
Normalisr

Lingfei Wang

Oct 31, 2020

CONTENTS:

1	Normalisr	1
1.1	Installation	1
1.2	Usage	1
1.3	Documentation	2
1.4	Examples and pipelines	2
1.5	Contact	2
1.6	References	2
1.7	FAQ	3
2	User API	5
2.1	Quality control	5
2.2	Normalization	5
2.3	Differential expression	5
2.4	Co-expression	5
2.5	API list	6
3	All API	15
3.1	normalisr.association	15
3.2	normalisr.binnet	19
3.3	normalisr.coex	20
3.4	normalisr.de	21
3.5	normalisr.gocovt	22
3.6	normalisr.lcpm	24
3.7	normalisr.norm	26
3.8	normalisr.parallel	28
3.9	normalisr.qc	28
4	Indices and tables	31
	Python Module Index	33
	Index	35

NORMALISR

pypi v0.4.2

DOI 10.5281/zenodo.3986428

Normalizr is a parameter-free normalization-association two-step inferential framework for scRNA-seq that solves case-control differential expression, co-expression, and pooled CRISPRi scRNA-seq screen under one umbrella of linear association testing. Normalizr addresses sparsity and technical confounding challenges of scRNA-seq with posterior mRNA abundances, nonlinear cellular summary covariates, and mean and variance normalization. All these enable linear association testing to achieve optimal sensitivity, specificity, and speed in all above scenarios.

Normalizr follows the conventional framework of normalization/imputation of scRNA-seq, and aims to recover the true, biological, but hidden expression levels which any analyses may then operate upon. Then, linear association testing provides a unified inferential framework with numerous advantages: (i) exact P-value estimation without permutation, (ii) native removal of covariates as fixed effects, (iii) non-parametric robustness, (iv) unbeatable time and memory complexities, and (v) extension potentials such as variations in genetic relatedness.

Normalizr is written in Python3 and provides a command-line and a python functional interface. You can read more about Normalizr from our preprint (See [References](#)).

1.1 Installation

Normalizr is on [PyPI](#) and can be installed with pip: `pip install normalizr`. You can also install Normalizr from github: `pip install git+https://github.com/lingfeiwang/normalizr.git`. Make sure you have added Normalizr's install path into PATH environment before using the command-line interface (See [FAQ](#)). Normalizr's installation should take less than a minute.

There are more advanced installation methods but if you want that, most likely you already know how to do it ;). If not, give me a shout (See [Contact](#)).

1.2 Usage

Normalizr provides a command-line and a python functional interface below. You can use the examples provided below to guide yourself through Normalizr's use. Sphinx-based documentation is underway.

- **Command-line interface** You can run Normalizr by typing `normalizr` on command-line. Normalizr uses submodules for different analysis steps. Type `normalizr` or `normalizr -h` for general help, and for example `normalizr de -h` for help on submodule 'de' of differential expression.

Normalizr uses tsv (tab separated values) file format for input and output matrices, and text file for row and column names, such as cells and genes, one per line. For initial input, Normalizr also accepts the sparse

mtx format (Cell Ranger output) for raw read count matrix. Gzipped input/output files are automatically recognized if file name suffix '.gz' is present.

- **Python functional interface** Normalisr's python functional interface is more flexible than command-line, but requires knowledge of python programming. Documentation of any function can be obtained with ? in ipython or jupyter notebook, such as:

```
import normalisr.normalisr as norm
?norm.de
```

The example jupyter notebooks also illustrate the scope of functions Normalisr provides.

1.3 Documentation

Documentations are available as [html](#) and [pdf](#).

1.4 Examples and pipelines

You can find several examples in the 'examples' folder, to cover all functions Normalisr currently provides. The example datasets have been scaled down to run on a 16GB-memory personal computer. Although they only serve as demonstrations of work here, the pipelines should be transferable to a full-scale, different dataset. Since Normalisr is non-parametric, the only adjustable parameters are for quality control and final cutoffs of differential or co-expression. You can change down-sampling parameters in the examples to run the full datasets on a larger computer.

You can find more details in the respective examples.

1.5 Contact

We look forward to your feedbacks or questions of any kind.

