

Package ‘msImpute’

August 1, 2020

Type Package

Title Peptide imputation in label-free proteomics

Version 1.2.0

Description msImpute provides tools for matrix completion in label-free proteomics quantification at the peptide-level. Currently, msImpute completes missing values by low-rank approximation of the underlying data matrix.

Imports softImpute

License MIT

Encoding UTF-8

LazyData true

BugReports <https://github.com/DavisLaboratory/msImpute/issues>

RoxygenNote 7.0.2

R topics documented:

| | |
|------------------------------------|----------|
| computeStructuralMetrics | 1 |
| CPD | 3 |
| KNC | 3 |
| KNN | 4 |
| msImpute | 4 |
| scaleData | 6 |
| selectFeatures | 7 |
| Index | 9 |

computeStructuralMetrics

Metrics for the assessment of post-imputation structural preservation

Description

For an imputed dataset, it computes within phenotype/experimental condition similarity (i.e. preservation of local structures), between phenotype distances (preservation of global structures), and the Gromov-Wasserstein (GW) distance between original and imputed data.

Usage

```
computeStructuralMetrics(x, group, xna = NULL)
```

Arguments

| | |
|--------------------|---|
| <code>x</code> | numeric matrix. An imputed data matrix. |
| <code>group</code> | factor. A vector of biological groups, experimental conditions or phenotypes (e.g. control, treatment). |
| <code>xna</code> | numeric matrix. Data matrix with missing values (i.e. the original intensity matrix with NAs) |

Details

For each group of experimental conditions (e.g. treatment and control), the group centroid is calculated as the average of observed peptide intensities. Withinness for each group is computed as sum of the squared distances between samples in that group and the group centroid. Betweenness is computed as sum of the squared distances between group centroids. When comparing imputation approaches, the optimal imputation strategy should minimize the within group distances, hence smaller withinness, and maximizes between group distances, hence larger betweenness. The GW metric considers preservation of both local and global structures simultaneously. A small GW distance suggests that imputation has introduced small distortions to global and local structures overall, whereas a large distance implies significant distortions. When comparing two or more imputation methods, the optimal method is the method with smallest GW distance. To compute the GW distance, the missing values in each column of `xna` are replaced by mean of observed values in that column. This is equivalent to imputation by KNN, where `k` is set to the total number of identified peptides (i.e. number of rows in the input matrix). GW distance estimation requires `python`. See example. All metrics are on log scale.

Value

list of three metrics: withinness (sum of squared distances within a phenotype group), betweenness (sum of squared distances between the phenotypes), and gromov-wasserstein distance (if `xna` is not NULL). All metrics are on log scale.

Examples

```
# To compute the GW distance you need to have python installed
# then install the reticulate R package from CRAN
# install.packages("reticulate")
library(reticulate)
# create a virtual environment
virtualenv_create('r-reticulate')
py_available() # if this returns TRUE, you've access to python from R. See reticulate if
# install scipy and POT python modules in this virtual environment
virtualenv_install("msImpute-reticulate", "scipy")
# if this runs successfully, the installation has been successful:
scipy <- import("scipy")
# You can now run the computeStructuralMetrics() function to compute GW distance.
# This setup should only be done for the first use. For all subsequent usages
# load the virtual environment that you've created using:
library(reticulate)
use_virtualenv("msImpute-reticulate")
# you can then run the computeStructuralMetrics() function. Note that the reticulate pack
set.seed(101)
```

```

n=200
p=100
J=50
np=n*p
missfrac=0.3
x=matrix(rnorm(n*J),n,J)%*%matrix(rnorm(J*p),J,p)+matrix(rnorm(np),n,p)/5
ix=seq(np)
imiss=sample(ix,np*missfrac,replace=FALSE)
xna=x
xna[imiss]=NA
y <- xna
xna <- scaleData(xna)
xcomplete <- msImpute(object=xna)
G <- as.factor(sample(1:5, 100, replace = TRUE))
computeStructuralMetrics(xcomplete, G, y)

```

CPD

*CPD***Description**

Spearman correlation between pairwise distances in the original data and imputed data. CPD quantifies preservation of the global structure after imputation. Requires complete datasets - for developers/use in benchmark studies only.

