

Package ‘SUMO’

April 13, 2025

Title Generating Multi-Omics Datasets for Testing and Benchmarking

Version 0.2.0

Description

Provides tools to simulate multi-omics datasets with predefined signal structures. The generated data can be used for testing, validating, and benchmarking integrative analysis methods such as factor models and clustering approaches. This version includes enhanced signal customization, visualization tools (scatter, histogram, 3D), MOFA-based analysis pipelines, PowerPoint export, and statistical profiling of datasets. Designed for both method development and teaching, SUMO supports real and synthetic data pipelines with interpretable outputs. Tini, Giulia, et al (2019) <[doi:10.1093/bib/bbx167](https://doi.org/10.1093/bib/bbx167)>.

License CC BY 4.0

Encoding UTF-8

RoxigenNote 7.3.1

Suggests testthat (>= 3.0.0), MOFADATA, MOFA2, rvg, fabia, tidyverse, grid, basilisk

Config/testthat.edition 3

Imports ggplot2, gridExtra, rlang, stats, graphics, utils, dplyr, readr, readxl, stringr, data.table, magrittr, officer

Collate 'divide_vector.R' 'divide_features_two.R'
'divide_features_one.R' 'feature_selection_two.R'
'feature_selection_one.R' 'divide_samples.R' 'OmixCraftHD.R'
'SUMO.R' 'compute_means_vars.R' 'demo_multomics_analysis.R'
'globals.R' 'plot_factor.R' 'plot_simData.R' 'plot_weights.R'

NeedsCompilation no

Author Bernard Isekah Osang'ir [aut, cre] (<<https://orcid.org/0000-0002-5557-3602>>),
Ziv Shkedy [ctb],
Surya Gupta [ctb],
Jürgen Claesen [ctb]

Maintainer Bernard Isekah Osang'ir <Bernard.Osangir@sckcen.be>

Repository CRAN

Date/Publication 2025-04-13 16:00:19 UTC

Contents

| | |
|-----------------------------------|----|
| compute_means_vars | 2 |
| demo_multomics_analysis | 3 |
| divide_features_one | 4 |
| divide_features_two | 5 |
| divide_samples | 5 |
| divide_vector | 6 |
| feature_selection_one | 6 |
| feature_selection_two | 7 |
| OmixCraftHD | 7 |
| plot_factor | 8 |
| plot_simData | 9 |
| plot_weights | 10 |
| SUMO | 11 |

| | |
|--------------|-----------|
| Index | 13 |
|--------------|-----------|

compute_means_vars *Compute Summary Statistics for a List of Datasets*

Description

Computes overall, row-wise, and column-wise means and standard deviations for each dataset in a list. Also provides average statistics across datasets.

Usage

```
compute_means_vars(data_list)
```

Arguments

| | |
|------------------|---|
| data_list | A list of numeric matrices or data frames. Each entry should be a matrix or data frame with numeric values. |
|------------------|---|

Value

A named list containing:

- Overall mean and SD for each dataset.
- Average row-wise mean and SD.
- Average column-wise mean and SD.
- **mean_smp**: Average row-wise mean across all datasets.
- **sd_smp**: Average row-wise SD across all datasets.

Examples

```
# Example using simulated matrices
set.seed(123)
dataset1 <- matrix(rnorm(100, mean = 5, sd = 2), nrow = 10, ncol = 10)
dataset2 <- matrix(rnorm(100, mean = 10, sd = 3), nrow = 10, ncol = 10)
data_list <- list(dataset1, dataset2)
results <- compute_means_vars(data_list)
print(results)

## Not run:
# Example using real experimental data (requires MOFADATA)
if (requireNamespace("MOFADATA", quietly = TRUE)) {
  utils::data("CLL_data", package = "MOFADATA")
  CLL_data2 <- CLL_data[c(2, 3)]
  results <- compute_means_vars(CLL_data2)
  print(results)
}

## End(Not run)
```

demo_multomics_analysis

Demonstration of SUMO Utility in Multi-Omics Analysis using MOFA2

Description

Run a complete MOFA2-based analysis pipeline using either SUMO-generated or real-world CLL multi-omics data. This function includes preprocessing, MOFA model training, variance decomposition visualization, and optional PowerPoint report generation.