- Regarding method and manuscript, please reach me by email (Lingfei.Wang.CN@gmail.com).
- Regarding Normalisr package, please raise an issue on [github](#) or reach me by email (Lingfei.Wang.github@outlook.com).

1.6 References

Please find the currently available materials below:

- [Abstract](#) submitted to ICML 2020 Workshop on Computational Biology

[Contact](#) me if you would like to request more details.

1.7 FAQ

- **What does Normalisr stand for?** Normalisr Offers Robust Modelling of Associations Linearly In Single-cell RNA-seq. Yes, it's a recursive acronym. See [GNU](#) and [pip](#).
- **I installed Normalisr but typing `normalisr` says 'command not found'.**
- **How do I use a specific python version for Normalisr's command-line interface?** You can always use the python command to run Normalisr, such as `python3 -m normalisr` to replace command `normalisr`. You can also use a specific path or version for python, such as `python3.7 -m normalisr` or `/usr/bin/python3.7 -m normalisr`. Make sure you have installed Normalisr for this python version.
- **Why don't the examples work?** Please make sure you followed every step in the README.md of the respective example folder with Internet connection, and then submit an issue report detailing at which executed line the error occurred with input and output.
- **Does Normalisr run on Windows?** I have not tested Normalisr on Windows. However, it is purely in python and should be able to function properly.

2.1 Quality control

<i>qc_reads</i>	Quality control by bounding read counts.
<i>qc_outlier</i>	Quality control by removing cell outliers by variance.

2.2 Normalization

<i>lcpm</i>	Computes Bayesian log CPM from raw read counts.
<i>normcov</i>	Normalizes each continuous covariate to 0 mean and unit variance.
<i>scaling_factor</i>	Computes scaling factor of variance normalization for every gene.
<i>compute_var</i>	Computes variance normalization scale for each cell.
<i>normvar</i>	Performs mean and variance normalizations.

2.3 Differential expression

<i>de</i>	Performs differential expression analyses for all genes against all groupings.
-----------	--

2.4 Co-expression

<i>coex</i>	Performs co-expression analyses for all gene pairs.
<i>binnet</i>	Binarizes P-value co-expression network to thresholded Q-value network.
<i>gotop</i>	Finds the top variable GO enrichment of top principal genes in the binary co-expression network.
<i>pccovt</i>	Introduces an extra covariate from the top principal component of given genes.

2.5 API list

`normalisr.normalisr.binnet` (*net*, *qcut*)

Binarizes P-value co-expression network to thresholded Q-value network.

Q-values are computed separately per row to account for differences in the number of genes co-expressed, especially by master regulators, using Benjamini–Hochberg procedure. Co-expression Q-value matrix is thresholded for return.

Parameters

- **net** (*numpy.ndarray* (*shape*=(*n_gene*, *n_gene*), *dtype*=*float*)) – Symmetric co-expression P-value matrix.
- **qcut** (*float*) – Cutoff for Q-value network.

Returns Binarized, assymmetric co-expression matrix.

Return type *numpy.ndarray*(*shape*=(*n_gene*, *n_gene*), *dtype*=*bool*)

`normalisr.normalisr.coex` (*dt*, *dc*, ***ka*)

Performs co-expression analyses for all gene pairs.

Performs parallel computation with multiple processes on the same machine.

Model for co-expression between genes i & j: $X_i = \gamma X_j + \alpha C + \epsilon$,

$\epsilon \sim \text{i.i.d. } N(0, \sigma^2)$.

Test statistic: conditional R^2 (or proportion of variance explained) between X_i and X_j .

Null hypothesis: $\gamma = 0$.

Parameters

- **dt** (*numpy.ndarray* (*shape*=(*n_gene*, *n_cell*), *dtype*=*float*)) – Normalized expression matrix X.
- **dc** (*numpy.ndarray* (*shape*=(*n_cov*, *n_cell*), *dtype*=*float*)) – Normalized covariate matrix C.
- **ka** (*dict*) – Keyword arguments passed to `normalisr.association.association_test_1`. See below.

