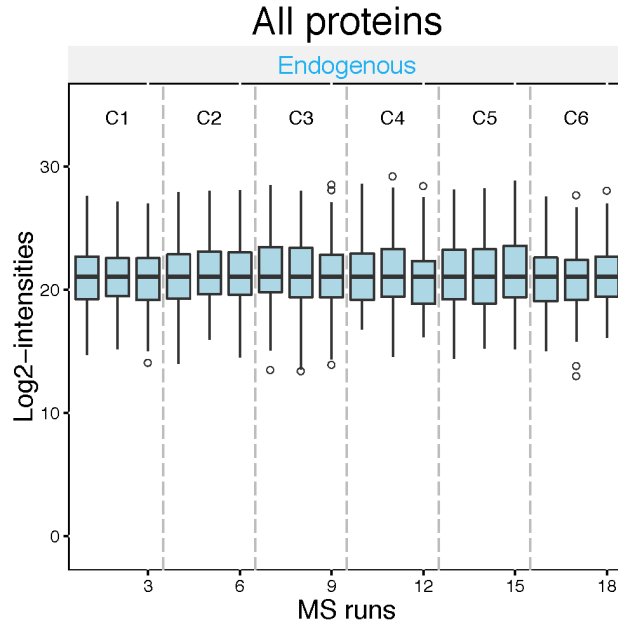
Quality control and normalization effects

QC plot visualizes potential systematic biases between mass spectrometry runs. Also it can be used to assess the effects of the normalization step. After constant normalization, the median intensities of reference transitions across all proteins should be equal between runs. After quantile normalization, the distribution of reference intensities across all proteins should be equal between runs. This step generates two types of QC plots: one for all the proteins combined, and the other separately for each protein (produced in a separate pdf file). These plots can be generated for either all proteins at once or each protein individually if we have a large dataset. The example below shows both options.

```

# use type="QCplot" with all proteins
# change the upper limit of y-axis=35
# set up the size of pdf
dataProcessPlots(data = DDA2009.proposed, type="QCplot", ylimUp=35,
                  width=5, height=5)

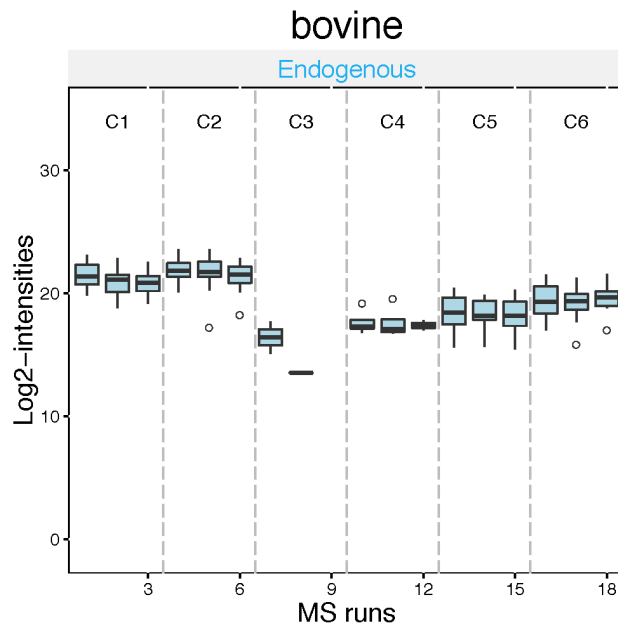
```



```

# use type="QCplot" for 1st protein only
# change the upper limit of y-axis=35
# set up the size of pdf
dataProcessPlots(data = DDA2009.proposed, type="QCplot", which.Protein=1,
                  ylimUp=35, width=5, height=5)

```



NOTE Don't worry about warning messages as below. It means NA values are not included in the plot,

which is a proper way for this case.

Warning messages:

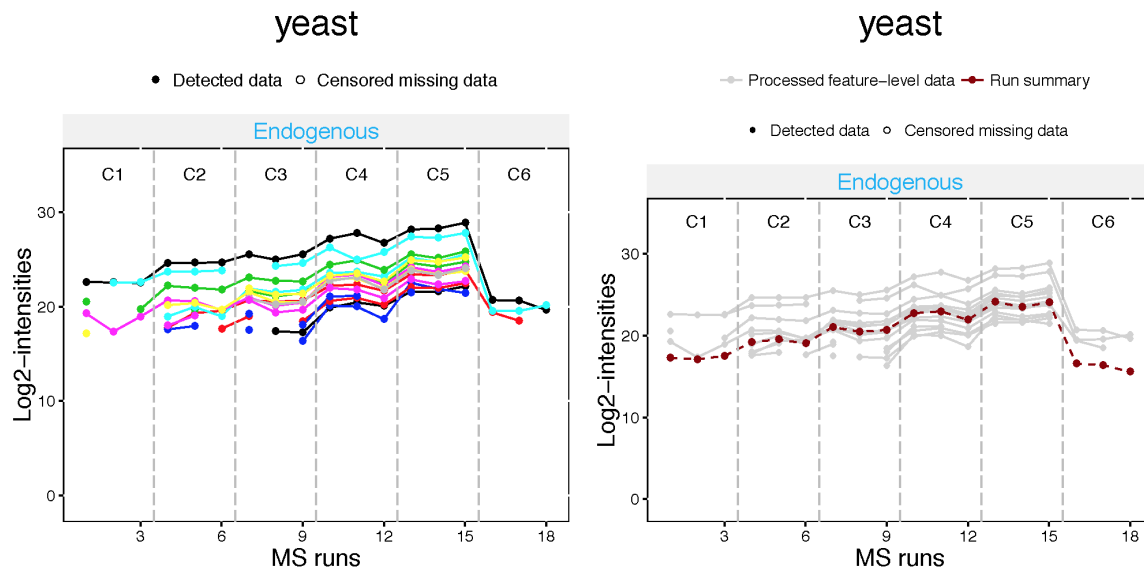
```
1: Removed 698 rows containing non-finite values
(stat_boxplot).
```

...

Profile plot

Profile plot helps identify potential sources of variation (both variation of interest and nuisance variation) for each protein. Such plots should be done after the normalization. Profile plots with summarization present the effects of the summarization step by showing all individual measurements of a protein and their summarized intensity. With `type="profileplot"`, two pdfs will be generated. The first pdf includes plots (per protein) to show individual measurement for each peptide (peptide for DDA, transition for SRM or DIA) across runs, grouped per condition. Each peptide has a different color/type layout. Disconnected lines show that there are missing value (NA). To ignore these plots, please use the option `originalPlot=FALSE`. The second pdf, which is named with 'wSummarization' suffix, shows run-level summarized data per protein. The same peptides (or transition) in the first plot are presented in grey, with the summarized values (by TMP, in this example) overlaid in red. To ignore these plots with summarization, please use the option `summaryPlot=FALSE`.

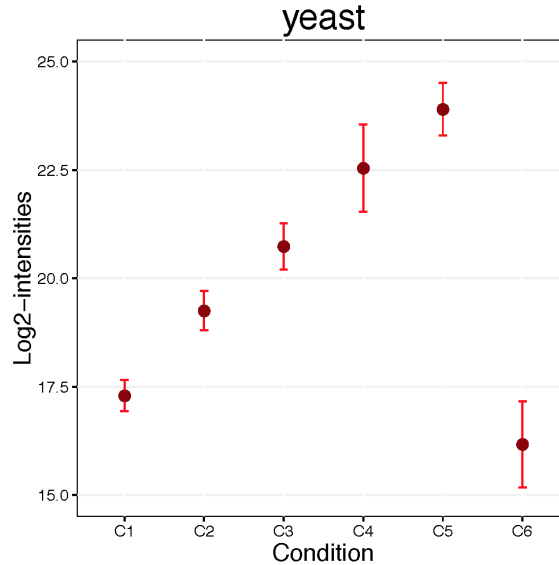
```
dataProcessPlots(data = DDA2009.proposed, type="Profileplot", ylimUp=35,
                 featureName="NA", width=5, height=5, address="DDA2009_proposed_")
```



Condition plot

Condition plot visualizes potential systematic differences in protein intensities between conditions. Dots indicate the mean of log2 intensities for each condition. With the option `interval='CI'` (default), error bars indicate the confidence interval with 0.95 significant level for each condition. With the option `interval='SD'`, error bars indicate the standard deviation among all feature intensities for each condition. **The intervals are for descriptive purposes only, as more refined model-based estimation is obtained as discussed below.** With the option `scale=TRUE`, the levels of conditions are scaled according to their labels. If `scale=FALSE` (default), the conditions on the x-axis are equally spaced.

```
dataProcessPlots(data = DDA2009.proposed, type="Conditionplot",
                 width=5, height=5, address="DDA2009_proposed_")
```

`dataProcessPlots` has a number of layout options, including size and description of axes labels, output file name etc for three types of plots above. The option `address` specifies the name of the folder storing pdf files with the plots. With the option `address=FALSE`, plots will be shown in the graphical window, but not saved in a file. If a file with this name already exists in working directory, a suffix with a number will be appended to the file name. In this way a record of all the analyses is kept.

For more details, visit the help file using the following code.

```
?dataProcessPlots
```

4.1.4 Different imputation options

Here is the summary of combinations for imputation options with `summaryMethod='TMP'`.

- `censoredInt=NULL` : It assumes that all intensities are missing at random and there is no action for missing values with `MBimpute=FALSE`. If you have `MBimpute=TRUE` with `censoredInt=NULL`, there will be error message to fix either `MBimpute` or `censoredInt` options.
- `censoredInt='NA' or '0' & MBimpute=TRUE` : AFT model-based imputation using `cutoffCensored` value in the AFT model.
- `censoredInt='NA' or '0' & MBimpute=FALSE` : censored intensities (here NA's) will be replaced with the value specified in `cutoffCensored`.

NOTE1 The default option for `cutoffCensored` is `minFeature`, taking the minimum value for the corresponding feature. With this option, those runs with substantial missing measurements will be biased by the cutoff value. In such case, you may remove the runs that have more than 50% missing values from the analysis with the option `remove50missing=TRUE`.

NOTE2 In case that there are completely missing measurements in a run for a protein, any imputation will not be performed. In addition, the condition, which has no measurement at all in a protein, will be not imputed.

Here is the example of `dataProcess` option without imputation, assuming that all missing values are random.

```
# No imputation
DDA2009.TMP <- dataProcess(raw = DDARawData,
                          normalization = 'equalizeMedians',
```

```
summaryMethod = 'TMP',
censoredInt = NULL, MBimpute=FALSE)
```

```
## * Use all features that the dataset originally has.
```

```
##
```

```
## Summary of Features :
```

```
##          count
## # of Protein      6
## # of Peptides/Protein  11-32
## # of Transitions/Peptide  1-1
```

```
##
```

```
## Summary of Samples :
```

```
##          C1 C2 C3 C4 C5 C6
## # of MS runs      3 3 3 3 3 3
## # of Biological Replicates  1 1 1 1 1 1
## # of Technical Replicates  3 3 3 3 3 3
```

```
##
```

```
## Summary of Missingness :
```

```
## # transitions are completely missing in one condition: 90
```

```
##   -> D.GPLTGTyr_23_23_NA_NA, F.HFHGSSDDQGSEHTVDR_402_402_NA_NA, G.PLTGTyr_8_8_NA_NA, H.SFNVEYDDSQ
```

```
##
```

```
## # run with 75% missing observations: 0
```

```
##
```

```
## == Start the summarization per subplot...
```

```
## Getting the summarization by Tukey's median polish per subplot for protein bovine ( 1 of 6 )
```

```
## Getting the summarization by Tukey's median polish per subplot for protein chicken ( 2 of 6 )
```

```
## Getting the summarization by Tukey's median polish per subplot for protein cyc_horse ( 3 of 6 )
```

```
## Getting the summarization by Tukey's median polish per subplot for protein myg_horse ( 4 of 6 )
```

```
## Getting the summarization by Tukey's median polish per subplot for protein rabbit ( 5 of 6 )
```

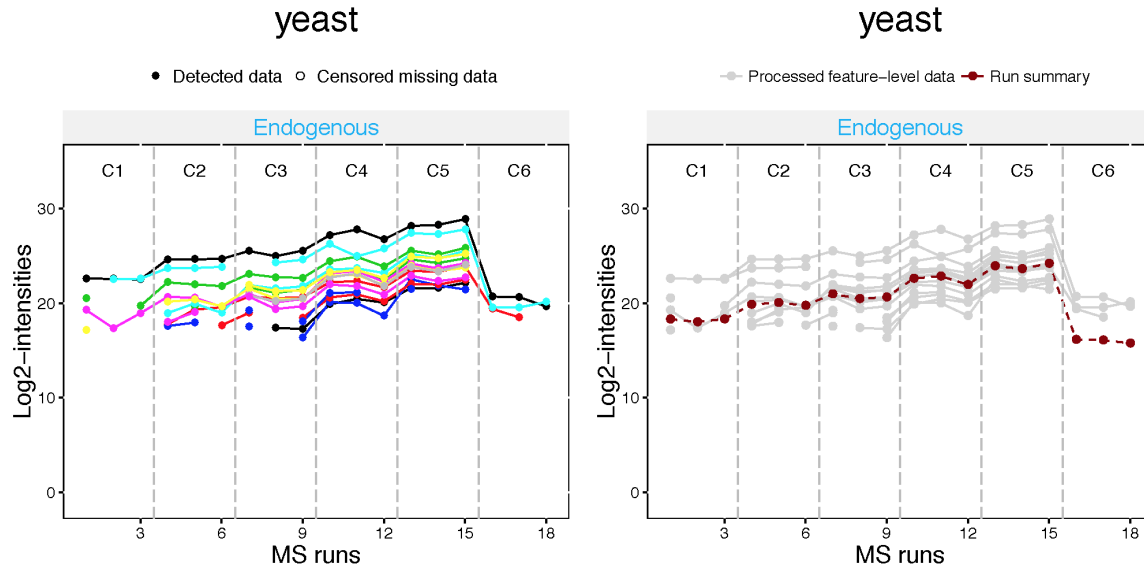
```
## Getting the summarization by Tukey's median polish per subplot for protein yeast ( 6 of 6 )
```

```
##
```

```
## == the summarization per subplot is done.
```

These plots can be used compare and select among different options for imputation (e.g., TMP with or without considering missing values for summarization in `dataProcess`).

```
dataProcessPlots(data = DDA2009.TMP, type="Profileplot", ylimUp=35,
                 featureName="NA", width=5, height=5, address="DDA2009_TMP_")
```



While original profile plots are the same, summarization plots reveal differences, especially for conditions 'C1' and 'C2' in 'yeast' protein, which have many missing values. Without imputation, summarized values in 'C1' group is higher than with imputation for missing values.