Usage

```
demo_multomics_analysis(
  data_type = c("SUMO", "real_world"),
  export_pptx = TRUE,
  verbose = TRUE
)
```

Arguments

| | |
|-------------|---|
| data_type | Character. Options are "SUMO" for synthetic data or "real_world" for the CLL dataset. |
| export_pptx | Logical. If TRUE, saves a PowerPoint summary of the analysis. Default is TRUE. |
| verbose | Logical. If TRUE, prints progress messages. Default is TRUE. |

Details

The function checks for required packages such as MOFA2, MOFADATA, officer, basilisk, and others using `requireNamespace()`. All downstream analysis is conditionally executed based on package availability.

Value

Invisibly returns the trained MOFA model object. Optionally saves visualizations and a PowerPoint report to disk.

See Also

[OmixCraftHD\(\)](#), [plot_factor\(\)](#), [plot_weights\(\)](#)

Examples

```
if (
  requireNamespace("MOFA2", quietly = TRUE) &&
  requireNamespace("MOFADATA", quietly = TRUE) &&
  identical(Sys.getenv("NOT_CRAN"), "true")
) {
  demo_multiomics_analysis("SUMO", export_pptx = FALSE)
  demo_multiomics_analysis("real_world", export_pptx = FALSE)
}
```

| | |
|----------------------------------|--|
| <code>divide_features_one</code> | <i>Dividing features to create vectors with signal in the first omic for single data</i> |
|----------------------------------|--|

Description

Dividing features to create vectors with signal in the first omic for single data

Usage

```
divide_features_one(n_features_one, num.factor)
```

Arguments

`n_features_one` number of features of first omic

`num.factor` number of factor = '1'

| | |
|---------------------|---|
| divide_features_two | <i>Dividing features to create vectors with signal in the second omic for single data</i> |
|---------------------|---|

Description

Dividing features to create vectors with signal in the second omic for single data

Usage

```
divide_features_two(n_features_two, num.factor)
```

Arguments

n_features_two number of features of second omic
num.factor type of factors - single or multiple

| | |
|----------------|------------------------|
| divide_samples | <i>Global Variable</i> |
|----------------|------------------------|

Description

A global variable used in multiple functions.

Usage

```
divide_samples(n_samples, num, min_size)
```

Arguments

n_samples number of samples
num number of factors
min_size Minimum length of any samples scores

| | |
|---------------|---|
| divide_vector | <i>#' Global Variable #' #' A global variable used in multiple functions. #' #'</i> |
|---------------|---|

Description

#' Global Variable #' #' A global variable used in multiple functions. #' #'

Usage

```
divide_vector(n_samples, num, min_size)
```

Arguments

| | |
|-----------|---|
| n_samples | number of samples |
| num | number of factors |
| min_size | Minimum length of any samples scores <i>#' ## ~~~~~ Updated IN USE (IN USE): Simulate the samples scores (IN USE) ~~~~~ ~~~~~ ~~~~~</i> ~~~~~ |

| | |
|-----------------------|--|
| feature_selection_one | <i>Dividing features to create vectors with signal in the first omic</i> |
|-----------------------|--|

Description

Dividing features to create vectors with signal in the first omic

Usage

```
feature_selection_one(n_features_one, num.factor, no_factor)
```

Arguments

| | |
|----------------|--------------------------------------|
| n_features_one | number of features of first omic |
| num.factor | type of factors - single or multiple |
| no_factor | number of factors |

feature_selection_two *Dividing features to create vectors with signal in the second omic*

Description

Dividing features to create vectors with signal in the second omic

Usage

```
feature_selection_two(n_features_two, num.factor, no_factor)
```

Arguments

| | |
|----------------|--------------------------------------|
| n_features_two | number of features of second omic |
| num.factor | type of factors - single or multiple |
| no_factor | number of factors |

OmixCraftHD

Simulation of omics with predefined single or multiple latent factors in multi-omics

Description

Simulates two high-dimensional omics datasets with customizable latent factor structures. Users can control the number and type of factors (shared, unique, mixed), the signal-to-noise ratio, and the distribution of signal-carrying samples and features. The function is flexible for benchmarking multi-omics integration methods under various controlled scenarios.