Returns

- **P-values** (*numpy.ndarray*(*shape*=(*n_gene*, *n_gene*))) – Co-expression P-value matrix.
- **dot** (*numpy.ndarray*(*shape*=(*n_gene*, *n_gene*))) – Inner product of expression between gene pairs, after removing covariates.
- **var** (*numpy.ndarray*(*shape*=(*n_gene*))) – Variance of gene expression after removing covariates. Pearson $R = ((\text{dot} / \text{numpy.sqrt}(\text{var})).T) / \text{numpy.sqrt}(\text{var})$.

Keyword Arguments

- **bs** (*int*) – Batch size, i.e. number of genes in each computing batch. Use 0 for default: Data transfer limited to 1GB, capped at `bs=500`.
- **nth** (*int*) – Number of parallel processes. Set to 0 for using automatically detected CPU counts.
- **dimreduce** (*numpy.ndarray* (*shape*=(*n_gene*,), *dtype*=*int*) or *int*) – If *dt* doesn't have full rank, such as due to prior covariate removal (although the recommended

method is to leave covariates in dc), this parameter allows to specify the loss of ranks/degrees of freedom to allow for accurate P-value computation. Default is 0, indicating no rank loss.

`normalisr.normalisr.compute_var(dt, dc, stepmax=1, eps=1e-06)`

Computes variance normalization scale for each cell.

Performs a log-linear fit of the variance of each cell with covariates. Optionally use Expectation-Maximization(EM)-like method to iteratively fit mean and variance. For EM-like method, early-stopping is suggested because of overfitting issues.

Parameters

- **dt** (`numpy.ndarray(shape=(n_gene, n_cell))`) – Bayesian LogCPM expression level matrix.
- **dc** (`numpy.ndarray(shape=(n_cov, n_cell))`) – Covariate matrix.
- **stepmax** (`int`) – Maximum number of EM-like iterations of mean and variance normalization. Stop of iteration is also possible when relative accuracy target is reached. Defaults to 1, indicating no iterative normalization.
- **eps** (`float`) – Relative accuracy target for early stopping. Constrains the maximum relative difference of fitted variance across cells compared to the last step. Defaults to 1E-6.

Returns Inverse sqrt of fitted variance for each cell, i.e. the multiplier for variance normalization. For iterative normalization, the optimal step will be returned, defined as having the minimal max relative change across cells.

Return type `numpy.ndarray(shape=(n_cell,))`

`normalisr.normalisr.de(dg, dt, dc, bs=0, nth=1, single=0, **ka)`

Performs differential expression analyses for all genes against all groupings.

Allows multiple options to treat other groupings when testing on one grouping.

Performs parallel computation with multiple processes on the same machine.

Model for differential expression between gene Y and grouping X: $Y = \gamma X + \alpha C + \epsilon$,
 $\epsilon \sim \text{i.i.d. } N(0, \sigma^2)$.

Test statistic: conditional R^2 (or proportion of variance explained) between Y and X.

Null hypothesis: $\gamma = 0$.

Parameters

- **dg** (`numpy.ndarray(shape=(n_group, n_cell))`) – Grouping matrix for a list of X to be tested, e.g. grouping by gene knock-out.
- **dt** (`numpy.ndarray(shape=(n_gene, n_cell), dtype=float)`) – Normalized expression matrix Y.
- **dc** (`numpy.ndarray(shape=(n_cov, n_cell), dtype=float)`) – Normalized covariate matrix C.
- **bs** (`int`) – Batch size, i.e. number of groupings and genes in each computing batch. For `single=0,1`, splits groupings & genes. Defaults to 500. For `single=4`, only allows `bs=0` (default automatic setting).
- **nth** (`int`) – Number of parallel processes. Set to 0 for using automatically detected CPU counts.
- **single** (`int`) – Option to deal with other groupings when testing one groupings v.s. gene expression.

- 0: Ignores other groupings.
- 1: Excludes all cells belonging to any other grouping (`value==1`), assuming `dg=0,1` only.
- 4: Treats other groupings as covariates for mean expression.
- **ka** (*dict*) – Keyword arguments passed to `normalisr.association.association_test_*`. See below.