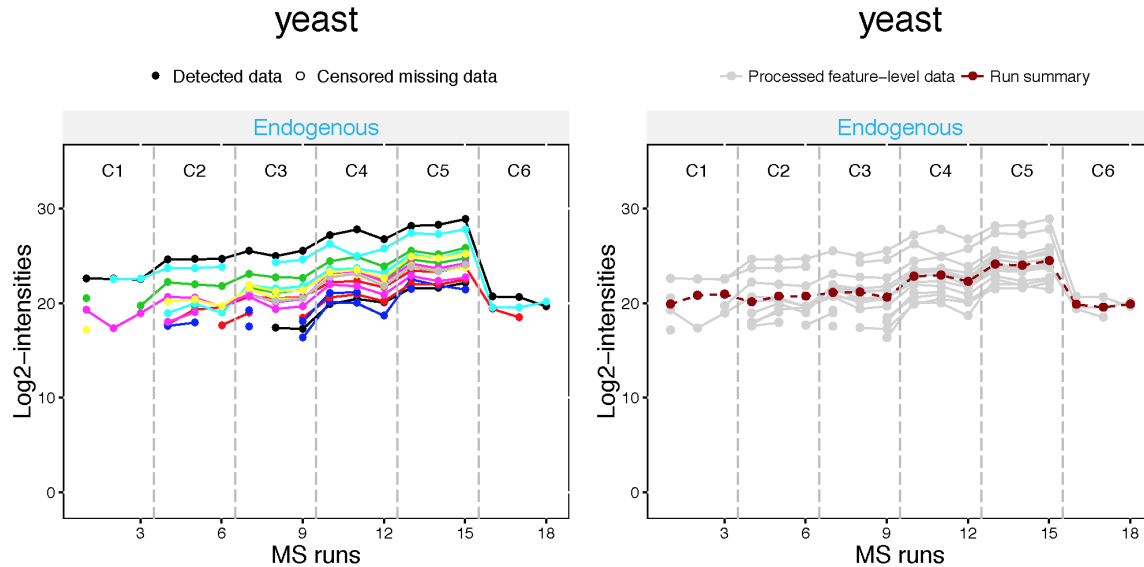
4.1.5 Different summarization options

Besides summarizing observations with the median polish method, MSstats also offers a summarization option using linear model with option `summaryMethod="linear"` with `censoredInt=NULL` assumes that all NA's are missing at random and uses `lm` for parameter estimation.

```
# linear model (lm) with run and feature
DDA2009.linear <- dataProcess(raw = DDARawData,
                             normalization = 'equalizeMedians',
                             summaryMethod = 'linear',
                             censoredInt = NULL,
                             MBimpute = FALSE)
```

Profile plots below can be used compare among different options for summarization (e.g., TMP with or without imputation vs linear for summarization in `dataProcess`).

```
dataProcessPlots(data = DDA2009.linear, type="Profileplot", ylimUp=35,
                 featureName="NA", width=5, height=5, address="DDA2009_linear_")
```



While original profile plots are the same, summarization plots reveal differences, especially for conditions ‘C1’, ‘C2’, and ‘C6’ in ‘yeast’ protein, which have many missing values. Summarized values with linear model in these groups are much higher than those with TMP considering missing values or not.

4.1.6 Finding differentially abundant proteins across conditions

Comparing conditions with groupComparison

With the normalized data and run-level summarized data obtained by applying one of the `dataProcess` summarization methods, it is of general interest to find proteins changing between groups of conditions. Within `MSstats` this can be done by using the `groupComparison` function, which takes the output of the `dataProcess` function as input.

```
?groupComparison
```

In addition to the processed data, the `groupComparison` function requires a contrast matrix to define the comparison to be made. The contrast matrix is created with each condition in column and each comparison in row. Note that the conditions are arranged **in alphabetical order**. The order of condition that `MSstats` recognizes can be shown by using `levels`:

```
levels(DDA2009.TMP$ProcessedData$GROUP_ORIGINAL)
```

```
## [1] "C1" "C2" "C3" "C4" "C5" "C6"
```

Entries in each row of the contrast matrix are filled in with 0, 1, or -1 to specify the comparison, where **0** is for conditions we would like to ignore, **1** is for conditions we would like to put in the numerator of the ratio or fold-change, and **-1** is for conditions we would like to put in the denominator of the ratio or fold-change.

For example, if you want to compare C2-C1, which means $\log(C2)-\log(C1)$ and the same as $\log(C2/C1)$, set ‘1’ for C2 and ‘-1’ for C1 in the row. Combining multiple groups for comparison is also possible. For example, if you want to compare between average of C2 and C3 and average of C1, $(C3+C2)/2-C1$ as formula, set ‘-1’ for C1, ‘0.5’ for C2 and ‘0.5’ for C3, and ‘0’ for rest of groups.

```
comparison1 <- matrix(c(-1,1,0,0,0,0),nrow=1)
comparison2 <- matrix(c(0,-1,1,0,0,0),nrow=1)
comparison3 <- matrix(c(0,0,-1,1,0,0),nrow=1)
comparison4 <- matrix(c(0,0,0,-1,1,0),nrow=1)
comparison5 <- matrix(c(0,0,0,0,-1,1),nrow=1)
```

```
comparison6 <- matrix(c(1,0,0,0,0,-1),nrow=1)
```

```
comparison<-rbind(comparison1,comparison2,comparison3,comparison4,comparison5,comparison6)  
row.names(comparison) <- c("C2-C1","C3-C2","C4-C3","C5-C4","C6-C5","C1-C6")
```

With the contrast matrix specified, group comparison can be performed as follows.

```
DDA2009.comparisons <- groupComparison(contrast.matrix = comparison, data = DDA2009.proposed)
```

Output of the groupComparison function contains three data frames:

```
# output from groupComparison function has three data frames  
names(DDA2009.comparisons)
```

```
## [1] "ComparisonResult" "ModelQC" "fittedmodel"
```

Results of the statistical comparison are stored in the data frame named ComparisonResult:

```
# name of columns in result data.frame  
head(DDA2009.comparisons$ComparisonResult)
```

```
##      Protein Label      log2FC      SE      Tvalue DF      pvalue  
## 1      bovine C2-C1  0.6048799 0.4245943  1.424607 11 1.820186e-01  
## 7      chicken C2-C1  0.7876884 0.2205455  3.571545 12 3.841470e-03  
## 13     cyc_horse C2-C1  1.1294964 0.1955787  5.775149 12 8.809149e-05  
## 19     myg_horse C2-C1 -7.9717333 0.2807086 -28.398612 12 2.254641e-12  
## 25      rabbit C2-C1  1.0617105 0.2209612  4.804963 12 4.298913e-04  
## 31      yeast C2-C1  1.9575344 0.3134408  6.245309 12 4.284464e-05  
##      adj.pvalue issue MissingPercentage ImputationPercentage  
## 1 1.820186e-01 NA 0.03571429 0.03571429  
## 7 4.609764e-03 NA 0.25757576 0.25757576  
## 13 1.761830e-04 NA 0.15104167 0.15104167  
## 19 1.352785e-11 NA 0.45833333 0.45833333  
## 25 6.448369e-04 NA 0.72043011 0.72043011  
## 31 1.285339e-04 NA 0.56666667 0.56666667
```

The result of the test for differential abundance is a table with columns **Protein**, **Label** (of the comparison), **log2 fold change (log2FC)**, standard error of the log2 fold change (**SE**), test statistic of the Student test (**Tvalue**), degree of freedom of the Student test (**DF**), raw p-values (**pvalue**), p-values adjusted among all the proteins in the specific comparison using the approach by Benjamini and Hochberg (**adj.pvalue**). The cutoff of the adjusted p-value corresponds to the cutoff of the False Discovery Rate (Benjamini and Hochberg 1955). The positive values of log2FC for **Label=C2-C1** indicate evidence in favor of $C2 > C1$ (i.e. proteins upregulated in C2), while the negative values indicate evidence in favor of $C2 < C1$ (i.e. proteins downregulated in C2), as compared to C1. The same model can be used to perform several comparisons of conditions simultaneously in the same protein.

NOTE **issue** column shows if there is any issue for inference in corresponding protein and comparison, for example, **OneConditionMissing** or **CompleteMissing**. If one of condition for comparison is completely missing, it would flag with **OneConditionMissing** with **adj.pvalue=0** and **log2FC=Inf** or **-Inf** even though **pvalue=NA**. For example, if you want to compare ‘Condition A - Condition B’, but condition B has complete missing, **log2FC=Inf** and **adj.pvalue=0**. **SE**, **Tvalue**, and **pvalue** will be **NA**. if you want to compare ‘Condition A - Condition B’, but condition A has complete missing, then **log2FC=-Inf** and **adj.pvalue=0**. But, please be careful for using this **log2FC** and **adj.pvalue**.

Based on the comparison results and desired significance level, a short list of the differentially abundant proteins can be obtained for further investigation:

```
# get only significant proteins and comparisons among all comparisons
# To simultaneously controll the overall FDR at the level, 0.05
SignificantProteins <- with(DDA2009.comparisons,
                           ComparisonResult[ComparisonResult$adj.pvalue < 0.05, ])
nrow(SignificantProteins)
```

```
## [1] 34
```

4.1.6 Verifying the assumption of the model

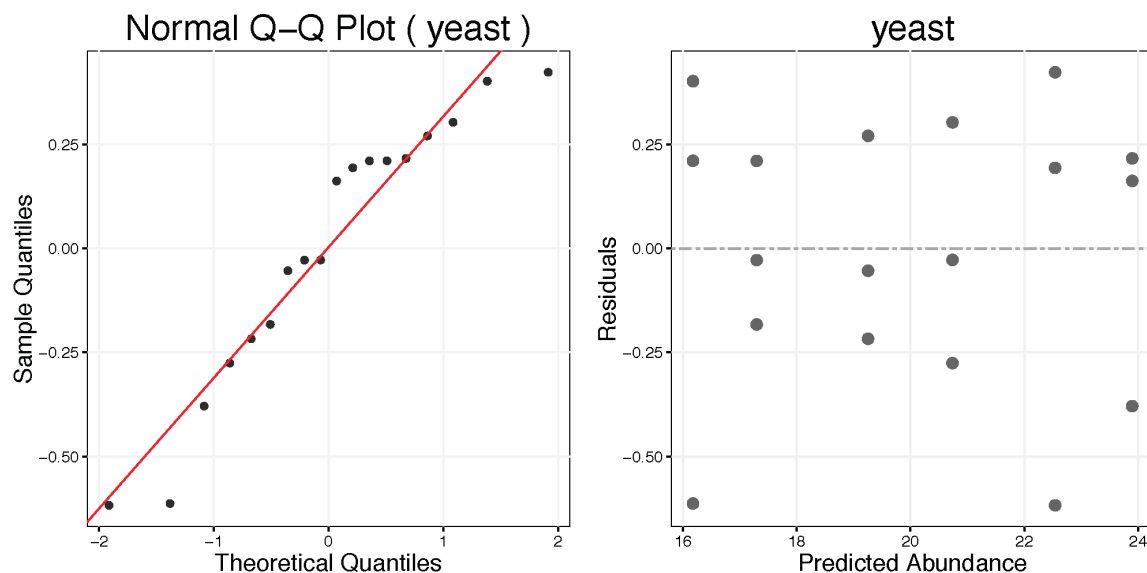
Results based on the statistical models are accurate as long as the assumptions of the models hold. Here we focus on the assumption of the Normal distribution of the measurement errors, and also on the assumption of constant variance of the measurement errors (if this option is specified in the model above). The assumptions can be checked by examining the residuals of the model fit (i.e., the deviations of the observed intensities of the transition from their model-based predictions).

`modelBasedQCPlots` function generates residual plots and Normal quantile-quantile plots for each protein, taking as input the results of model fitting and testing in `groupComparison`. Normal quantile-quantile plot with the option `type='QQPlots'` illustrates that such deviations from constant variance can be mistaken for deviations from Normality. Only large deviations of transition intensities from the straight line are problematic.

Residual plot with the option `type='ResidualPlots'` shows variance of the residuals that is associated with the mean feature intensity. Any specific pattern, such as increasing or decreasing by predicted abundance, is problematic.

```
# normal quantile-quantile plots
modelBasedQCPlots(data=DDA2009.comparisons, type="QQPlots",
                  width=5, height=5, address="DDA2009_proposed_")

# residual plots
modelBasedQCPlots(data=DDA2009.comparisons, type="ResidualPlots",
                  width=5, height=5, address="DDA2009_proposed_")
```



For more details, visit the help file using the following code.

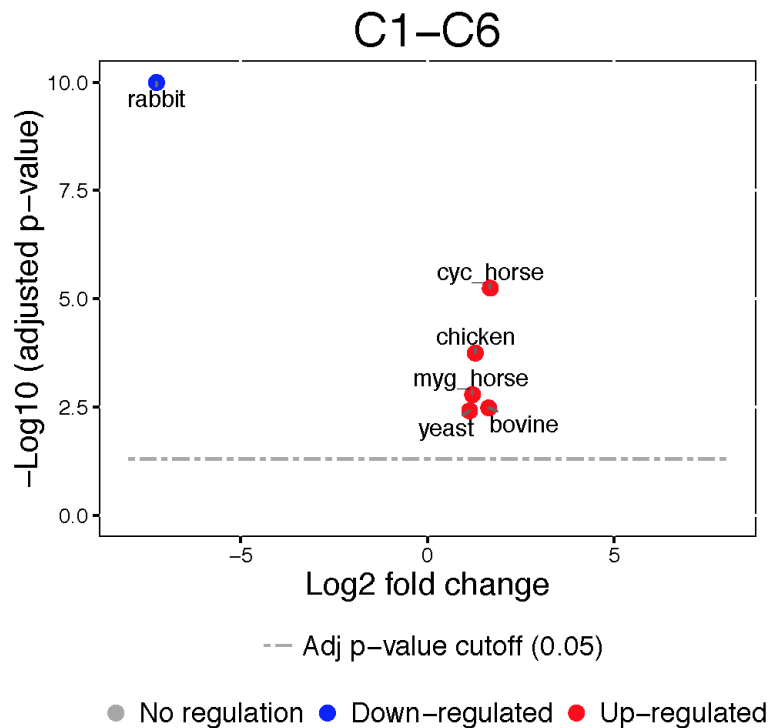
4.1.7 Visualization of differentially abundant proteins

Volcano plots

Volcano plots visualize the outcome of one comparison between conditions for all the proteins, and combine the information on statistical and practical significance. The y-axis displays the FDR-adjusted p-values on the negative log10 scale, representing statistical significance. The horizontal dashed line shows the FDR cutoff. The points above the FDR cutoff line are statistically significant proteins that are differentially abundant across conditions. These points are colored in red and blue for upregulated and downregulated proteins, respectively. The x-axis is the model-based estimate of fold change on log scale (the base of logarithm transform is the same as specified in the `logTrans` option of the `dataProcess` function), and represents practical significance. It is possible to specify a practical significance cutoff based on the estimate of fold change in addition to the statistical significance cutoff. If the fold change cutoff is specified, the points above the horizontal cutoff line but within the vertical cutoff line will be considered as not differentially abundant (and will be colored in black). The practical significance cutoff should only be applied in addition to the statistical significance cutoff (i.e., the fold change alone does not present enough evidence for differential abundance).

```
groupComparisonPlots(data = DDA2009.comparisons$ComparisonResult, type = 'VolcanoPlot',  
                    width=5, height=5, address="DDA2009_proposed_")
```

'VolcanoPlot.pdf' will be saved under the folder you assigned. It has the plots per comparison defined in `contrast.matrix`. Please check `?groupComparisonPlots` for detail, such as labelling protein names, size of dots, font sizes, etc. Below is one of volcano plots, for comparison 'C1-C6' including protein name labelling. Protein name will be shown for significant proteins, without overlapping protein names each other.