Usage

```
OmixCraftHD(
  vector_features = c(2000, 2000),
  n_samples = 50,
  n_factors = 3,
  signal.samples = NULL,
  signal.features.one = NULL,
  signal.features.two = NULL,
  num.factor = "multiple",
  snr = 1,
  advanced_dist = NULL,
  ...
)
```

Arguments

| | |
|----------------------------------|---|
| <code>vector_features</code> | A numeric vector of length two, specifying the number of features in the first and second omics datasets, respectively. |
| <code>n_samples</code> | Integer. The number of samples shared between both omics datasets. |
| <code>n_factors</code> | Integer. Number of latent factors to simulate. |
| <code>signal.samples</code> | Optional numeric vector of length two: the first element is the mean, and the second is the variance of the number of signal-carrying samples per factor. If <code>NULL</code> , signal assignment is inferred from <code>snr</code> . |
| <code>signal.features.one</code> | Optional numeric vector of length two: the first element is the mean, and the second is the variance of the number of signal-carrying features per factor in the first omic. |
| <code>signal.features.two</code> | Optional numeric vector of length two: the first element is the mean, and the second is the variance of the number of signal-carrying features per factor in the second omic. |
| <code>num.factor</code> | Character string. Either ' <code>'single'</code> ' or ' <code>'multiple'</code> '. Determines whether to simulate a single latent factor or multiple factors. |
| <code>snr</code> | Numeric. Signal-to-noise ratio used to estimate the background noise. The function uses this value to infer the proportion of signal versus noise in the simulated datasets. |
| <code>advanced_dist</code> | Character string. Specifies how latent factors are distributed when <code>num.factor = 'multiple'</code> . Options include: <code>"</code> , <code>NULL</code> , <code>'mixed'</code> , <code>'omic.one'</code> , <code>'omic.two'</code> , or <code>'exclusive'</code> . |
| <code>...</code> | Additional arguments (not currently used). |

`plot_factor`

Visualization of factor scores (ground truth)

Description

Scatter or histogram plots of sample-level factor scores from simulated multi-omics data, using scores from `list_alphas` and `list_gammas`.

Usage

```
plot_factor(
  sim_object = NULL,
  factor_num = NULL,
  type = "scatter",
  show.legend = TRUE
)
```

Arguments

| | |
|-------------|---|
| sim_object | R object containing simulated data output from OmixCraftHD |
| factor_num | Integer or "all". Which factor(s) to plot. |
| type | Character. Either "scatter" (default) or "histogram" for plot type. |
| show.legend | Logical. Whether to show legend in plots. Default is TRUE. |

Examples

```
output_obj <- OmixCraftHD(  
  vector_features = c(4000, 3000),  
  n_samples = 100,  
  n_factors = 2,  
  snr = 2.5,  
  num.factor = 'multiple',  
  advanced_dist = 'mixed')  
  
plot_factor(sim_object = output_obj, factor_num = 1)  
plot_factor(sim_object = output_obj, factor_num = 'all', type = 'histogram')
```

plot_simData

Visualizing the simulated data using heatmap or 3D surface plot

Description

Generates a visual representation of the simulated omics data either as a heatmap or a 3D surface plot. You can select which dataset to visualize: the merged/concatenated matrix, omic one, or omic two.

Usage

```
plot_simData(sim_object, data = "merged", type = "heatmap")
```

Arguments

| | |
|------------|---|
| sim_object | R object containing simulated data as returned by OmixCraftHD. |
| data | Character. Specifies which data matrix to visualize. Options are "merged" (or "concatenated"), "omic.one", or "omic.two". |
| type | Character. Type of plot: either "heatmap" for a 2D image plot or "3D" for a 3D perspective surface plot. |

Examples

```
output_obj <- OmixCraftHD(
  vector_features = c(4000,3000),
  n_samples=100,
  n_factors=2,
  signal.samples = NULL,
  signal.features.one = NULL,
  signal.features.two = NULL,
  snr = 2.5,
  num.factor='multiple',
  advanced_dist='mixed')

plot_simData(sim_object = output_obj, data = "merged", type = "heatmap")
plot_simData(sim_object = output_obj, data = "omic.one", type = "3D")
```

plot_weights

Visualizing the raw loading/weights of the features

Description

Generates scatter or histogram plots of feature loadings (weights) from simulated multi-omics data. Supports plotting for omic.one, omic.two, or integrated views.