Returns

- **P-values** (*numpy.ndarray(shape=(n_group,n_gene))*) – Differential expression P-value matrix.
- **gamma** (*numpy.ndarray(shape=(n_group,n_gene))*) – Differential expression log fold change matrix.
- **alpha** (*numpy.ndarray(shape=(n_group,n_gene,n_cov))*) – Maximum likelihood estimators of alpha, separately tested for each grouping.
- **varg** (*numpy.ndarray(shape=(n_group))*) – Variance of grouping after removing covariates.
- **vart** (*numpy.ndarray(shape=(n_group,n_gene))*) – Variance of gene expression after removing covariates. It can depend on the grouping being tested depending on parameter `single`.

Keyword Arguments

- **dimreduce** (*numpy.ndarray(shape=(n_gene,), dtype=int) or int*) – If `dt` doesn't have full rank, such as due to prior covariate removal (although the recommended method is to leave covariates in `dc`), this parameter allows to specify the loss of ranks/degrees of freedom to allow for accurate P-value computation. Default is 0, indicating no rank loss.
- **method** (*str, only for single=4*) – Method to compute eigenvalues in SVD-based matrix inverse (for removal of covariates):
 - `auto`: Uses `scipy` for `n_matrix < mpc` or `mpc == 0` and `sklearn` otherwise. Default.
 - `scipy`: Uses `scipy.linalg.svd`.
 - `scipys`: NOT IMPLEMENTED. Uses `scipy.sparse.linalg.svds`.
 - `sklearn`: Uses `sklearn.decomposition.TruncatedSVD`.
- **mpc** (*int, only for single=4*) – Uses only the top `mpc` singular values as non-zero in SVD-based matrix inverse. Here effectively reduces covariates to their top principal components. This reduction is performed after including other groupings as additional covariates. Defaults to 0 to disable dimension reduction. For very large grouping matrix, use a small value (e.g. 100) to save time at the cost of accuracy.
- **qr** (*int, only for single=4*) – Whether to use QR decomposition method for SVD in matrix inverse. Only effective when `method=sklearn`, or `=auto` and defaults to `sklearn`. Takes the following values:
 - 0: No (default).
 - 1: Yes with default settings.
 - 2+: Yes with `n_iter=qr` for `sklearn.utils.extmath.randomized_svd`.
- **tol** (*float, only for single=4*) – Eigenvalues $< \text{tol} * (\text{maximum eigenvalue})$ are treated as zero in SVD-based matrix inverse. Default is $1\text{E-}8$.

`normalisr.normalisr.gotop(net, namet, go_file, goa_file, n=100, **ka)`

Finds the top variable GO enrichment of top principal genes in the binary co-expression network.

Principal genes are those with most co-expressed genes. They reflect the most variable pathways in the dataset. When the variable pathways are housekeeping related, they may conceal cell-type-specific co-expression patterns from being observed and understood. This function identifies the most variable pathway with gene ontology enrichment study of the top principal genes. Background genes are all genes provided.

Parameters

- **net** (*numpy.ndarray* (*shape*=(*n_gene*, *n_gene*), *dtype*=*bool*)) – Binary co-expression network matrix.
- **namet** (*list of str*) – Gene names matching the rows and columns of net.
- **go_file** (*str*) – File path for GO DAG (downloadable at <http://geneontology.org/docs/download-ontology/>).
- **goa_file** (*str*) – File path of GO annotation. See parameter **conversion** in `normalisr.gocovt.goe`.
- **n** (*int*) – Number of top principal genes to include for GO enrichment. Default is 100, giving good performance in general.
- **ka** (*dict*) – **IMPORTANT**: Keyword arguments passed to `normalisr.gocovt.goe` to determine how to perform GO enrichment study. If you see no gene mapped, check your gene name conversion rule in **conversion** parameter of `normalisr.gocovt.goe`. GO annotation have a specific gene ID system.

Returns

- **principals** (*list of str*) – List of principal genes.
- **goe** (*pandas.DataFrame*) – GO enrichment results.
- **gotop** (*str*) – Top enriched GO ID.
- **genes** (*list of str*) – List of genes in the gotop GO ID.