Heatmap

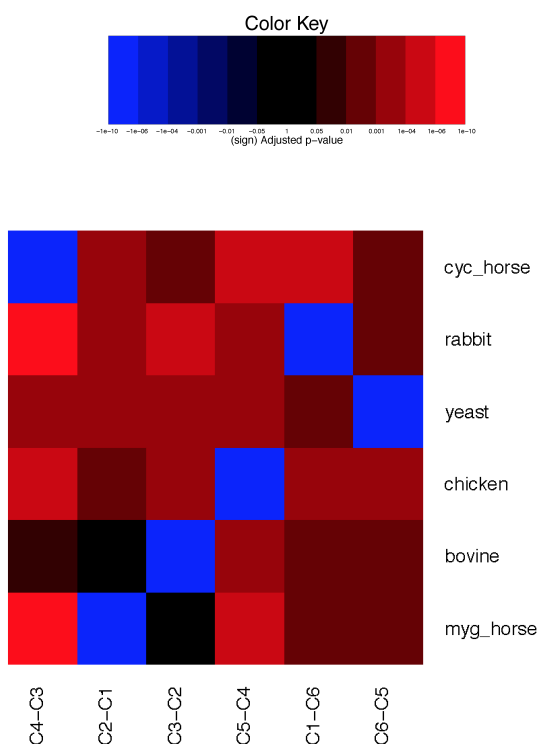
Heatmaps illustrate the patterns of up- and down-regulation of proteins in several comparisons. Columns in the heatmaps are comparison of conditions assigned in `contrast.matrix`, and rows are proteins. The heatmaps display signed FDR-adjusted p-values of the tests, colored in red/blue for significantly up-/down-regulated proteins, while taking into account the specified FDR cutoff and the additional optional fold change cutoff. Brighter colors indicate stronger evidence in favor of differential abundance. Black color represents proteins that are not significantly differentially abundant.

NOTE To draw heatmap, at least two comparisons are needed.

The rows and columns of the heatmaps can be ordered with the option `clustering`, which performs hierarchical clustering with the Ward method (minimum variance). The option `clustering='protein'` (default) clusters the rows (proteins) in the space of comparisons, based on the values of $(\text{sign of comparison}) \cdot (-\log_2(\text{adjusted p-values}))$. The option `clustering='comparison'` clusters the columns in the space of proteins, based on the values of $(\text{sign of comparison}) \cdot (-\log_2(\text{adjusted p-value}))$. The option `clustering='both'` reorders both columns and rows.

```
groupComparisonPlots(data = DDA2009.comparisons$ComparisonResult, type = 'Heatmap')
```

'Heatmap.pdf' will be saved under the folder you assigned. Below is one example, showing the results for several comparisons simultaneously.

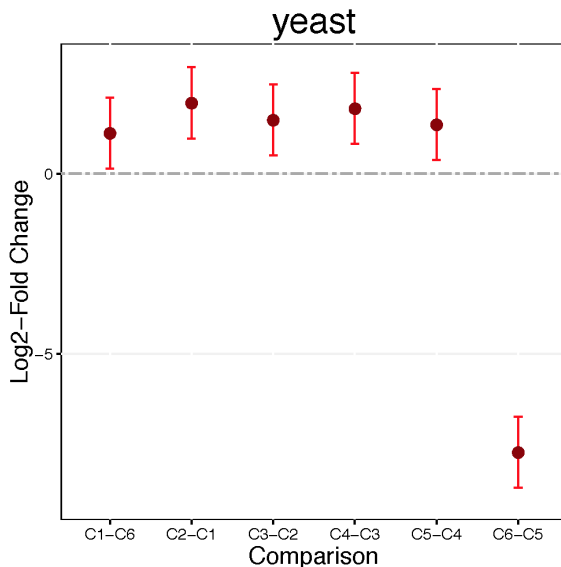


Comparison plots

Comparison plots illustrate model-based estimates of log-fold changes, and the associated uncertainty, in several comparisons of conditions for one protein. X-axis is the comparison of interest. Y-axis is the log fold change. The dots are the model-based estimates of log-fold change, and the error bars are the model-based 95% confidence intervals (the option `sig` can be used to change the significance level of significance). For

simplicity, the confidence intervals are adjusted for multiple comparisons within protein only, using the Bonferroni approach. For proteins with N comparisons, the individual confidence intervals are at the level of $1-\text{sig}/N$.

```
groupComparisonPlots(data=DDA2009.comparisons$ComparisonResult, type="ComparisonPlot",
                    width=5, height=5, address="DDA2009_proposed_")
```



For further details, such as labelling protein names, size of dots, font sizes, etc., visit the help file using the following code.

```
?groupComparisonPlots
```

4.1.8 Sample size calculation for a future experiment

This last analysis step views the dataset as a pilot study of a future experiment, utilizes its variance components, and calculates the minimal number of replicates required in a future experiment to achieve the desired statistical power. The calculation is performed by the function `designSampleSize`, which takes as input the fitted model in `groupComparison`. Sample size calculation assumes same experimental design (i.e. group comparison, time course or paired design) as in the current dataset, and uses the model fit to estimate the median variance components across all the proteins. Finally, sample size calculation assumes that a large proportion of proteins (specifically, 99%) will not change in abundance in the future experiment. This assumption also provides conservative results. Using the estimated variance components, the function relates the number of biological replicates per condition (`numSample`, rounded to 0 decimal), average statistical power across all the proteins (`power`), minimal fold change that we would like to detect (can be specified as a range, e.g. `desiredFC=c(1.1, 2)`), and the False Discovery Rate (FDR). The user should specify all these quantities but one, and the function will solve for the remainder. The quantity to solve for should be set to `= TRUE`.

```
# Minimal number of biological replicates per condition
result.sample <- designSampleSize(data=DDA2009.comparisons$fittedmodel, numSample=TRUE,
                                desiredFC=c(1.25, 3), FDR=0.05, power=0.8)
result.sample
```

```
##      desiredFC numSample  FDR power   CV
## 1         1.250         35 0.05  0.8 0.004
## 2         1.275         30 0.05  0.8 0.005
## 3         1.300         25 0.05  0.8 0.006
```

## 4	1.325	22 0.05	0.8 0.007
## 5	1.350	19 0.05	0.8 0.007
## 6	1.375	17 0.05	0.8 0.008
## 7	1.400	16 0.05	0.8 0.009
## 8	1.425	14 0.05	0.8 0.010
## 9	1.450	13 0.05	0.8 0.010
## 10	1.475	12 0.05	0.8 0.011
## 11	1.500	11 0.05	0.8 0.012
## 12	1.525	10 0.05	0.8 0.013
## 13	1.550	9 0.05	0.8 0.014
## 14	1.575	9 0.05	0.8 0.014
## 15	1.600	8 0.05	0.8 0.015
## 16	1.625	7 0.05	0.8 0.017
## 17	1.650	7 0.05	0.8 0.017
## 18	1.675	7 0.05	0.8 0.016
## 19	1.700	6 0.05	0.8 0.019
## 20	1.725	6 0.05	0.8 0.018
## 21	1.750	6 0.05	0.8 0.018
## 22	1.775	5 0.05	0.8 0.022
## 23	1.800	5 0.05	0.8 0.021
## 24	1.825	5 0.05	0.8 0.021
## 25	1.850	5 0.05	0.8 0.021
## 26	1.875	4 0.05	0.8 0.026
## 27	1.900	4 0.05	0.8 0.025
## 28	1.925	4 0.05	0.8 0.025
## 29	1.950	4 0.05	0.8 0.025
## 30	1.975	4 0.05	0.8 0.024
## 31	2.000	4 0.05	0.8 0.024
## 32	2.025	4 0.05	0.8 0.024
## 33	2.050	3 0.05	0.8 0.031
## 34	2.075	3 0.05	0.8 0.031
## 35	2.100	3 0.05	0.8 0.030
## 36	2.125	3 0.05	0.8 0.030
## 37	2.150	3 0.05	0.8 0.030
## 38	2.175	3 0.05	0.8 0.029
## 39	2.200	3 0.05	0.8 0.029
## 40	2.225	3 0.05	0.8 0.029
## 41	2.250	3 0.05	0.8 0.028
## 42	2.275	3 0.05	0.8 0.028
## 43	2.300	3 0.05	0.8 0.028
## 44	2.325	2 0.05	0.8 0.041
## 45	2.350	2 0.05	0.8 0.041
## 46	2.375	2 0.05	0.8 0.040
## 47	2.400	2 0.05	0.8 0.040
## 48	2.425	2 0.05	0.8 0.039
## 49	2.450	2 0.05	0.8 0.039
## 50	2.475	2 0.05	0.8 0.039
## 51	2.500	2 0.05	0.8 0.038
## 52	2.525	2 0.05	0.8 0.038
## 53	2.550	2 0.05	0.8 0.038
## 54	2.575	2 0.05	0.8 0.037
## 55	2.600	2 0.05	0.8 0.037
## 56	2.625	2 0.05	0.8 0.036
## 57	2.650	2 0.05	0.8 0.036

```

## 58      2.675      2 0.05   0.8 0.036
## 59      2.700      2 0.05   0.8 0.035
## 60      2.725      2 0.05   0.8 0.035
## 61      2.750      2 0.05   0.8 0.035
## 62      2.775      2 0.05   0.8 0.034
## 63      2.800      2 0.05   0.8 0.034
## 64      2.825      2 0.05   0.8 0.034
## 65      2.850      2 0.05   0.8 0.034
## 66      2.875      2 0.05   0.8 0.033
## 67      2.900      2 0.05   0.8 0.033
## 68      2.925      2 0.05   0.8 0.033
## 69      2.950      1 0.05   0.8 0.065
## 70      2.975      1 0.05   0.8 0.064
## 71      3.000      1 0.05   0.8 0.064

```

```
# Power calculation
```

```

result.power <- designSampleSize(data=DDA2009.comparisons$fittedmodel, numSample=3,
                                desiredFC=c(1.25, 3), FDR=0.05, power=TRUE)

```

```
result.power
```

```

##      desiredFC numSample  FDR power    CV
## 1          1.250         3 0.05  0.01 0.051
## 2          1.275         3 0.05  0.01 0.050
## 3          1.300         3 0.05  0.01 0.049
## 4          1.325         3 0.05  0.01 0.048
## 5          1.350         3 0.05  0.01 0.047
## 6          1.375         3 0.05  0.01 0.046
## 7          1.400         3 0.05  0.01 0.046
## 8          1.425         3 0.05  0.01 0.045
## 9          1.450         3 0.05  0.01 0.044
## 10         1.475         3 0.05  0.01 0.043
## 11         1.500         3 0.05  0.02 0.043
## 12         1.525         3 0.05  0.03 0.042
## 13         1.550         3 0.05  0.04 0.041
## 14         1.575         3 0.05  0.06 0.041
## 15         1.600         3 0.05  0.08 0.040
## 16         1.625         3 0.05  0.10 0.039
## 17         1.650         3 0.05  0.13 0.039
## 18         1.675         3 0.05  0.16 0.038
## 19         1.700         3 0.05  0.20 0.038
## 20         1.725         3 0.05  0.23 0.037
## 21         1.750         3 0.05  0.27 0.036
## 22         1.775         3 0.05  0.31 0.036
## 23         1.800         3 0.05  0.35 0.035
## 24         1.825         3 0.05  0.39 0.035
## 25         1.850         3 0.05  0.43 0.034
## 26         1.875         3 0.05  0.47 0.034
## 27         1.900         3 0.05  0.51 0.034
## 28         1.925         3 0.05  0.55 0.033
## 29         1.950         3 0.05  0.58 0.033
## 30         1.975         3 0.05  0.61 0.032
## 31         2.000         3 0.05  0.65 0.032
## 32         2.025         3 0.05  0.68 0.032
## 33         2.050         3 0.05  0.71 0.031
## 34         2.075         3 0.05  0.73 0.031

```

```

## 35      2.100          3 0.05  0.76 0.030
## 36      2.125          3 0.05  0.78 0.030
## 37      2.150          3 0.05  0.80 0.030
## 38      2.175          3 0.05  0.82 0.029
## 39      2.200          3 0.05  0.84 0.029
## 40      2.225          3 0.05  0.86 0.029
## 41      2.250          3 0.05  0.87 0.028
## 42      2.275          3 0.05  0.88 0.028
## 43      2.300          3 0.05  0.90 0.028
## 44      2.325          3 0.05  0.91 0.027
## 45      2.350          3 0.05  0.92 0.027
## 46      2.375          3 0.05  0.93 0.027
## 47      2.400          3 0.05  0.93 0.027
## 48      2.425          3 0.05  0.94 0.026
## 49      2.450          3 0.05  0.95 0.026
## 50      2.475          3 0.05  0.95 0.026
## 51      2.500          3 0.05  0.96 0.026
## 52      2.525          3 0.05  0.96 0.025
## 53      2.550          3 0.05  0.97 0.025
## 54      2.575          3 0.05  0.97 0.025
## 55      2.600          3 0.05  0.98 0.025
## 56      2.625          3 0.05  0.98 0.024
## 57      2.650          3 0.05  0.98 0.024
## 58      2.675          3 0.05  0.98 0.024
## 59      2.700          3 0.05  0.99 0.024
## 60      2.725          3 0.05  0.99 0.023
## 61      2.750          3 0.05  0.99 0.023
## 62      2.775          3 0.05  0.99 0.023
## 63      2.800          3 0.05  0.99 0.023
## 64      2.825          3 0.05  0.99 0.023
## 65      2.850          3 0.05  0.99 0.022
## 66      2.875          3 0.05  0.99 0.022
## 67      2.900          3 0.05  0.99 0.022
## 68      2.925          3 0.05  0.99 0.022
## 69      2.950          3 0.05  0.99 0.022
## 70      2.975          3 0.05  0.99 0.021
## 71      3.000          3 0.05  0.99 0.021

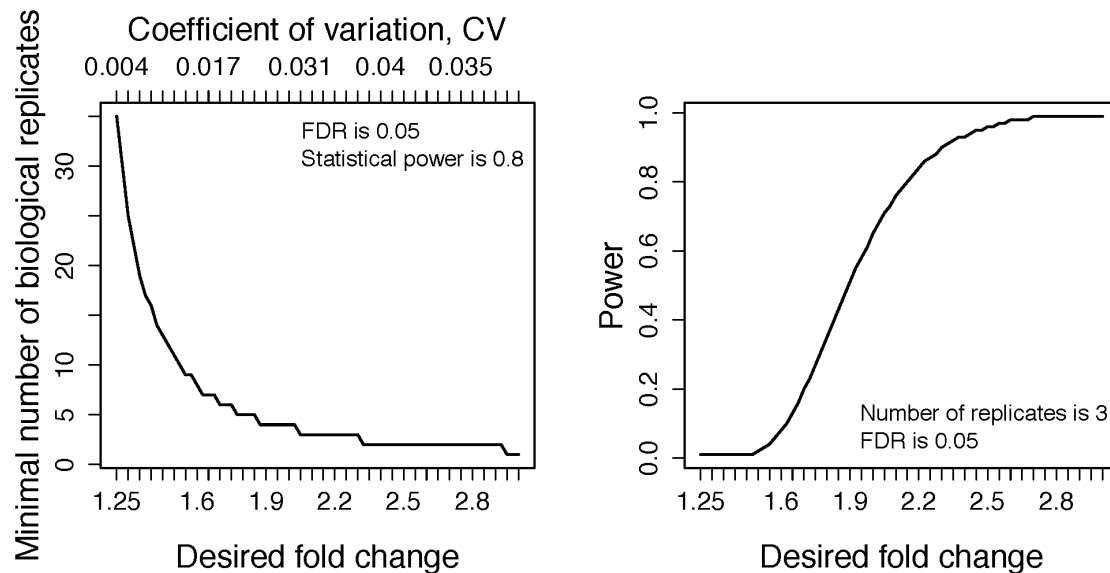
```

For further details, visit the help file using the following code.

```
?designSampleSize
```

Visualization of sample size calculations

The calculated relationship between the number of biological replicates per condition (`numSample`), average statistical power across all the proteins (`power`), minimal fold change that we would like to detect (`desiredFC`), and the False Discovery Rate (FDR) can be visualized using the function `designSampleSizePlots`. The function takes as input the output of `designSampleSize`.