Usage

```
CPD(xorigin, ximputed)
```

Arguments

| | |
|-----------------------|--|
| <code>xorigin</code> | numeric matrix. The original data. Can not contain missing values. |
| <code>ximputed</code> | numeric matrix. The imputed data. Can not contain missing values. |

Value

numeric

KNC

*k-nearest class means (KNC)***Description**

The fraction of k-nearest class means in the original data that are preserved as k-nearest class means in imputed data. KNC quantifies preservation of the mesoscopic structure after imputation. Requires complete datasets - for developers/use in benchmark studies only.

Usage

```
KNC(xorigin, ximputed, class, k = 3)
```

Arguments

| | |
|----------|---|
| xorigin | numeric matrix. The original data. Can contain missing values. |
| ximputed | numeric matrix. The imputed data. |
| class | factor. A vector of length number of columns (samples) in the data specifying the class/label (i.e. experimental group) of each sample. |
| k | number of nearest class means. default to k=3. |

Value

numeric The proportion of preserved k-nearest class means in imputed data.

| | |
|-----|----------------------------------|
| KNN | <i>k-nearest neighbour (KNN)</i> |
|-----|----------------------------------|

Description

The fraction of k-nearest neighbours in the original data that are preserved as k-nearest neighbours in imputed data. KNN quantifies preservation of the local, or microscopic structure. Requires complete datasets - for developers/use in benchmark studies only.

Usage

```
KNN(xorigin, ximputed, k = 3)
```

Arguments

| | |
|----------|--|
| xorigin | numeric matrix. The original data. Can not contain missing values. |
| ximputed | numeric matrix. The imputed data. Can not contain missing values. |
| k | number of nearest neighbours. default to k=3. |

Value

numeric The proportion of preserved k-nearest neighbours in imputed data.

| | |
|----------|--|
| msImpute | <i>Peptide-level imputation in mass spectrometry label-free proteomics by low-rank approximation</i> |
|----------|--|

Description

Returns a completed peptide intensity matrix where missing values (NAs) are imputed by low-rank approximation of the input matrix. Non-NA entries remain unmodified. msImpute requires at least 4 non-missing measurements per peptide across all samples. It is assumed that peptide intensities (DDA), or MS1/MS2 normalised peak areas (DIA), are log2-transformed and normalised (e.g. quantile normalisation).

Usage

```
msImpute(object, rank.max = NULL, lambda = NULL, thresh = 1e-05,
          maxit = 100, trace.it = FALSE, warm.start = NULL,
          final.svd = TRUE)
```

Arguments

| | |
|-------------------------|--|
| <code>object</code> | Numeric matrix where missing values are denoted by NA. Rows are peptides, columns are samples. |
| <code>rank.max</code> | Numeric. This restricts the rank of the solution. is set to <code>min(dim(object)-1)</code> by default. |
| <code>lambda</code> | Numeric. Nuclear-norm regularization parameter. Controls the low-rank property of the solution to the matrix completion problem. By default, it is determined at the scaling step. If set to zero the algorithm reverts to "hardImputation", where the convergence will be slower. |
| <code>thresh</code> | Numeric. Convergence threshold. Set to 1e-05, by default. |
| <code>maxit</code> | Numeric. Maximum number of iterations of the algorithm before the algorithm is converged. 100 by default. |
| <code>trace.it</code> | Logical. Prints traces of progress of the algorithm. |
| <code>warm.start</code> | List. A SVD object can be used to initialize the algorithm instead of random initialization. |
| <code>final.svd</code> | Logical. Shall final SVD object be saved? The solutions to the matrix completion problems are computed from U, D and V components of final SVD. |

Details

`msImpute` operates on the `softImpute`-ALS algorithm. For more details on the underlying algorithm, please see [softImpute](#) package.

Value

Missing values are imputed by low-rank approximation of the input matrix. If input is a numeric matrix, a numeric matrix of identical dimensions is returned. If `x` is a `MAList` object, the E component is replaced with the completed matrix, and the updated `MAList` object is returned. Non-NA entries remain unmodified.

See Also

`selectFeatures`, `scaleData`

Examples

```
set.seed(101)
n=200
p=100
J=50
np=n*p
missfrac=0.3
x=matrix(rnorm(n*J),n,J)%*%matrix(rnorm(J*p),J,p)+matrix(rnorm(np),n,p)/5
ix=seq(np)
imiss=sample(ix,np*missfrac,replace=FALSE)
xna=x
```

```
xna[imiss]=NA
xna <- scaleData(xna)
xcomplete <- msImpute(object=xna)
```

| | |
|-----------|--|
| scaleData | <i>Standardize a matrix to have optionally row means zero and variances one, and/or column means zero and variances one.</i> |
|-----------|--|

Description

Standardize a matrix to have optionally row means zero and variances one, and/or column means zero and variances one.