`normalisr.normalisr.lcpm` (*reads*, *normalize*=*True*, *nth*=0, *ntot*=*None*, *varscale*=0, *seed*=*None*, *lowmem*=*True*)

Computes Bayesian log CPM from raw read counts.

The technical sampling process is modelled as a Binomial distribution. The logCPM given read counts is a Bayesian inference problem and follows (shifted) Beta distribution. We use the expectation of posterior logCPM as the estimated expression levels. Resampling function is also provided to account for variances in the posterior distribution.

Warning: Modifying keyword arguments other than *nth* or *seed* is neither recommended nor supported for function ‘`lcpm`’. Do so at your own risk.

Parameters

- **reads** (*numpy.ndarray* (*shape*=(*n_gene*, *n_cell*), *dtype*='uint')) – Read count matrix.
- **normalize** (*bool*) – Whether to normalize output to logCPM per cell. Default: *True*.
- **nth** (*int*) – Number of threads to use. Defaults to 0 to use all cores automatically detected.
- **ntot** (*int*) – Manually sets value of total number of reads in binomial distribution. Since the posterior distribution stabilizes quickly as *ntot* increases, a large number, e.g. 1E9 is good for general use. Defaults to *None* to disable manual value.
- **varscale** (*float*) – Resamples estimated expression using the posterior Beta distribution. *varscale* sets the scale of variance than its actual value from the posterior distribution. Defaults to 0, to compute expectation with no variance.

- **seed** (*int*) – Initial random seed if set.
- **lowmem** (*bool*) – Low memory mode disable mean and var in Returns and therefore saves memory.

Returns

- **lcpm** (*numpy.ndarray(shape=(n_gene, n_cell))*) – Estimated expression as logCPM from read counts.
- **mean** (*numpy.ndarray(shape=(n_gene, n_cell)) or None*) – Mean/Expectation of lcpm's every entry's posterior distribution. None if lowmem=True.
- **var** (*numpy.ndarray(shape=(n_gene, n_cell)) or None*) – Variance of lcpm's every entry's posterior distribution. None if lowmem=True.
- **cov** (*numpy.ndarray(shape=(3, n_cell))*) – Cellular summary covariates computed from read count matrix that may confound lcpm. Contains:
 - cov[0]: Log total read count per cell
 - cov[1]: Number of 0-read genes per cell
 - cov[2]: cov[0]**2

`normalisr.normalisr.normcov(dc, c=True)`

Normalizes each continuous covariate to 0 mean and unit variance.

Optionally introduces constant 1 covariate as intercept. Categorical covariates should be in binary/one-hot form, and will be left unchanged.

Parameters

- **dc** (*numpy.ndarray(shape=(n_cov, n_cell))*) – Current covariate matrix. Use empty matrix with n_cov=0 if no covariate.
- **c** (*bool*) – Whether to introduce a constant 1 covariate.

Returns Processed covariate matrix.

Return type `numpy.ndarray(shape=(n_cov+1 if c else n_cov, n_cell))`

`normalisr.normalisr.normvar(dt, dc, w, wt, dextra=None, cat=1, nth=1, bs=500, keepvar=True)`

Performs mean and variance normalizations.

Expression levels are normalized at mean and then at variance levels. Effectively each gene x is multiplied by $w**wt[x]$ before removing covariates as $dc*(w**wt[x])$. Continuous covariates are normalized at variance levels. Effectively covariates are transformed to $dc*w$. Therefore, variance normalization for expression are scaled differently for each gene.

Parameters

- **dt** (*numpy.ndarray(shape=(n_gene, n_cell))*) – Bayesian logCPM matrix.
- **dc** (*numpy.ndarray(shape=(n_cov, n_cell))*) – Covariate matrix.
- **w** (*numpy.ndarray(shape=(n_cell,))*) – Computed variance normalization multiplier.
- **wt** (*numpy.ndarray(shape=(n_gene,))*) – Computed scaling factor for each gene.
- **dextra** (*numpy.ndarray(shape=(n_extra, n_cell))*) – Extra data matrix also to be normalized like continuous covariates.
- **cat** (*int*) – Whether to normalize categorical/binary covariates (those with only 0 or 1s). Defaults to 1.