For further details, visit the help file using the following code.

```
?designSampleSizePlots
```

4.1.9 Quantification of protein abundance in individual samples or conditions

Many downstream analysis steps (such as clustering or classification of individual samples in the space of their protein profiles) require summary values of protein abundance in each biological replicate or in each condition, on a relative scale that is comparable between runs.

`dataProcess` function performs model-based run-level summarization. `quantification` function enables subject-level summarization or group-level summarization with the run-level summarization from `dataProcess`.

The option, `type='sample'` (default), performs sample quantification, i.e. it outputs the estimates of relative protein abundance in each biological replicate. If there are technical replicates for biological replicates, sample quantification will be the median among technical replicates. If there is no technical replicate for biological replicate (sample), sample quantification will be the same as run-level summarization. In presence of completely missing values in biological replicate, the estimates will be zero.

The option `type='group'` performs group quantification, i.e. it outputs the estimates of relative protein abundance in each condition, summarized over the biological replicates (median among sample quantification). In presence of completely missing values in a condition, the estimates will be zero.

MSstats supports two output formats. The option `format='matrix'` (default) outputs an array where rows are `proteins`, and columns are `conditions` (for group quantification), or combinations of biological replicate and condition ids (for sample quantification). The option `format='long'` produces an array where each row corresponding to relative protein abundances, and columns are `Protein`, `Condition`, `LogIntensities` (and `BioReplicate` in the case of sample quantification).

```
subQuant <- quantification(DDA2009.proposed)
head(subQuant)
```

```
##   Protein   C1_1   C2_1   C3_1   C4_1   C5_1   C6_1
## 1  bovine 20.85653 21.60443 14.32690 16.10441 17.63141 19.27802
## 2  chicken 18.48792 19.43204 20.41274 22.42284 15.92462 17.09803
## 3  cyc_horse 20.25927 21.33967 22.22028 15.85252 17.62720 18.45536
## 4  myg_horse 22.66495 14.73701 14.99667 18.61740 20.26392 21.52022
```

```
## 5    rabbit 14.89507 15.88492 17.43767 20.19014 21.27964 22.07550
## 6     yeast 17.26792 19.19987 20.71073 22.73666 24.06156 16.38660
```

```
groupQuant <- quantification(DDA2009.proposed, type='group')
head(groupQuant)
```

```
##      Protein      C1      C2      C3      C4      C5      C6
## 1    bovine 20.85653 21.60443 14.32690 16.10441 17.63141 19.27802
## 2   chicken 18.48792 19.43204 20.41274 22.42284 15.92462 17.09803
## 3 cyc_horse 20.25927 21.33967 22.22028 15.85252 17.62720 18.45536
## 4 myg_horse 22.66495 14.73701 14.99667 18.61740 20.26392 21.52022
## 5    rabbit 14.89507 15.88492 17.43767 20.19014 21.27964 22.07550
## 6     yeast 17.26792 19.19987 20.71073 22.73666 24.06156 16.38660
```

For further details, visit the help file using the following code.

```
?quantification
```

4.2 Suggested workflow with Skyline output for DDA

This section describes steps and considerations to properly format data processed by Skyline, prior to the MSstats analysis. In the following example, the raw files for Dynamic benchmark dataset (J. Cox et al. 2014) are used, searched by Andromeda in MaxQuant, but quantified by Skyline.

4.2.1 Load Skyline output

This required input data is generated automatically when using MSstats report format in Skyline. We first load and access the dataset processed by Skyline. The name of saved file from Skyline using MSstats report format is 'Cox.Skyline.csv'.

```
# Read output from skyline : Cox.skyline.csv
raw <- read.csv("Cox.skyline.csv")
```

We can read csv file. Here we will load R data file which is the exactly same data in Cox.skyline.csv file.

```
# Load R data, which is converted from csv file, output from skyline : Cox.skyline.csv
load("Cox.skyline.RData")
raw <- Cox.skyline
head(raw)
```

```
##      ProteinName PeptideSequence PeptideModifiedSequence
## 73 sp|P02768|ALBU_HUMAN      DLGEENFK      DLGEENFK
## 74 sp|P02768|ALBU_HUMAN      DLGEENFK      DLGEENFK
## 75 sp|P02768|ALBU_HUMAN      DLGEENFK      DLGEENFK
## 76 sp|P02768|ALBU_HUMAN      DLGEENFK      DLGEENFK
## 77 sp|P02768|ALBU_HUMAN      DLGEENFK      DLGEENFK
## 78 sp|P02768|ALBU_HUMAN      DLGEENFK      DLGEENFK
##      PrecursorCharge PrecursorMz FragmentIon ProductCharge ProductMz
## 73                2    476.2245 precursor          2    476.2245
## 74                2    476.2245 precursor          2    476.2245
## 75                2    476.2245 precursor          2    476.2245
## 76                2    476.2245 precursor          2    476.2245
## 77                2    476.2245 precursor          2    476.2245
## 78                2    476.2245 precursor          2    476.2245
##      IsotopeLabelType Condition BioReplicate
## 73                light      NA      NA
```

```

## 74          light          NA          NA
## 75          light          NA          NA
## 76          light          NA          NA
## 77          light          NA          NA
## 78          light          NA          NA
##              FileName          Area StandardType Truncated
## 73 20130510_EXQ1_IgPa_QC_UPS1_01.raw 1222621440          NA      False
## 74 20130510_EXQ1_IgPa_QC_UPS1_02.raw 1641793792          NA      False
## 75 20130510_EXQ1_IgPa_QC_UPS1_03.raw 1225490048          NA      False
## 76 20130510_EXQ1_IgPa_QC_UPS1_04.raw 1777631616          NA      False
## 77 20130510_EXQ1_IgPa_QC_UPS2_01.raw 7395562496          NA      False
## 78 20130510_EXQ1_IgPa_QC_UPS2_02.raw 6193937408          NA      False
##      DetectionQValue
## 73              #N/A
## 74              #N/A
## 75              #N/A
## 76              #N/A
## 77              #N/A
## 78              #N/A

```

Annotation information is required to fill in `Condition` and `BioReplicate` for corresponding `Run` information. Users have to prepare as csv or txt file like 'Cox_skyline_annotation.csv', which includes `Run`, `Condition`, and `BioReplicate` information, and load it in R.

```

annot <- read.csv("Cox_skyline_annotation.csv", header=TRUE)
annot

```

```

##              Run Condition BioReplicate
## 1 20130510_EXQ1_IgPa_QC_UPS1_01.raw      UPS1          1
## 2 20130510_EXQ1_IgPa_QC_UPS1_02.raw      UPS1          2
## 3 20130510_EXQ1_IgPa_QC_UPS1_03.raw      UPS1          3
## 4 20130510_EXQ1_IgPa_QC_UPS1_04.raw      UPS1          4
## 5 20130510_EXQ1_IgPa_QC_UPS2_01.raw      UPS2          5
## 6 20130510_EXQ1_IgPa_QC_UPS2_02.raw      UPS2          6
## 7 20130510_EXQ1_IgPa_QC_UPS2_03.raw      UPS2          7
## 8 20130510_EXQ1_IgPa_QC_UPS2_04.raw      UPS2          8

```

4.2.2 Preprocessing with DDA experiment from Skyline output

The input data for `MSstats` is required to contain variables of `ProteinName`, `PeptideSequence`, `PrecursorCharge`, `FragmentIon`, `ProductCharge`, `IsotopeLabelType`, `Condition`, `BioReplicate`, `Run`, `Intensity`. These variable names should be fixed. `MSstats` input from Skyline adapts the column scheme of the dataset so that it fits `MSstats` input format. However there are several extra column names and also some of them need to be changed. `SkylinetoMSstatsFormat` function helps pre-processing for making right format of `MSstats` input from Skyline output. For example, it renames some column name, and replace truncated peak intensities with NA. Another important step is to handle isotopic peaks before using `dataProcess`. The output from Skyline for DDA experiment has several measurements of peak area from the monoisotopic, M+1 and M+2 peaks. To get a robust measure of peptide intensity, we can sum over isotopic peaks per peptide or use the highest peak. Here we take a summation per peptide ion.

Here is the summary of pre-processing steps in `SkylinetoMSstatsFormat` function (in orange box below).

In Skyline

Remove duplicated rows : exactly same values in some rows
Remove decoy proteins
Remove protein which has only one peptide per protein

In MSstats

SkylinetoMSstatsFormat

- Rename column names
 - Replace NA for truncated rows
 - Sum of isotopic peaks per peptide and charge
-
- Replace intensity = 1 with zero if intensity < 1
 - Log 2 transform for intensity
 - Normalization
 - **Extra step for imputation** : Distinguish missing at random and censored missing
 - Decide cutoff for censored missing values among all $\log_2(\text{intensity}) > 0$.
 - If $\log_2(\text{intensity}) < \text{cutoff}$, $\log_2(\text{intensity})$ replaces with zero and is considered as censored missing values (**censoredInt='0'**), then, will be imputed.
 - NA will be remained as NA, which are missing at random.

```
# reformatting and pre-processing for Skyline output.  
quant <- SkylinetoMSstatsFormat(raw, annotation=annot)
```

```
## Peptides, that are used in more than one proteins, are removed.  
## Warning in SkylinetoMSstatsFormat(raw, annotation = annot): NAs introduced  
## by coercion  
## ** Truncated peaks are replaced with NA.  
## ** For DDA datasets, three isotopic peaks per feature and run are summed.
```

```
head(quant)
```

```
##           ProteinName  PeptideSequence PrecursorCharge FragmentIon  
## 1 P00915ups|CAH1_HUMAN_UPS  A[+42]SPDWGYDDK           2         sum  
## 2 sp|P0AC41|SDHA_ECOLI      IYQRPFGGQSK           2         sum  
## 3 sp|P04825|AMPN_ECOLI      EVALELYVDR           2         sum  
## 4 sp|P60240|RAPA_ECOLI      IGQAHDIIQIHVPYLEK      4         sum  
## 5 sp|P25524|CODA_ECOLI      FVETVAALAHHEGMGAR      3         sum  
## 6 sp|P04983|RBSA_ECOLI      GLMTRPK                2         sum  
##   ProductCharge IsotopeLabelType Condition BioReplicate  
## 1           NA           L           UPS1           1  
## 2           NA           L           UPS1           1  
## 3           NA           L           UPS1           1  
## 4           NA           L           UPS1           1  
## 5           NA           L           UPS1           1  
## 6           NA           L           UPS1           1  
##           Run Intensity  
## 1 20130510_EXQ1_IgPa_QC_UPS1_01.raw 41232179  
## 2 20130510_EXQ1_IgPa_QC_UPS1_01.raw 3011359328  
## 3 20130510_EXQ1_IgPa_QC_UPS1_01.raw 5848428  
## 4 20130510_EXQ1_IgPa_QC_UPS1_01.raw 203377208  
## 5 20130510_EXQ1_IgPa_QC_UPS1_01.raw 170613020  
## 6 20130510_EXQ1_IgPa_QC_UPS1_01.raw 202657636
```


For further details, visit the help file using the following code.

```
?SkylinetoMSstatsFormat
```

4.2.3 Different options for Skyline in dataProcess

The difference between output from Skyline and other spectral processing tool is that Skyline distinguishes random missing (NA) by technical issues and low noisy intensity due to less than limit of detection. The output from Skyline can have both NA (expect small number of NAs or none of them) and very small intensity close to zero (less than 1 in intensity) and those should be treated different types of missing. In `dataProcess`, users need to use `censoredInt='0'` for Skyline output, which means to distinguish between NA as random missing and 0 as censored missing.

```
cox.skyline.proposed <- dataProcess(quant,
  normalization='equalizeMedian',
  summaryMethod="TMP",
  cutoffCensored="minFeature",
  censoredInt="0",
  MBimpute=TRUE,
  maxQuantileforCensored=0.999)
```

Further steps is the same as in general workflow for DDA.

4.3 Suggested workflow with MaxQuant output for DDA

The following R code chunks show steps to format a MaxQuant output for analysis by `MSstats`. Here a controlled mixture dataset with dynamic range benchmark (J. Cox et al. 2014) is used for demonstration. This dataset is available in `MSstats` material GitHub under the folder named 'example dataset/DDA_controlledMixture20014'.