Usage

```
scaleData(object, maxit = 20, thresh = 1e-09, row.center = TRUE,
  row.scale = TRUE, col.center = TRUE, col.scale = TRUE,
  trace = FALSE)
```

Arguments

| | |
|------------|--|
| object | numeric matrix where missing values are denoted by NA. Rows are peptides, columns are samples. |
| maxit | numeric. maximum iteration for the algorithm to converge (default to 20). When both row and column centering/scaling is requested, iteration may be necessary. |
| thresh | numeric. Convergence threshold (default to 1e-09). |
| row.center | logical. if row.center==TRUE (the default), row centering will be performed resulting in a matrix with row means zero. If row.center is a vector, it will be used to center the rows. If row.center=FALSE nothing is done. |
| row.scale | if row.scale==TRUE, the rows are scaled (after possibly centering, to have variance one. Alternatively, if a positive vector is supplied, it is used for row centering. |
| col.center | Similar to row.center |
| col.scale | Similar to row.scale |
| trace | logical. With trace=TRUE, convergence progress is reported, when iteration is needed. |

Details

Standardizes rows and/or columns of a matrix with missing values, according to the biScale algorithm in Hastie et al. 2015.

Value

A list of two components: E and E.scaled. E contains the input matrix, E.scaled contains the scaled data

See Also

selectFeatures, msImpute

Examples

```

set.seed(101)
n=200
p=100
J=50
np=n*p
missfrac=0.3
x=matrix(rnorm(n*J),n,J)%*%matrix(rnorm(J*p),J,p)+matrix(rnorm(np),n,p)/5
ix=seq(np)
imiss=sample(ix,np*missfrac,replace=FALSE)
xna=x
xna[imiss]=NA
xna <- scaleData(xna)

```

selectFeatures

*Select features with high biological dropout rate***Description**

Fits a linear model to peptide dropout rate against peptide abundance. The selected features (peptides) can be used to determine if data is Missing Not At Random (MNAR). Users should note that msImpute assumes peptides are Missing At Random (MAR).

Usage

```
selectFeatures(object, n_features = 500, suppress_plot = FALSE)
```

Arguments

| | |
|---------------|--|
| object | Numeric matrix where missing values are denoted by NA. Rows are peptides, columns are samples. |
| n_features | Numeric, number of features with high dropout rate. 500 by default. |
| suppress_plot | Logical show plot of dropouts vs abundances. |

Value

A data frame with a logical column denoting the selected features

See Also

scaleData, msImpute

Examples

```

set.seed(101)
n=800
p=100
J=50
np=n*p
missfrac=0.3
x=matrix(rnorm(n*J),n,J)%*%matrix(rnorm(J*p),J,p)+matrix(rnorm(np),n,p)/5

```

```

ix=seq(np)
imiss=sample(ix,np*missfrac,replace=FALSE)
xna=x
xna[imiss]=NA
rownames(xna) <- 1:nrow(xna)
hdp <- selectFeatures(xna, n_features=500, suppress_plot=FALSE)
# construct matrix M to capture missing entries
M <- ifelse(is.na(xna),1,0)
M <- M[hdp$msImpute_feature,]
# plot a heatmap of missingness patterns for the selected peptides
library(ComplexHeatmap)
hm <- Heatmap(M,
column_title = "dropout pattern, columns ordered by dropout similarity",
  name = "",
  col = c("#8FBC8F", "#FFEFDB"),
  show_row_names = FALSE,
  show_column_names = TRUE,
  cluster_rows = TRUE,
  cluster_columns = TRUE,
  show_column_dend = FALSE,
  show_row_dend = FALSE,
  row_names_gp = gpar(fontsize = 7),
  column_names_gp = gpar(fontsize = 8),
  heatmap_legend_param = list(#direction = "horizontal",
  heatmap_legend_side = "bottom",
  labels = c("observed","missing"),
  legend_width = unit(6, "cm")),
)
hm <- draw(hm, heatmap_legend_side = "left")

```


Index

`computeStructuralMetrics`, [1](#)
`CPD`, [3](#)

`KNC`, [3](#)
`KNN`, [4](#)

`msImpute`, [4](#)

`scaleData`, [6](#)
`selectFeatures`, [7](#)
`softImpute`, [5](#)