- 0: No
- 1: No except constant-1 covariate (intercept)
- 2: Yes
- **nth** (*int*) – Number of parallel threads.
- **bs** (*int*) – Batch size for each job.
- **keepvar** (*bool*) – Whether to maintain the variance of each gene invariant in mean normalization step. If so, expression variances are scaled back to original after mean normalization and before variance normalization. This function only affects overall variance level and its downstreams (e.g. differential expression log fold change). This function would not affect P-value computation. Default: True.

Returns

- (*dtn, dcn*) or (*dtn, dcn, dextran*) if *dextra* is not None
- **dtn** (*numpy.ndarray(shape=(n_gene, n_cell))*) – Normalized gene expression matrix.
- **dcn** (*numpy.ndarray(shape=(n_cov, n_cell))*) – Normalized covariate matrix.
- **dextran** (*numpy.ndarray(shape=(n_extra, n_cell))*) – Normalized extra data matrix.

`normalisr.normalisr.pccovt(dt, dc, namet, genes, condcov=True)`

Introduces an extra covariate from the top principal component of given genes.

The extra covariate is the top principal component of normalized expressions of the selected genes. Adding a covariate from housekeeping pathway can reveal cell-type-specific activities in co-expression networks.

Parameters

- **dt** (*numpy.ndarray(shape=(n_gene, n_cell))*) – Normalized expression matrix.
- **dc** (*numpy.ndarray(shape=(n_cov, n_cell))*) – Existing normalized covariate matrix.
- **namet** (*list of str*) – List of gene names for rows in dt.
- **genes** (*list of str*) – List of gene names to include in finding top PC of their expression as an extra covariate.
- **condcov** (*bool*) – Whether to condition on existing covariates before computing top PC. Default: True.

Returns New normalized covariate matrix.

Return type `numpy.ndarray(shape=(n_cov+1, n_cell))`

`normalisr.normalisr.qc_outlier(dw, pcut=1e-10, outrate=0.02)`

Quality control by removing cell outliers by variance.

Fit normal distribution on the inverse sqrt variance to detect outliers. This is performed by iterative estimation of normal distribution with non-outliers and then determination of outliers with the normal distribution.

Parameters

- **dw** (*numpy.ndarray(shape=(n_cell,))*) – Fitted inverse sqrt cell variance.
- **pcut** (*float*) – Bonferroni P-value cutoff for asserting outliers in a normal distribution of fitted cell variance. Default: 1E-10.
- **outrate** (*float*) – Upper bound of proportion of outliers on either side of variance distribution. Used for initial outlier assignment and final validity check. Default: 0.02.

Returns Whether each cell passed QC

Return type `numpy.ndarray(shape=(n_cell,),dtype=bool)`

`normalisr.normalisr.qc_reads (reads, n_gene, nc_gene, ncp_gene, n_cell, nt_cell, ntp_cell)`

Quality control by bounding read counts.

Quality control is performed separately on genes based on their cell statistics and on cells based on their gene statistics, iteratively until dataset remains unchanged. A gene or cell is removed if any of the QC criteria is violated at any time in the iteration. All QC parameters can be set to 0 to disable QC filtering for that criterion.

Parameters

- **reads** (`numpy.ndarray((n_gene, n_cell), dtype='uint')`) – Read count matrix.
- **n_gene** (`int`) – Lower bound on total read counts for gene QC.
- **nc_gene** (`int`) – Lower bound on number of expressed cells for gene QC.
- **ncp_gene** (`float`) – Lower bound on proportion of expressed cells for gene QC.
- **n_cell** (`int`) – Lower bound on total read counts for cell QC.
- **nt_cell** (`int`) – Lower bound on number of expressed genes for cell QC.
- **ntp_cell** (`float`) – Lower bound on proportion of expressed genes for cell QC.

Returns

- **genes_select** (`numpy.ndarray(dtype='uint')`) – Array of indices of genes passed QC.
- **cells_select** (`numpy.ndarray(dtype='uint')`) – Array of indices of cells passed QC.

`normalisr.normalisr.scaling_factor (dt, varname='nt0mean', v0=0, v1='max')`

Computes scaling factor of variance normalization for every gene.