4.3.1 Load MaxQuant outputs

Three files should be prepared before `MSstats`. Two files, 'proteinGroups.txt' and 'evidence.txt' are outputs from MaxQuant.

```
## First, get protein ID information
proteinGroups <- read.table("Cox_maxquant_proteinGroups.txt", sep = "\t", header = TRUE)
```

```
## Read in MaxQuant file: evidence.txt
infile <- read.table("evidence.txt", sep = "\t", header = TRUE)
```

One file is for annotation information, required to fill in `Condition` and `BioReplicate` for corresponding Run information. Users have to prepare as csv or txt file like 'Cox_maxquant_annotation.csv', which includes Run, Condition, and BioReplicate information, and load it in R.

```
## Read in annotation including condition and biological replicates: annotation.csv
annot <- read.csv("Cox_maxquant_annotation.csv", header = TRUE)
annot
```

```
##           Raw.file Condition BioReplicate Experiment
## 1 20130510_EXQ1_IgPa_QC_UPS1_01      UPS1           1     UPS1_01
## 2 20130510_EXQ1_IgPa_QC_UPS1_02      UPS1           1     UPS1_02
## 3 20130510_EXQ1_IgPa_QC_UPS1_03      UPS1           1     UPS1_03
## 4 20130510_EXQ1_IgPa_QC_UPS1_04      UPS1           1     UPS1_04
## 5 20130510_EXQ1_IgPa_QC_UPS2_01      UPS2           2     UPS2_01
```

```
## 6 20130510_EXQ1_IgPa_QC_UPS2_02      UPS2      2      UPS2_02
## 7 20130510_EXQ1_IgPa_QC_UPS2_03      UPS2      2      UPS2_03
## 8 20130510_EXQ1_IgPa_QC_UPS2_04      UPS2      2      UPS2_04
##   IsotopeLabelType
## 1                      L
## 2                      L
## 3                      L
## 4                      L
## 5                      L
## 6                      L
## 7                      L
## 8                      L
```

4.3.2 Preprocessing with DDA experiment from MaxQuant output

`MaxQtoMSstatsFormat` function helps pre-processing for making right format of MSstats input from MaxQuant output. Basically, this function gets peptide ion intensity from 'evidence.txt' file. In addition, there are several steps to filter out or to modify the data in order to get required information.

Here is the summary of pre-processing steps in `MaxQtoMSstatsFormat` function (in orange box below).

In MSstats

MaxQtoMSstatsFormat

- Remove 'Contaminant', 'Potential.contaminant', 'Reverse', or 'Only.identified.by.site' protein ID.
 - Match protein group ID with 'Proteins' column in evidence.txt
 - Remove not unique peptides (assigned for more than one protein)
 - Use Maximum(highest) intensity if there are multiple measurements for feature and run.
 - Remove features that have 1 or 2 measurements across runs.
 - Remove peptides including 'M' sequence.
 - Remove proteins which have one peptide and charge in a protein.
-
- Replace intensity = 1 with zero if intensity < 1
 - Log 2 transform for intensity
 - Normalization
 - **Extra step for imputation**
 - Decide cutoff for censored missing values among all $\log_2(\text{intensity})$.
 - If $\log_2(\text{intensity}) < \text{cutoff}$, replaces with zero and is considered as censored missing values (`censoredInt='NA'`), then, will be imputed.

```
## check options for converting format
?MaxQtoMSstatsFormat
```

```
quant <- MaxQtoMSstatsFormat(evidence=infile, annotation=annot, proteinGroups=proteinGroups)
```

```
## + Contaminant, + Reverse, + Only.identified.by.site, proteins are removed.
```

```
## Peptides, that are used in more than one proteins, are removed.
```

```
## Peptide and charge, that have 1 or 2 measurements across runs, are removed.
```

```
## now 'quant' is ready for MSstats
head(quant)
```

```

## ProteinName PeptideSequence PrecursorCharge FragmentIon
## 1 A5A614 QVAESTPDIPK 2 NA
## 2 000762ups DPAATSVAAAR 2 NA
## 3 000762ups FLTPCYHPNVDTQGNICLDILK 2 NA
## 4 000762ups FLTPCYHPNVDTQGNICLDILK 3 NA
## 5 000762ups GAEPSSGAAR 2 NA
## 6 000762ups GISAFPESDNLFK 2 NA
## ProductCharge IsotopeLabelType Condition BioReplicate
## 1 NA L UPS1 1
## 2 NA L UPS1 1
## 3 NA L UPS1 1
## 4 NA L UPS1 1
## 5 NA L UPS1 1
## 6 NA L UPS1 1
## Run Intensity
## 1 20130510_EXQ1_IgPa_QC_UPS1_01 NA
## 2 20130510_EXQ1_IgPa_QC_UPS1_01 1144800000
## 3 20130510_EXQ1_IgPa_QC_UPS1_01 32793000
## 4 20130510_EXQ1_IgPa_QC_UPS1_01 566960000
## 5 20130510_EXQ1_IgPa_QC_UPS1_01 58709000
## 6 20130510_EXQ1_IgPa_QC_UPS1_01 861090000

```

4.3.3 Different options for MaxQuant in dataProcess

MaxQuant has certain or fixed threshold for intensity value internally as an parameter. Intensities less than the threshold are reported as NA. All missing values are NA in output from MaxQuant. In `dataProcess`, users need to use `censoredInt='NA'`. Users can use the same choice for other options.

```

cox.maxquant.proposed <- dataProcess(quant,
                                     normalization='equalizeMedian',
                                     summaryMethod="TMP",
                                     cutoffCensored="minFeature", censoredInt="NA",
                                     MBimpute=TRUE,
                                     maxQuantileforCensored=0.999)

```

Further steps is the same as in general workflow for DDA.

4.4 Suggested workflow with Progenesis output for DDA

This section describes steps and considerations to properly format data processed by Progenesis, prior to the `MSstats` analysis. In the following example, the same raw dataset from the previous section is used, but it is processed by Progenesis.

4.4.1 Load Progenesis output

Here is the expected input for `MSstats`, which is output of Progenesis.

```

## First, read output of Progenesis
raw <- read.csv("Cox_Progenesis.csv")
head(raw)

```

```

## X X.1 X.2 X.3 X.4
## 1

```

```

## 2 # Retention time (min) Charge m/z Measured mass
## 3 897 57.93265 2 459.27857242872 916.542591923679
## 4 1281 118.816733333333 2 1002.03351355306 2002.05247417235
## 5 3867 114.552433333333 2 1002.01296785898 2002.0113827842
## 6 1660 31.1707666666667 2 502.277007963734 1002.53946299371
## X.5 X.6 X.7 X.8
## 1
## 2 Mass error (u) Mass error (ppm) Score Sequence
## 3 0.0043889236793575 4.7885878242628 0.9992 VPYGAVLAK
## 4 0.0555781723510336 27.7613678932665 0.9956 LVITPVDGSDPYEEMIPK
## 5 0.0144867842016083 7.23616716417142 1 LVITPVDGSDPYEEMIPK
## 6 0.00488799370805282 4.87563604282956 1 AAAESSIQVK
## X.9 X.10
## 1
## 2 Modifications Accession
## 3 sp|P0A8T7|RPOC_ECOLI
## 4 sp|P0A8T7|RPOC_ECOLI
## 5 sp|P0A8T7|RPOC_ECOLI
## 6 sp|P0A8T7|RPOC_ECOLI
## X.11
## 1
## 2 Grouped accessions (for this sequence)
## 3
## 4
## 5
## 6
## X.12
## 1
## 2 Description
## 3 DNA-directed RNA polymerase subunit beta' OS=Escherichia coli (strain K12) GN=rpoC PE=1 SV=1
## 4 DNA-directed RNA polymerase subunit beta' OS=Escherichia coli (strain K12) GN=rpoC PE=1 SV=1
## 5 DNA-directed RNA polymerase subunit beta' OS=Escherichia coli (strain K12) GN=rpoC PE=1 SV=1
## 6 DNA-directed RNA polymerase subunit beta' OS=Escherichia coli (strain K12) GN=rpoC PE=1 SV=1
## X.13 X.14 X.15
## 1
## 2 Use in quantitation Max fold change Highest mean condition
## 3 False 1.07024295028376 Condition 1
## 4 False 1.26408361729762 Condition 2
## 5 True 1.12493122091942 Condition 2
## 6 False 2.08957830009021 Condition 1
## X.16 X.17 X.18
## 1
## 2 Lowest mean condition Anova Maximum CV
## 3 Condition 2 0.456859930184477 20.1184072001903
## 4 Condition 1 0.159939805145712 23.8614566196111
## 5 Condition 1 0.349651536742781 17.4104032083395
## 6 Condition 2 0.113361006195836 96.5523112559413
## Normalized.abundance X.19
## 1 Condition 1
## 2 20130510_EXQ1_IgPa_QC_UPS1_01 20130510_EXQ1_IgPa_QC_UPS1_02
## 3 20931810.5655776 21680597.2888134
## 4 160825738.322531 204844686.21804
## 5 73462123.4527211 95179635.7807487
## 6 35548647.7029835 30254160.104057

```

##		X.20		X.21
##	1			
##	2	20130510_EXQ1_IgPa_QC_UPS1_03	20130510_EXQ1_IgPa_QC_UPS1_04	
##	3	19010264.4603025	19287017.2266766	
##	4	159221868.441427	198325215.213719	
##	5	66486528.3804692	93200193.8860612	
##	6	21703274.487964	21391131.0188304	
##		X.22		X.23
##	1	Condition 2		
##	2	20130510_EXQ1_IgPa_QC_UPS2_01	20130510_EXQ1_IgPa_QC_UPS2_02	
##	3	16362526.9389495	15818549.690736	
##	4	232357617.95467	167457335.555413	
##	5	103977057.096685	74305266.8499236	
##	6	25651156.0484212	22037815.8546054	
##		X.24		X.25
##	1			
##	2	20130510_EXQ1_IgPa_QC_UPS2_03	20130510_EXQ1_IgPa_QC_UPS2_04	
##	3	24123356.4583506	19294933.8780514	
##	4	299234521.249582	215157929.093349	
##	5	88240964.864741	102823670.745066	
##	6	2576827.30283144	1848645.75613902	
##		Raw.abundance		X.26
##	1	Condition 1		
##	2	20130510_EXQ1_IgPa_QC_UPS1_01	20130510_EXQ1_IgPa_QC_UPS1_02	
##	3	15105416.9228044	19732517.8634299	
##	4	116059708.340506	186438656.471507	
##	5	53013856.5563334	86627404.1374054	
##	6	25653640.5636365	27535715.3100392	
##		X.27		X.28
##	1			
##	2	20130510_EXQ1_IgPa_QC_UPS1_03	20130510_EXQ1_IgPa_QC_UPS1_04	
##	3	13810320.3071721	19287017.2266766	
##	4	115669353.662823	198325215.213719	
##	5	48300235.6418326	93200193.8860612	
##	6	15766701.8793535	21391131.0188304	
##		X.29		X.30
##	1	Condition 2		
##	2	20130510_EXQ1_IgPa_QC_UPS2_01	20130510_EXQ1_IgPa_QC_UPS2_02	
##	3	16908752.5099689	12362738.4578786	
##	4	240114346.057948	130873644.09502	
##	5	107448093.544628	58072111.4178035	
##	6	26507461.2763444	17223308.0098951	
##		X.31		X.32
##	1			
##	2	20130510_EXQ1_IgPa_QC_UPS2_03	20130510_EXQ1_IgPa_QC_UPS2_04	
##	3	20505908.0138557	16847582.9106248	
##	4	254362429.950747	187867503.05488	
##	5	75008679.3143534	89781614.6456385	
##	6	2190415.9038019	1614165.29570777	
##		Spectral.counts		X.33
##	1	Condition 1		
##	2	20130510_EXQ1_IgPa_QC_UPS1_01	20130510_EXQ1_IgPa_QC_UPS1_02	
##	3	1	0	
##	4	0	0	

```

## 5          4          3
## 6          1          1
##          X.34          X.35
## 1
## 2 20130510_EXQ1_IgPa_QC_UPS1_03 20130510_EXQ1_IgPa_QC_UPS1_04
## 3          1          0
## 4          0          0
## 5          7          3
## 6          1          1
##          X.36          X.37
## 1          Condition 2
## 2 20130510_EXQ1_IgPa_QC_UPS2_01 20130510_EXQ1_IgPa_QC_UPS2_02
## 3          0          0
## 4          0          1
## 5          3          6
## 6          1          1
##          X.38          X.39
## 1
## 2 20130510_EXQ1_IgPa_QC_UPS2_03 20130510_EXQ1_IgPa_QC_UPS2_04
## 3          0          0
## 4          0          0
## 5          2          3
## 6          0          0

```

One file is for annotation information, required to fill in `Condition` and `BioReplicate` for corresponding `Run` information. Users have to prepare as csv or txt file like ‘`Cox_progenesis_annotation.csv`’, which includes `Run`, `Condition`, and `BioReplicate` information, and load it in R.

```

## Read in annotation including condition and biological replicates: annotation.csv
annot <- read.csv("Cox_Progenesis_annotation.csv", header = TRUE)
annot

```

```

##          Run Condition BioReplicate
## 1 20130510_EXQ1_IgPa_QC_UPS1_01      UPS1      1
## 2 20130510_EXQ1_IgPa_QC_UPS1_02      UPS1      2
## 3 20130510_EXQ1_IgPa_QC_UPS1_03      UPS1      3
## 4 20130510_EXQ1_IgPa_QC_UPS1_04      UPS1      4
## 5 20130510_EXQ1_IgPa_QC_UPS2_01      UPS2      5
## 6 20130510_EXQ1_IgPa_QC_UPS2_02      UPS2      6
## 7 20130510_EXQ1_IgPa_QC_UPS2_03      UPS2      7
## 8 20130510_EXQ1_IgPa_QC_UPS2_04      UPS2      8

```

4.4.2 Preprocessing with DDA experiment from progenesis output

The output from Progenesis includes peptide ion-level quantification for each MS runs. `ProgenisitoMSstatsFormat` function helps pre-processing for making right format of MSstats input from Progenesis output. Basically, this function reformats wide format to long format. It provide ‘`Raw.abundance`’, ‘`Normalized.abundance`’ and ‘`Spectral count`’ columns. This converter uses ‘`Raw.abundance`’ columns for Intensity values. In addition, there are several steps to filter out or to modify the data in order to get required information.

Here is the summary of pre-processing steps in `ProgenisitoMSstatsFormat` function (in orange box below).