Usage

```
plot_weights(
  sim_object,
  factor_num = 1,
  data = "omic.one",
  type = "scatter",
  show.legend = TRUE
)
```

Arguments

| | |
|--------------------------|---|
| <code>sim_object</code> | R object containing data to be plotted. |
| <code>factor_num</code> | Integer or "all". Specifies which factor(s) to visualize. |
| <code>data</code> | Character. Section of the data to visualize: "omic.one", "omic.two", or "integrated". |
| <code>type</code> | Character. Type of plot: "scatter" or "histogram". |
| <code>show.legend</code> | Logical. Whether to show the legend in the plot. Default is TRUE. |

Value

A ggplot object or a combined grid of plots.

Examples

```

output_obj <- OmixCraftHD(
  vector_features = c(4000, 3000),
  n_samples = 100,
  n_factors = 2,
  signal.samples = NULL,
  signal.features.one = NULL,
  signal.features.two = NULL,
  snr = 2.5,
  num.factor = 'multiple',
  advanced_dist = 'mixed'
)

plot_weights(
  sim_object = output_obj,
  factor_num = 1,
  data = 'omic.one',
  type = 'scatter',
  show.legend = FALSE
)

plot_weights(
  sim_object = output_obj,
  factor_num = 2,
  data = 'omic.two',
  type = 'histogram'
)

```

Description

It provides tools for simulating complex multi-omics datasets, enabling researchers to generate data that mirrors the biological intricacies observed in real-world omics studies. This package addresses a critical gap in current bioinformatics by offering flexible and customizable methods for synthetic multi-omics data generation, supporting method development, validation, and benchmarking.

Details

Key Features:

- **Multi-Omics Simulation:** Generate multi-layered datasets with shared and modality-specific structures.
- **Flexible Generation Engine:** Fine control over samples, noise levels, signal distributions, and latent factor structures.
- **Pipeline-Friendly Design:** Seamlessly integrates with existing multi-omics analysis workflows and packages (e.g., `SummarizedExperiment`, `MultiAssayExperiment`).

- **Visualization Tools:** Built-in plotting functions for exploring synthetic signals, factor structures, and noise.

Main Functions:

- `OmixCraftHD()`: Simulates synthetic high-dimensional multi-omics datasets.
- `plot_simData()`: Visualizes generated data at different levels.
- `plot_factor()`: Displays factor scores across samples for signal inspection.
- `plot_weights()`: Visualizes feature loadings to assess signal versus noise.
- `demo_multiomics_analysis()`: Full demo function for applying MOFA on SUMO-generated or real-world data.
- `compute_means_vars()`: Estimate parameters from the real experimental dataset.

Author(s)

Maintainer: Bernard Isekah Osang'ir <Bernard.Osangir@sckcen.be> ([ORCID](#))

Other contributors:

- Ziv Shkedy [contributor]
- Surya Gupta [contributor]
- Jürgen Claesen [contributor]

Index

- * **MOFA**
 - demo_multiomics_analysis, 3
- * **benchmarking**
 - SUMO, 11
- * **demo**
 - demo_multiomics_analysis, 3
- * **models**
 - SUMO, 11
- * **multi-omics**
 - demo_multiomics_analysis, 3
 - SUMO, 11
- * **synthetic-data**
 - demo_multiomics_analysis, 3
- compute_means_vars, 2
- demo_multiomics_analysis, 3
- divide_features_one, 4
- divide_features_two, 5
- divide_samples, 5
- divide_vector, 6
- feature_selection_one, 6
- feature_selection_two, 7
- OmixCraftHD, 7
- OmixCraftHD(), 4
- plot_factor, 8
- plot_factor(), 4
- plot_simData, 9
- plot_weights, 10
- plot_weights(), 4
- SUMO, 11
- SUMO-package (SUMO), 11