Lowly expressed genes need full variance normalization because of technical confounding from sequencing depth. Highly expressed genes do not need variance normalization because they are already accurately measured. The scaling factor operates as an exponential factor on the variance normalization scale for each gene. It should be maximum/minimum for genes with lowest/highest expression.

Warning: Modifying keyword arguments is neither recommended nor supported for function 'scaling_factor'. Do so at your own risk.

Parameters

- **dt** (`numpy.ndarray(shape=(n_gene, n_cell), dtype='uint')`) – Read count matrix.
- **varname** (`str`) – Variable used to compute scaling factor for each gene. Can be:
 - `logtpropmean`: `log(dt.mean(axis=1)/dt.mean(axis=1).sum())`
 - `logtmeanprop`: `log((dt/dt.sum(axis=0)).mean(axis=1))`
 - `nt0mean`: `(dt==0).mean(axis=1)`
 - `lognt0mean`: `log((dt==0).mean(axis=1))`
 - `log1-nt0mean`: `log(1-(dt==0).mean(axis=1))`Defaults to `nt0mean`.
- **v0** (`float`) – Variable values to set scaling factor to 0 (for `v0`) and 1 (for `v1`). Linear assignment is applied for values inbetween. Can be:
 - `max`: `max`

- min: min
- any float: that float
- **v1** (*float*) – Variable values to set scaling factor to 0 (for v0) and 1 (for v1). Linear assignment is applied for values inbetween. Can be:
 - max: max
 - min: min
 - any float: that float

Returns Scaling factor of variance normalization for each gene

Return type `numpy.ndarray(shape=(n_gene,))`

3.1 normalizr.association

`normalizr.association.association_test_1(vx, vy, dx, dy, dc, dci, dcr, dimreduce=0)`

Fast linear association testing in single-cell non-cohort settings with covariates.

Single threaded version to allow for parallel computing wrapper. Mainly used for naive differential expression and co-exp

$y = \gamma x + \alpha C + \epsilon$,

$\epsilon \sim \text{i.i.d. } N(0, \sigma^2)$.

Test statistic: conditional R^2 (or proportion of variance explained) between x and y .

Null hypothesis: $\gamma = 0$.

Parameters

- **vx** (*any*) – Starting indices of dx . Only used for information passing.
- **vy** (*any*) – Starting indices of dy . Only used for information passing.
- **dx** (`numpy.ndarray(shape=(n_x, n_cell))`) – Predictor matrix for a list of vector x to be tested, e.g. gene expression or grouping.
- **dy** (`numpy.ndarray(shape=(n_y, n_cell))`) – Target matrix for a list of vector y to be tested, e.g. gene expression.
- **dc** (`numpy.ndarray(shape=(n_cov, n_cell))`) – Covariate matrix as C .
- **dci** (`numpy.ndarray(shape=(n_cov, n_cov))`) – Low-rank inverse matrix of $dc \cdot dc.T$.
- **dcr** (*int*) – Rank of dci .
- **dimreduce** (`numpy.ndarray(shape=(ny,), dtype='uint')` or *int*.) – If each vector y doesn't have full rank in the first place, this parameter is the loss of degree of freedom to allow for accurate P-value computation.

Returns

- **vx** (*any*) – vx from input for information passing.
- **vy** (*any*) – vy from input for information passing.
- **pv** (`numpy.ndarray(shape=(n_x, n_y))`) – P-values of association testing ($\gamma == 0$).
- **gamma** (`numpy.ndarray(shape=(n_x, n_y))`) – Maximum likelihood estimator of γ in model.
- **alpha** (`numpy.ndarray(shape=(n_x, n_y, n_cov))`) – Maximum likelihood estimator of α in model.

- **var_x** (*numpy.ndarray(shape=(n_x,))*) – Variance of dx unexplained by covariates C.
- **var_y** (*numpy.ndarray(shape=(n_y,))*) – Variance of dy unexplained by covariates C.

`normalisr.association.association_test_2(vx, vy, dx, dy, dc, sselectx, dimreduce=0)`

Like `association_test_1`, but takes a different subset of samples for each x.