In MSstats

ProgenesistoMSstatsFormat

- Use 'Raw.abundance'
 - Remove intensities with 'Use.in.quantitation' = FALSE.
 - Remove not unique peptides (assigned for more than one protein)
 - Use Maximum(highest) intensity if there are multiple measurements for feature and run.
 - Option : Remove features that have 1 or 2 measurements across runs.
 - Option : Remove proteins which have one peptide and charge in a protein.
 - Missing values are left with zero value.
-
- Replace intensity = 1 with zero if intensity < 1
 - Log 2 transform for intensity
 - Normalization
 - **Extra step for imputation**
 - Decide cutoff for censored missing values among all $\log_2(\text{intensity}) > 0$
 - If $\log_2(\text{intensity}) < \text{cutoff}$, $\log_2(\text{intensity})$ replaces with zero and is considered as censored missing values (`censoredInt='0'`), then, will be imputed.

```
## check options for converting format
?ProgenesistoMSstatsFormat
```

```
quant <- ProgenesistoMSstatsFormat(raw, annotation=annot)
```

```
## now 'quant' is ready for MSstats
head(quant)
```

```
##           ProteinName
## 1 000762ups|UBE2C_HUMAN_UPS
## 2 000762ups|UBE2C_HUMAN_UPS
## 3 000762ups|UBE2C_HUMAN_UPS
## 4 000762ups|UBE2C_HUMAN_UPS
## 5 000762ups|UBE2C_HUMAN_UPS
## 6 000762ups|UBE2C_HUMAN_UPS
##           PeptideModifiedSequence PrecursorCharge
## 1 FLTPCYHPNVDTQGNICLDILK [5] C+57.0215| [17] C+57.0215      2
## 2 FLTPCYHPNVDTQGNICLDILK [5] C+57.0215| [17] C+57.0215      3
## 3 FLTPCYHPNVDTQGNICLDILK [5] C+57.0215| [17] C+57.0215      4
## 4                               GISAFPESDNLFK          2
## 5                               WVGTIHGAAGTVYEDLR        2
## 6                               WVGTIHGAAGTVYEDLR        3
##   FragmentIon ProductCharge IsotopeLabelType Condition BioReplicate
## 1           NA           NA                L      UPS1           1
## 2           NA           NA                L      UPS1           1
## 3           NA           NA                L      UPS1           1
## 4           NA           NA                L      UPS1           1
## 5           NA           NA                L      UPS1           1
## 6           NA           NA                L      UPS1           1
##           Run Intensity
## 1 20130510_EXQ1_IgPa_QC_UPS1_01 3790142.3
## 2 20130510_EXQ1_IgPa_QC_UPS1_01 63386703.5
## 3 20130510_EXQ1_IgPa_QC_UPS1_01 165145.8
## 4 20130510_EXQ1_IgPa_QC_UPS1_01 98738397.1
```

```
## 5 20130510_EXQ1_IgPa_QC_UPS1_01 9624505.6
## 6 20130510_EXQ1_IgPa_QC_UPS1_01 12633723.8
```

4.4.3 Different options for Progenesis in dataProcess

Progenesis reports 0(zero) for missing values and does not have NA. Therefore, in `dataProcess`, users need to use `censoredInt='0'`. Users can use the same choice for other options.

```
cox.progenesis.proposed <- dataProcess(quant,
                                       normalization='equalizeMedian',
                                       summaryMethod="TMP",
                                       cutoffCensored="minFeature",
                                       censoredInt="0",
                                       MBimpute=TRUE,
                                       maxQuantileforCensored=0.999)
```

Further steps is the same as in general workflow for DDA.

4.5 Suggested workflow with Proteome Discoverer output for DDA

This section describes steps and considerations to properly format data processed by Proteome Discoverer, prior to the `MSstats` analysis. In the following example, another spike-in dataset processed by Proteome Discoverer is used to demonstrate.

4.5.1 Load Proteome Discoverer output

The output from Proteome Discoverer includes several level of datasets. PSM sheet should be saved as csv as below. Here is the expected input for `MSstats`.

```
## Read PSM-level data
raw <- read.csv("spikein_PD_psm.csv")
head(raw)
```

```
## Confidence.Level Search.ID Processing.Node.No Sequence
## 1 High A 4 AALGVLR
## 2 High A 4 NLLLVK
## 3 High A 4 LIVVEK
## 4 High A 4 LLVDLK
## 5 High A 4 IITLLK
## 6 High A 4 HEFLR
## Unique.Sequence.ID PSM.Ambiguity
## 1 2 Unambiguous
## 2 4 Unambiguous
## 3 5 Unambiguous
## 4 6 Unambiguous
## 5 9 Unambiguous
## 6 10 Unambiguous
##
## Protein.Descrip
## 1 Glycine--tRNA ligase beta subunit OS=Escherichia coli (strain K12) GN=glyS PE=1 SV=4 - [SYGB_E
## 2 50S ribosomal protein L3 OS=Escherichia coli (strain K12) GN=rp1C PE=1 SV=1 - [RL3_E
## 3 50S ribosomal protein L4 OS=Escherichia coli (strain K12) GN=rp1D PE=1 SV=1 - [RL4_E
## 4 Peptidyl-prolyl cis-trans isomerase D OS=Escherichia coli (strain K12) GN=ppiD PE=1 SV=1 - [PPID_E
## 5 3-dehydroquinate synthase OS=Escherichia coli (strain K12) GN=aroB PE=1 SV=1 - [AROB_E
```



```

## 6          GTP cyclohydrolase 1 OS=Escherichia coli (strain K12) GN=fo1E PE=1 SV=2 - [GCH1_E
## X..Proteins X..Protein.Groups Protein.Group.Accessions Modifications
## 1          1          1          P00961
## 2          1          1          P60438
## 3          1          1          P60723
## 4          1          1          POADY1
## 5          1          1          P07639
## 6          1          1          POA6T5
## Activation.Type DeltaScore DeltaCn Rank Search.Engine.Rank
## 1          CID      1.0000      0   1          1
## 2          CID      0.5455      0   1          1
## 3          CID      0.0000      0   1          1
## 4          CID      0.4062      0   1          1
## 5          CID      1.0000      0   1          1
## 6          CID      1.0000      0   1          1
## Precursor.Area QuanResultID Decoy.Peptides.Matched Exp.Value
## 1          3.77e+07      NA          11  0.00033
## 2          6.59e+08      NA           6  0.00940
## 3          3.83e+08      NA          17  0.20000
## 4          1.42e+07      NA           4  0.01300
## 5          3.93e+07      NA          NA  0.00860
## 6          2.80e+07      NA           7  0.27000
## Homology.Threshold Identity.High Identity.Middle IonScore
## 1          13          13          13      48
## 2          13          13          13      33
## 3          13          13          13      20
## 4          13          13          13      32
## 5          13          13          13      34
## 6          13          13          13      19
## Peptides.Matched X..Missed.Cleavages Isolation.Interference....
## 1          5          0          53
## 2          11         0           8
## 3          19         0          38
## 4          6          0          34
## 5          5          0          13
## 6          4          0          41
## Ion.Inject.Time..ms. Intensity Charge m.z..Da. MH...Da. Delta.Mass..Da.
## 1          4  1700000      2 350.2295 699.4517      0
## 2          2  2520000      2 350.2417 699.4761      0
## 3          5   739000      2 350.7340 700.4607      0
## 4          3  1520000      2 350.7342 700.4611      0
## 5          2  2480000      2 350.7520 700.4968      0
## 6          70   53500      2 351.1900 701.3728      0
## Delta.Mass..PPM. RT..min. First.Scan Last.Scan MS.Order Ions.Matched
## 1          0.68   32.17      8180      8180      MS2      Jun-50
## 2          -0.44   38.77     10907     10907      MS2      May-52
## 3          0.41   27.49      6221      6221      MS2      May-40
## 4          0.93   43.27     12766     12766      MS2      May-40
## 5          -0.03   42.75     12552     12552      MS2      Apr-40
## 6          -0.25   17.39      2693      2693      MS2      Apr-32
## Matched.Ions Total.Ions          Spectrum.File
## 1          6          50 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw
## 2          5          52 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw
## 3          5          40 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw

```

```
## 4          5          40 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw
## 5          4          40 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw
## 6          4          32 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw
## Annotation
## 1          NA
## 2          NA
## 3          NA
## 4          NA
## 5          NA
## 6          NA
```

One file is for annotation information, required to fill in `Condition` and `BioReplicate` for corresponding Run information. Users have to prepare as csv or txt file like 'spikein_PD_annotation.csv', which includes Run, Condition, and BioReplicate information, and load it in R.

```
## Read in annotation including condition and biological replicates: annotation.csv
annot <- read.csv("spikein_PD_annotation.csv", header = TRUE)
annot
```

```
##              Run Condition BioReplicate
## 1 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw Condition1      1
## 2 121219_S_CCES_01_02_LysC_Try_1to10_Mixt_1_2.raw Condition1      2
## 3 121219_S_CCES_01_03_LysC_Try_1to10_Mixt_1_3.raw Condition1      3
## 4 121219_S_CCES_01_04_LysC_Try_1to10_Mixt_2_1.raw Condition2      4
## 5 121219_S_CCES_01_05_LysC_Try_1to10_Mixt_2_2.raw Condition2      5
## 6 121219_S_CCES_01_06_LysC_Try_1to10_Mixt_2_3.raw Condition2      6
## 7 121219_S_CCES_01_07_LysC_Try_1to10_Mixt_3_1.raw Condition3      7
## 8 121219_S_CCES_01_08_LysC_Try_1to10_Mixt_3_2.raw Condition3      8
## 9 121219_S_CCES_01_09_LysC_Try_1to10_Mixt_3_3.raw Condition3      9
## 10 121219_S_CCES_01_10_LysC_Try_1to10_Mixt_4_1.raw Condition4     10
## 11 121219_S_CCES_01_11_LysC_Try_1to10_Mixt_4_2.raw Condition4     11
## 12 121219_S_CCES_01_12_LysC_Try_1to10_Mixt_4_3.raw Condition4     12
## 13 121219_S_CCES_01_13_LysC_Try_1to10_Mixt_5_1.raw Condition5     13
## 14 121219_S_CCES_01_14_LysC_Try_1to10_Mixt_5_2.raw Condition5     14
## 15 121219_S_CCES_01_15_LysC_Try_1to10_Mixt_5_3.raw Condition5     15
```

4.5.2 Preprocessing with DDA experiment from Proteome Discoverer output

`PDtoMSstatsFormat` function helps pre-processing for making right format of MSstats input from Proteome Discoverer output. `Protein.Group.Accessions` is used for `ProteinName`. The combination of `Sequence` and `Modifications` is used for `PeptideSequence`. `Charge` is used for `PrecursorCharge`. `Precursor.Area` is used for `Intensity`. In addition, there are several steps to filter out or to modify the data in order to get required information.

Here is the summary of pre-processing steps in `PDtoMSstatsFormat` function (in orange box below).

In MSstats

PDtoMSstatsFormat

- Get subset of useful columns
 - Remove peptides which are used in more than one protein (Use only unique peptides.)
 - Remove not unique peptides (assigned for more than one protein)
 - Use Maximum(highest) intensity if there are multiple measurements for feature and run.
 - Option : Remove features that have 1 or 2 measurements across runs.
 - Option : Remove proteins which have one peptide and charge in a protein.
 - Missing values are left with zero value.
-
- Replace intensity = 1 with zero if intensity < 1
 - Log 2 transform for intensity
 - Normalization
 - **Extra step for imputation**
 - Decide cutoff for censored missing values among all $\log_2(\text{intensity})$.
 - If $\log_2(\text{intensity}) < \text{cutoff}$, replaces with zero and is considered as censored missing values (`censoredInt='NA'`), then, will be imputed.

```
## check options for converting format
?PDtoMSstatsFormat
```

```
quant <- PDtoMSstatsFormat(raw, annotation=annot)
```

```
## ** Rows with #Proteins, which are not equal to 1, are removed.
```

```
## ** Peptides, that are used in more than one proteins, are removed.
```

```
## ** Multiple measurements in a feature and a run are summarized by summaryforMultipleRows.
```

```
## now 'quant' is ready for MSstats
```

```
head(quant)
```

```
## ProteinName PeptideModifiedSequence PrecursorCharge FragmentIon
## 1 P00961 AALGVLR_ 2 NA
## 2 P60438 NLLLVK_ 2 NA
## 3 P60723 LIVVEK_ 2 NA
## 4 POADY1 LLVDLK_ 2 NA
## 5 PO7639 IITLLK_ 2 NA
## 6 POA6T5 HEFLR_ 2 NA
## ProductCharge IsotopeLabelType Condition BioReplicate
## 1 NA L Condition1 1
## 2 NA L Condition1 1
## 3 NA L Condition1 1
## 4 NA L Condition1 1
## 5 NA L Condition1 1
## 6 NA L Condition1 1
## Run Intensity
## 1 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw 3.77e+07
## 2 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw 6.59e+08
## 3 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw 3.83e+08
## 4 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw 1.42e+07
## 5 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw 3.93e+07
## 6 121219_S_CCES_01_01_LysC_Try_1to10_Mixt_1_1.raw 2.80e+07
```

4.5.3 Different options for Proteome Discoverer in dataProcess

Progenesis reports NA for missing values. Therefore, in `dataProcess`, users need to use `censoredInt='NA'`. Users can use the same choice for other options.

```
cox.progenesis.proposed <- dataProcess(quant,
  normalization='equalizeMedian',
  summaryMethod="TMP",
  cutoffCensored="minFeature",
  censoredInt="NA",
  MBimpute=TRUE,
  maxQuantileforCensored=0.999)
```

Further steps is the same as in general workflow for DDA.

5. DIA analysis with MSstats

5.1 Suggested workflow with Skyline output for DIA

This section describes steps and considerations to properly format data processed by Skyline for SWATH/DIA experiments, prior to the `MSstats` analysis. In the following example, the raw files for profiling standard sample set (Bruderer et al. 2015) are quantified by Skyline.

5.1.1 Load Skyline output

This required input data is generated automatically when using `MSstats` report format in Skyline. We first load and access the dataset processed by Skyline. The name of saved file from Skyline using `MSstats` report format is 'Cox.Skyline.csv'.

```
# Read output from skyline : Bruderer.skyline.csv
raw <- read.csv("Bruderer.skyline.csv")
```

We can read csv file. Here we will load R data file which is the exactly same data in Cox.skyline.csv file.

```
# Load R data, which is converted from csv file, output from skyline : Bruderer.skyline.csv
load("Bruderer.skyline.RData")
raw <- Bruderer.skyline
```

Annotation information in `Condition` and `BioReplicate` for corresponding `Run` was already filled in Skyline > Result grid. If not, users have to prepare as csv or txt file, which includes `Run`, `Condition`, and `BioReplicate` information, and load it in R as section 4.2.1.

5.1.2 Preprocessing with DIA experiment from Skyline output

The input data for `MSstats` is required to contain variables of `ProteinName`, `PeptideSequence`, `PrecursorCharge`, `FragmentIon`, `ProductCharge`, `IsotopeLabelType`, `Condition`, `BioReplicate`, `Run`, `Intensity`. These variable names should be fixed. `MSstats` input from Skyline adapts the column scheme of the dataset so that it fits `MSstats` input format. However there are several extra column names and also some of them need to be changed.

```
## ProteinName PeptideSequence PeptideModifiedSequence
## 1 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK ELASQPDVDGFLVGGASLKPEFVDIINAK
```

```

## 2      P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK ELASQPDVDGFLVGGASLKPEFVDIINAK
## 3      P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK ELASQPDVDGFLVGGASLKPEFVDIINAK
## 4      P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK ELASQPDVDGFLVGGASLKPEFVDIINAK
## 5      P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK ELASQPDVDGFLVGGASLKPEFVDIINAK
## 6      P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK ELASQPDVDGFLVGGASLKPEFVDIINAK
##      PrecursorCharge PrecursorMz FragmentIon ProductCharge ProductMz
## 1          3      1010.533          y10          1      1145.62
## 2          3      1010.533          y10          1      1145.62
## 3          3      1010.533          y10          1      1145.62
## 4          3      1010.533          y10          1      1145.62
## 5          3      1010.533          y10          1      1145.62
## 6          3      1010.533          y10          1      1145.62
##      IsotopeLabelType Condition BioReplicate
## 1          light          S1          S1
## 2          light          S1          S1
## 3          light          S1          S1
## 4          light          S2          S2
## 5          light          S2          S2
## 6          light          S2          S2
##
##      FileName          Area StandardType Truncated
## 1 B_D140314_SGSDSsample1_R01_MHRM_T0.raw 17578982          False
## 2 B_D140314_SGSDSsample1_R02_MHRM_T0.raw 19800498          False
## 3 B_D140314_SGSDSsample1_R03_MHRM_T0.raw 16162569          False
## 4 B_D140314_SGSDSsample2_R01_MHRM_T0.raw 19254086          False
## 5 B_D140314_SGSDSsample2_R02_MHRM_T0.raw 16377574          False
## 6 B_D140314_SGSDSsample2_R03_MHRM_T0.raw 15045770          False
##
##      DetectionQValue
## 1 3.5156850231032877E-07
## 2 2.2968222879171662E-07
## 3 1.0004539490182651E-06
## 4 3.2378503078689391E-07
## 5 2.3675326588090684E-07
## 6 9.7241468210995663E-07

```

SkylinetoMSstatsFormat function helps pre-processing for making right format of MSstats input from Skyline output. For example, it removes iRT protein, renames some column name, and replace truncated peak intensities with NA. Another important step for SWATH/DIA experiment is to use q-value (column named DetectionQValue) for filtering data before using dataProcess. The option, filter_with_qvalue=TRUE, will replace Intensity value with zero for the rows with DetectionQValue column value greater than qvalue_cutoff option value in SkylinetoMSstatsFormat function.

Here is the summary of pre-processing steps for SWATH/DIA experiment in SkylinetoMSstatsFormat function (in orange box below).

In Skyline

Remove duplicated rows : exactly same values in some rows
Remove protein which has only one peptide per protein

In MSstats

SkylinetoMSstatsFormat

- Rename column names
 - Remove iRT proteins
 - Replace NA for truncated rows
 - Replace with zero if q-value > 0.01
-
- Replace intensity = 1 with zero if intensity < 1
 - Log 2 transform for intensity
 - Normalization
 - **Extra step for imputation** : Distinguish missing at random and censored missing
 - Decide cutoff for censored missing values among all $\log_2(\text{intensity}) > 0$.
 - If $\log_2(\text{intensity}) < \text{cutoff}$, $\log_2(\text{intensity})$ replaces with zero and is considered as censored missing values (**censoredInt='0'**), then, will be imputed.
 - NA will be remained as NA, which are missing at random.

```
## check options for converting format  
?SkylinetoMSstatsFormat
```

```
quant <- SkylinetoMSstatsFormat(raw)
```

```
## ** iRT proteins/peptides are removed.
```

```
## Peptides, that are used in more than one proteins, are removed.
```

```
## Warning in SkylinetoMSstatsFormat(raw): NAs introduced by coercion
```

```
## ** Truncated peaks are replaced with NA.
```

```
## Warning in SkylinetoMSstatsFormat(raw): NAs introduced by coercion
```

```
## Intensities with great than 0.01 in DetectionQValue are replaced with zero.
```

```
## now 'quant' is ready for MSstats
```

```
head(quant)
```

```
## ProteinName PeptideSequence PrecursorCharge PrecursorMz  
## 1 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK 3 1010.533  
## 2 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK 3 1010.533  
## 3 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK 3 1010.533  
## 4 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK 3 1010.533  
## 5 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK 3 1010.533  
## 6 P60174 ELASQPDVDGFLVGGASLKPEFVDIINAK 3 1010.533  
## FragmentIon ProductCharge ProductMz IsotopeLabelType Condition  
## 1 y10 1 1145.62 light S1  
## 2 y10 1 1145.62 light S1  
## 3 y10 1 1145.62 light S1  
## 4 y10 1 1145.62 light S2  
## 5 y10 1 1145.62 light S2  
## 6 y10 1 1145.62 light S2  
## BioReplicate Run Intensity  
## 1 S1 B_D140314_SGSDSsample1_R01_MHRM_TO.raw 17578982  
## 2 S1 B_D140314_SGSDSsample1_R02_MHRM_TO.raw 19800498
```

```
## 3          S1 B_D140314_SGSDSsample1_R03_MHRM_TO.raw 16162569
## 4          S2 B_D140314_SGSDSsample2_R01_MHRM_TO.raw 19254086
## 5          S2 B_D140314_SGSDSsample2_R02_MHRM_TO.raw 16377574
## 6          S2 B_D140314_SGSDSsample2_R03_MHRM_TO.raw 15045770
## StandardType Truncated DetectionQValue
## 1                False 3.515685e-07
## 2                False 2.296822e-07
## 3                False 1.000454e-06
## 4                False 3.237850e-07
## 5                False 2.367533e-07
## 6                False 9.724147e-07
```

5.1.3 Different options for Skyline output of DIA experiment in dataProcess

In `dataProcess`, users need to use `censoredInt='0'` for Skyline output, which means to distinguish between NA as random missing and 0 as censored missing as described in section 4.2.3. The same options of summarization method and imputation for DDA experiments (section 4.2.3) are recommended for SWATH/DIA experiments. `featureSubset` option for using subset of features can be used for SWATH/DIA experiments, which have relatively large number of features in each protein.

```
bruderer.skyline.proposed <- dataProcess(quant,
                                         normalization='equalizeMedian',
                                         summaryMethod="TMP",
                                         cutoffCensored="minFeature", censoredInt="0",
                                         MBimpute=TRUE,
                                         maxQuantileforCensored=0.999)
```

Further steps is the same as in general workflow for DDA.

5.2 Suggested workflow with Spectronaut output for DIA

This section describes steps and considerations to properly format data processed by Spectronaut for SWATH/DIA experiments, prior to the `MSstats` analysis. In the following example, the same raw files in section 5.2 for profiling standard sample set (Bruderer et al. 2015) are quantified by Spectronaut.

5.2.1 Load Spectronaut output

We first load and access the dataset processed by Spectronaut.

```
# Read output from skyline : Bruderer.spectronaut.xls
raw <- read.table("Bruderer.spectronaut.xls", header=TRUE)
```

We can read the file as above. Here instead, we will load R data file which is the exactly same data in `Bruderer.spectronaut.xls` file.

```
# Load R data, which is converted, output from spectronaut : Bruderer.SN.RData
load("Bruderer.SN.RData")
raw <- Bruderer.SN
```

Annotation information should be filled by Spectronaut.

5.2.2 Preprocessing with DIA experiment from Spectronaut output

The output from Spectronaut should look like below.

head(raw)

```
##      R.Condition                R.FileName R.Replicate
## 1 SGSDSsample1_B_D140314_SGSDSsample1_RO1_MHRM      1
## 2 SGSDSsample1_B_D140314_SGSDSsample1_RO1_MHRM      1
## 3 SGSDSsample1_B_D140314_SGSDSsample1_RO1_MHRM      1
## 4 SGSDSsample1_B_D140314_SGSDSsample1_RO1_MHRM      1
## 5 SGSDSsample1_B_D140314_SGSDSsample1_RO1_MHRM      1
## 6 SGSDSsample1_B_D140314_SGSDSsample1_RO1_MHRM      1
##      PG.ProteinAccessions PG.ProteinGroups PG.Quantity PEP.GroupingKey
## 1                A0AVT1                A0AVT1    13667.06      VVQTDE TAR
## 2                A0AVT1                A0AVT1    13667.06      VVQTDE TAR
## 3                A0AVT1                A0AVT1    13667.06      VVQTDE TAR
## 4                A0AVT1                A0AVT1    13667.06      VVQTDE TAR
## 5                A0AVT1                A0AVT1    13667.06      VVQTDE TAR
## 6                A0AVT1                A0AVT1    13667.06      VVQTDE TAR
##      PEP.StrippedSequence PEP.Quantity EG.iRTPredicted
## 1                VVQTDE TAR          24986.3          -44.6755
## 2                VVQTDE TAR          24986.3          -44.6755
## 3                VVQTDE TAR          24986.3          -44.6755
## 4                VVQTDE TAR          24986.3          -44.6755
## 5                VVQTDE TAR          24986.3          -44.6755
## 6                VVQTDE TAR          24986.3          -44.6755
##      EG.Library EG.ModifiedSequence EG.PrecursorId
## 1 GSDS HEK293 - 24 Runs (SN10) Fixed      _VVQTDE TAR_      _VVQTDE TAR_.2
## 2 GSDS HEK293 - 24 Runs (SN10) Fixed      _VVQTDE TAR_      _VVQTDE TAR_.2
## 3 GSDS HEK293 - 24 Runs (SN10) Fixed      _VVQTDE TAR_      _VVQTDE TAR_.2
## 4 GSDS HEK293 - 24 Runs (SN10) Fixed      _VVQTDE TAR_      _VVQTDE TAR_.2
## 5 GSDS HEK293 - 24 Runs (SN10) Fixed      _VVQTDE TAR_      _VVQTDE TAR_.2
## 6 GSDS HEK293 - 24 Runs (SN10) Fixed      _VVQTDE TAR_      _VVQTDE TAR_.2
##      EG.Qvalue FG.Charge      FG.Id FG.PrecMz FG.Quantity F.Charge
## 1 0.004752379          2 _VVQTDE TAR_.2 509.7618    24986.3      1
## 2 0.004752379          2 _VVQTDE TAR_.2 509.7618    24986.3      1
## 3 0.004752379          2 _VVQTDE TAR_.2 509.7618    24986.3      1
## 4 0.004752379          2 _VVQTDE TAR_.2 509.7618    24986.3      1
## 5 0.004752379          2 _VVQTDE TAR_.2 509.7618    24986.3      1
## 6 0.004752379          2 _VVQTDE TAR_.2 509.7618    24986.3      1
##      F.FrgIon F.FrgLossType F.FrgMz F.FrgNum F.FrgType
## 1          y7          no loss 820.3795          7          y
## 2          y6          no loss 692.3210          6          y
## 3          y4          no loss 476.2463          4          y
## 4          y7           NH3 803.3530          7          y
## 5          y4           H2O 458.2358          4          y
## 6          y5          no loss 591.2733          5          y
##      F.ExcludedFromQuantification F.NormalizedPeakArea F.NormalizedPeakHeight
## 1                False          20440.797593          116995.869
## 2                True          1699.219836          5008.674
## 3                True           1.778585          1.000
## 4                False          575.263595          2806.576
## 5                True          2138.069289          11731.250
## 6                False          3970.243087          16882.633
##      F.PeakArea F.PeakHeight
## 1 15882.360352 9.090499e+04
## 2 1320.282227 3.891706e+03
```



```
## 3      1.381949 6.769939e-01
## 4     446.975891 2.180690e+03
## 5     1661.265259 9.115102e+03
## 6     3084.851807 1.311769e+04
```

The input data for MSstats is required to contain variables of ProteinName, PeptideSequence, PrecursorCharge, FragmentIon, ProductCharge, IsotopeLabelType, Condition, BioReplicate, Run, Intensity. These variable names should be fixed. Therefore, we need to get subset of useful columns and to rename them. Also several filtering steps are required. SpectronauttoMSstatsFormat function helps pre-processing for making right format of MSstats input from Spectronaut output. First, it uses only noLoss from F.FrgLossType. If not, multiple measurements for each feature and run can be happen. Spectronaut provides the column named F.ExcludedFromQuantification based on XIC quality such as interference between chromatographies. Only features with F.ExcludedFromQuantification == 'False' should be used. PG.ProteinGroups is used for ProteinName. EG.ModifiedSequence is used for PeptideSequence. FG.Charge is used for PrecursorCharge. F.FrgIon is used for FragmentIon. F.Charge is used for ProductCharge. F.PeakArea with default option is used for Intensity. Then several filtering steps will be performed.

Here is the summary of pre-processing steps for SWATH/DIA experiment in SpectronauttoMSstatsFormat function (in orange box below).

In MSstats

SpectronauttoMSstatsFormat

- Use only 'F.FrgLossType = 'noLoss'
- Remove 'F.ExcludedFromQuantification=True'.
- Use 'F.PeakArea' for Intensity value
- Replace with zero if q-value > 0.01
- Remove features that include all intensities = NA or zero across all MS runs.
- Option : Remove features that have 1 or 2 measurements across runs.
- Option : Remove proteins which have one peptide and charge in a protein.
- Use Maximum(highest) intensity if there are multiple measurements for feature and run.

- Replace intensity = 1 with zero if intensity < 1
- Log 2 transform for intensity
- Normalization
- **Extra step for imputation**
 - Decide cutoff for censored missing values among all log2(intensity).
 - If log2(intensity) < cutoff, replaces with zero and is considered as censored missing values (censoredInt='0'), then, will be imputed.

```
## check options for converting format
?SpectronauttoMSstatsFormat

quant <- SpectronauttoMSstatsFormat(raw)

## ** Intensities with great than 0.01 in EG.Qvalue are replaced with zero.
## ** All peptides are unique peptides in proteins.
## ** No multiple measurements in a feature and a run.

## now 'quant' is ready for MSstats
head(quant)

##      ProteinName      PeptideSequence PrecursorCharge FragmentIon
## 1      A0AVT1          _VVQTDETAR_                2             y7
```

```

## 6      AOAVT1      _VVQTDE TAR_      2      y5
## 10     AOAVT1      _LATS ISETLEEK_    2      y8
## 11     AOAVT1      _LATS ISETLEEK_    2      y4
## 13     AOAVT1 _VC[+57]PTTETIYNDEFYTK_  2      y8
## 14     AOAVT1 _VC[+57]PTTETIYNDEFYTK_  2      y14
##      ProductCharge      Condition      Run BioReplicate
## 1      1      SGSDSsample1 B_D140314_SGSDSsample1_R01_MHRM      1
## 6      1      SGSDSsample1 B_D140314_SGSDSsample1_R01_MHRM      1
## 10     1      SGSDSsample1 B_D140314_SGSDSsample1_R01_MHRM      1
## 11     1      SGSDSsample1 B_D140314_SGSDSsample1_R01_MHRM      1
## 13     1      SGSDSsample1 B_D140314_SGSDSsample1_R01_MHRM      1
## 14     2      SGSDSsample1 B_D140314_SGSDSsample1_R01_MHRM      1
##      Intensity IsotopeLabelType
## 1      15882.3604      L
## 6      3084.8518      L
## 10     182.7725      L
## 11     667.2112      L
## 13     2555.1445      L
## 14     1747.8701      L

```

5.2.3 Different options for Spectronaut output of DIA experiment in dataProcess

In `dataProcess`, users need to use `censoredInt='0'` for Spectronaut output. Spectronaut output generates very few number of NA. After applying Qvalue, zero intensities will be generated and those should be imputed. The same options of summarization method and imputation for DDA experiments (section 4.2.3) are recommended for SWATH/DIA experiments. `featureSubset` option for using subset of features can be used for SWATH/DIA experiments, which have relatively large number of features in each protein.

```

bruderer.spectronaut.proposed <- dataProcess(quant,
      normalization='equalizeMedian',
      summaryMethod="TMP",
      cutoffCensored="minFeature",
      censoredInt="0",
      MBimpute=TRUE,
      maxQuantileforCensored=0.999)

```

Further steps is the same as in general workflow for DDA.

5.3 Suggested workflow with OpenSWATH output for SWATH

This section describes steps and considerations to properly format data processed by OpenSWATH for SWATH experiments, prior to the `MSstats` analysis. In the following example, the dataset processed and quantified by OpenSWATH and available as supplementary in (Röst et al. 2014) is used.

5.3.1 Load OpenSWATH output and reformat using SWATH2stats package

R package, `SWATH2stats`, in Bioconductor reformats SWATH data from OpenSWATH software for `MSstats` input format. (Blattmann, Heusel, and Aebersold 2016)

```

library(SWATH2stats)

# Read the data
raw <- read.table('OpenSWATH_SM3_GoldStandardAutomatedResults_human_peakgroups.txt',

```

```

sep="\t", header=TRUE)

# Users should prepare the file including, FileName, Condition, BioReplicate and Run.
# and Read it.
annot <- read.csv('DIA_Rost2014_annotation.csv')
head(annot)

```

```

##                               Filename Condition BioReplicate
## 1 split_napedro_L120417_001_SW_combined.featureXML          1          1
## 2 split_napedro_L120417_002_SW_combined.featureXML          2          2
## 3 split_napedro_L120417_003_SW_combined.featureXML          3          3
## 4 split_napedro_L120417_004_SW_combined.featureXML          4          4
## 5 split_napedro_L120417_005_SW_combined.featureXML          5          5
## 6 split_napedro_L120417_006_SW_combined.featureXML          6          6
##                               Run
## 1 split_napedro_L120417_001_SW_combined.featureXML
## 2 split_napedro_L120417_002_SW_combined.featureXML
## 3 split_napedro_L120417_003_SW_combined.featureXML
## 4 split_napedro_L120417_004_SW_combined.featureXML
## 5 split_napedro_L120417_005_SW_combined.featureXML
## 6 split_napedro_L120417_006_SW_combined.featureXML

```

```
colnames(raw)[colnames(raw) == 'filename'] <- 'Filename'
```

```

# Fist, match annotation information.
raw2 <- sample_annotation(data=raw,
                           sample.annotation=annot,
                           data.type='OpenSWATH',
                           column.file='Filename')

# Filter with mscore threshold 0.01 and remove decoys
data.filtered <- filter_mscore(raw2, 0.01)

```

```
## Dimension difference: 6339, 0
```

```
# Let's check the first three rows.
```

```
data.filtered[1:3, ]
```

```

##           ProteinName           FullPeptideName Charge
## 11    AQUA4SWATH_HMLangeC    LDASLPALLLIR(UniMod:267)      2
## 101   AQUA4SWATH_MouseSabido  QEPAAPSLSPAVSAK(UniMod:259)      2
## 170   AQUA4SWATH_Lepto       AIAEEVPK(UniMod:259)      2
##
## 11    AQUA4SWATH_HMLangeC_LDASLPALLLIR(UniMod:267)/2_y9;AQUA4SWATH_HMLangeC_
## 101   AQUA4SWATH_MouseSabido_QEPAAPSLSPAVSAK(UniMod:259)/2_y10;AQUA4SWATH_MouseSabido_QEPAAPSLSPAVSAK(
## 170   AQUA4SWATH_Lepto_AIAEEVPK(UniMod:259)/2_y
##           agr_Peak_Area Condition
## 11    0.000000;0.000000;0.000000;130.000000      1
## 101   150.000000;0.000000;50.000000;40.000000      1
## 170   4239.000000;15066.000000;44784.000000;4753.000000      1
##           BioReplicate           Run
## 11    1 split_napedro_L120417_001_SW_combined.featureXML
## 101   1 split_napedro_L120417_001_SW_combined.featureXML
## 170   1 split_napedro_L120417_001_SW_combined.featureXML
##

```

transition

```

## 11      AQUA4SWATH_HMLangeC_LDASLPALLLIR(UniMod:267)/2_run0_split_napedro_L120417_001_SW_combined.
## 101 AQUA4SWATH_MouseSabido_QEPAAPSLSPAVSAK(UniMod:259)/2_run0_split_napedro_L120417_001_SW_combined.
## 170      AQUA4SWATH_Lepto_AIAEEVPK(UniMod:259)/2_run0_split_napedro_L120417_001_SW_combined.
##      decoy main_var_xx_swath_prelim_score var_bseries_score
## 11 FALSE                1.053124                1
## 101 FALSE                1.407187                1
## 170 FALSE                2.592268                4
##      var_elution_model_fit_score var_intensity_score
## 11                -0.5072533                0.022184300
## 101                0.4999997                0.007779326
## 170                0.9660140                0.106216673
##      var_isotope_correlation_score var_isotope_overlap_score
## 11                0.8245458                0.00000000
## 101                0.7895422                0.00000000
## 170                0.9110748                0.06904215
##      var_library_corr var_library_rmsd var_log_sn_score var_massdev_score
## 11                0.9989755                0.06604379                2.014903                4.633929
## 101                0.9768712                0.04136264                1.230344                2.806702
## 170                0.9688958                0.12375189                2.521098                6.274398
##      var_massdev_score_weighted var_norm_rt_score var_xcorr_coelution
## 11                16.087377                0.02200757                3.648683
## 101                3.424393                0.05639886                1.832796
## 170                6.682840                0.03075860                1.859502
##      var_xcorr_coelution_weighted var_xcorr_shape var_xcorr_shape_weighted
## 11                0.7401841                0.1000000                0.7532720
## 101                0.2552277                0.6000000                0.8723862
## 170                0.2428034                0.8764185                0.9650843
##      var_yseries_score
## 11                1
## 101                1
## 170                3
##
##
## 11      AQUA4SWATH_HMLangeC_LDASLPALLLIR(UniMod:267)/2_run0_split_napedro_L120417_001_SW_combined.
## 101 AQUA4SWATH_MouseSabido_QEPAAPSLSPAVSAK(UniMod:259)/2_run0_split_napedro_L120417_001_SW_combined.
## 170      AQUA4SWATH_Lepto_AIAEEVPK(UniMod:259)/2_run0_split_napedro_L120417_001_SW_combined.
##
##      run_id
## 11  1_1_split_napedro_L120417_001_SW_combined.featureXML
## 101 1_1_split_napedro_L120417_001_SW_combined.featureXML
## 170 1_1_split_napedro_L120417_001_SW_combined.featureXML
##
##      Filename      RT
## 11  split_napedro_L120417_001_SW_combined.featureXML 5706.057
## 101 split_napedro_L120417_001_SW_combined.featureXML 2130.896
## 170 split_napedro_L120417_001_SW_combined.featureXML 1632.214
##
##      id      Sequence      m.z Intensity assay_rt
## 11  f_11766711863966126076      LDASLPALLLIR 652.912      130 5620.726
## 101 f_7536148822812640906 QEPAAPSLSPAVSAK 730.895      240 2326.711
## 170 f_6525552705944724733      AIAEEVPK 432.749      68842 1530.658
##
##      delta_rt leftWidth      norm_RT nr_peaks peak_apices_sum rightWidth
## 11      85.33057      5702.47 116.80076      4      70      5716.13
## 101 -195.81477      2128.42 12.96011      4      150      2135.25
## 170 101.55614      1618.53 -1.52414      4      14452      1652.67
##
##      rt_score sn_ratio total_xic dotprod_score library_dotprod
## 11  2.200757  7.500000      5860      0.4210968      0.9316182
## 101 5.639886  3.422408      30851      0.7415526      0.9619197

```

```
## 170 3.075860 12.442253 648128 0.8172930 0.9471398
## library_manhattan manhatt_score xx_lda_prelim_score
## 11 0.7964739 1.2636574 2.647497
## 101 0.2837612 0.7505515 3.030819
## 170 0.3700136 0.6605900 4.604896
## xx_swath_prelim_score aggr_Peak_Apex log10_total_xic LD1
## 11 0 NA;NA;NA;NA 3.767898 1.978064
## 101 0 NA;NA;NA;NA 4.489269 1.245393
## 170 0 NA;NA;NA;NA 5.811661 1.149967
## peak_group_rank d_score m_score
## 11 1 3.399867 0.0007289805
## 101 1 2.649140 0.0071091755
## 170 1 2.551362 0.0093207975
```

```
data.transition <- disaggregate(data.filtered)
```

```
## The library contains between 3 and 4 transitions per precursor.
```

```
## The data table was transformed into a table containing one row per transition.
```

```
## 210 row(s) was/were removed because they did not contain data due to different number of transitions
```

```
# convert4MSstats : the function to convert into MSstats required format
```

```
MSstats.input <- convert4MSstats(data.transition)
```

```
## One or several columns required by MSstats were not in the data. The columns were created and filled
```

```
## Missing columns: ProductCharge, IsotopeLabelType
```

```
## IsotopeLabelType was filled with light.
```

```
## Warning in convert4MSstats(data.transition): Intensity values that were 0,
## were replaced by NA
```

```
## now 'MSstats.input' is ready for MSstats
```

```
head(MSstats.input)
```

```
## ProteinName PeptideSequence PrecursorCharge
## 1 AQUA4SWATH_HMLangeC LDASLPALLLIR(UniMod_267) 2
## 2 AQUA4SWATH_MouseSabido QEPAAPSLSPAVSAK(UniMod_259) 2
## 3 AQUA4SWATH_Lepto AIAEEVPK(UniMod_259) 2
## 4 AQUA4SWATH_HMLangeC LDASLPALLLIR(UniMod_267) 2
## 5 AQUA4SWATH_Lepto ILELPTEVDSEK(UniMod_259) 2
## 6 AQUA4SWATH_Lepto AIAEEVPK(UniMod_259) 2
## FragmentIon ProductCharge
## 1 AQUA4SWATH_HMLangeC_LDASLPALLLIR(UniMod_267)/2_y9 NA
## 2 AQUA4SWATH_MouseSabido_QEPAAPSLSPAVSAK(UniMod_259)/2_y10 NA
## 3 AQUA4SWATH_Lepto_AIAEEVPK(UniMod_259)/2_y7 NA
## 4 AQUA4SWATH_HMLangeC_LDASLPALLLIR(UniMod_267)/2_y9 NA
## 5 AQUA4SWATH_Lepto_ILELPTEVDSEK(UniMod_259)/2_b2 NA
## 6 AQUA4SWATH_Lepto_AIAEEVPK(UniMod_259)/2_y7 NA
## IsotopeLabelType Intensity BioReplicate Condition
## 1 light NA 1 1
## 2 light 150 1 1
## 3 light 4239 1 1
## 4 light NA 2 2
## 5 light 14185 2 2
## 6 light 4092 2 2
## Run
## 1 split_napedro_L120417_001_SW_combined.featureXML
```

```
## 2 split_napedro_L120417_001_SW_combined.featureXML
## 3 split_napedro_L120417_001_SW_combined.featureXML
## 4 split_napedro_L120417_002_SW_combined.featureXML
## 5 split_napedro_L120417_002_SW_combined.featureXML
## 6 split_napedro_L120417_002_SW_combined.featureXML
```

5.3.2 Different options for OpenSWATH output of DIA experiment in dataProcess

In `dataProcess`, users need to use `censoredInt='NA'` for OpenSWATH output. The same options of summarization method and imputation for DDA experiments (section 4.2.3) are recommended for SWATH/DIA experiments. `featureSubset` option for using subset of features can be used for SWATH/DIA experiments, which have relatively large number of features in each protein.

```
goldstandard.proposed <- dataProcess(MSstats.input,
                                     normalization='equalizeMedian',
                                     summaryMethod="TMP",
                                     cutoffCensored="minFeature", censoredInt="NA",
                                     MBimpute=TRUE,
                                     maxQuantileforCensored=0.999)
```

Further steps is the same as in general workflow for DDA.

Reference

- Benjamini, Y., and Y. Hochberg. 1955. "Controlling the false discovery rate: a practical and powerful approach to multiple testing." *J.R. Statist. Soc. B* 57 (1): 289–300.
- Blattmann, P., M. Heusel, and R. Aebersold. 2016. "SWATH2stats: An R/Bioconductor Package to Process and Convert Quantitative SWATH-MS Proteomics Data for Downstream Analysis Tools." *PLoS ONE* 11 (4). doi:10.1371/journal.pone.0153160.
- Bruderer, R., O. M. Bernhardt, T. Gandhi, S. M. Miladinović, L.-Y. Cheng, S. Messner, T. Ehrenberger, et al. 2015. "Extending the limits of quantitative proteome profiling with Data-Independent Acquisition and application to acetaminophen-treated three-dimensional liver microtissues." *Mol. Cell. Proteomics* 14 (5): 1400–1410.
- Cox, J., M. Y. Hein, C. A. Lubner, I. Paron, N. Nagaraj, and M. Mann. 2014. "Accurate proteome-wide label-free quantification by delayed normalization and maximal peptide ratio extraction, termed MaxLFQ." *Mol. Cell. Proteomics* 13 (9): 2513–26.
- Cox, Jürgen, and Matthias Mann. 2008. "MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification." *Nature Biotechnology* 26 (12): 1367–72.
- Gatto, L., and K. S. Lilley. 2012. "MSnbase-an R/Bioconductor package for isobaric tagged mass spectrometry data visualization, processing and quantitation." *Bioinformatics* 28: 288–89.
- MacLean, B., D. M. Tomazela, N. Shulman, M. Chambers, G. Finney, B. Frewen, R. Kern, D. L. Tabb, D. C. Liebler, and M. J. MacCoss. 2010. "Skyline: An open source document editor for creating and analyzing targeted proteomics experiments." *Bioinformatics* 26–27: 966.
- Mueller, L. N., O. Rinner, A. Schmidt, S. Letarte, B. Bodenmiller, M.-Y. Brusniak, O. Vitek, R. Aebersold, and M. Müller. 2007. "SuperHirn - a novel tool for high resolution LC-MS-based peptide/protein profiling."

Proteomics 7: 3470–80.

Röst, H. L., G. Rosenberger, P. Navarro, L. Gillet, S. M. Miladinović, O. T. Schubert, W. Wolski, et al. 2014. “OpenSWATH enables automated, targeted analysis of data-independent acquisition MS data.” *Nat. Biotechnol.* 32 (3): 219–23.

Sturm, Marc, Andreas Bertsch, Clemens Gröpl, Andreas Hildebrandt, Rene Hussong, Eva Lange, Nico Pfeifer, et al. 2008. “OpenMS – An open-source software framework for mass spectrometry.” *BMC Bioinformatics* 9 (1): 163.