See `association_test_1` for additional details.

Parameters

- **vx** (*any*) – Starting indices of dx. Only used for information passing.
- **vy** (*any*) – Starting indices of dy. Only used for information passing.
- **dx** (*numpy.ndarray(shape=(n_x, n_cell))*) – Predictor matrix for a list of vector x to be tested, e.g. gene expression or grouping.
- **dy** (*numpy.ndarray(shape=(n_y, n_cell))*) – Target matrix for a list of vector y to be tested, e.g. gene expression.
- **dc** (*numpy.ndarray(shape=(n_cov, n_cell))*) – Covariate matrix as C.
- **sselectx** (*numpy.ndarray(shape=(n_x, n_cell), dtype=bool)*) – Subset of samples to use for each x.
- **dimreduce** (*numpy.ndarray(shape=(ny,), dtype='uint') or int.*) – If each vector y doesn't have full rank in the first place, this parameter is the loss of degree of freedom to allow for accurate P-value computation.

Returns

- **vx** (*any*) – vx from input for information passing.
- **vy** (*any*) – vy from input for information passing.
- **p** (*numpy.ndarray(shape=(n_x, n_y))*) – P-values of association testing (`gamma==0`).
- **gamma** (*numpy.ndarray(shape=(n_x, n_y))*) – Maximum likelihood estimator of gamma in model.
- **alpha** (*numpy.ndarray(shape=(n_x, n_y, n_cov))*) – Maximum likelihood estimator of alpha in model.
- **var_x** (*numpy.ndarray(shape=(n_x,))*) – Variance of dx unexplained by covariates C.
- **var_y** (*numpy.ndarray(shape=(n_x, n_y))*) – Variance of dy unexplained by covariates C.

`normalisr.association.association_test_4(vx, vy, prod, prody, proddy, na, dimreduce=0, **ka)`

Like `association_test_1`, but regards all other x's as covariates when testing each x.

See `association_test_1` for additional details. Other x's are treated as covariates but their coefficients (alpha) would not be returned to reduce memory footprint.

Parameters

- **vx** (*any*) – Starting indices of dx.
- **vy** (*any*) – Starting indices of dy. Only used for information passing.
- **prod** (*numpy.ndarray(shape=(n_x+n_cov, n_x+n_cov))*) – $A @ A.T$, where $A = \text{numpy.block}([dx, dc])$.
- **prody** (*numpy.ndarray(shape=(n_x+n_cov, n_y))*) – $A @ dy.T$, where $A = \text{numpy.block}([dx, dc])$.
- **proddy** (*numpy.ndarray(shape=(n_y,))*) – $(dy ** 2).sum(axis=1)$

- **na** (*tuple*) – (n_x, n_y, n_cov, n_cell, lenx). Numbers of (x's, y's, covariates, cells, x's to compute association for)
- **dimreduce** (*numpy.ndarray (shape=(ny,), dtype='uint') or int.*) – If each vector y doesn't have full rank in the first place, this parameter is the loss of degree of freedom to allow for accurate P-value computation.
- **ka** (*dict*) – Keyword arguments passed to inv_rank.

Returns

- **vx** (*any*) – vx from input for information passing.
- **vy** (*any*) – vy from input for information passing.
- **pv** (*numpy.ndarray(shape=(n_x, n_y))*) – P-values of association testing (gamma==0).
- **gamma** (*numpy.ndarray(shape=(n_x, n_y))*) – Maximum likelihood estimator of gamma in model.
- **alpha** (*numpy.ndarray(shape=(n_x, n_y, n_cov))*) – Maximum likelihood estimator of alpha in model.
- **var_x** (*numpy.ndarray(shape=(lenx,))*) – Variance of dx unexplained by covariates C.
- **var_y** (*numpy.ndarray(shape=(lenx, n_y))*) – Variance of dy unexplained by covariates C.

`normalisr.association.association_tests_1(dx, dy, dc, bsx=0, bsy=0, nth=1, lowmem=True, **ka)`

Performs association tests between all pairs of two (or one) variables. Performs parallel computation with multiple processes on the same machine.

Model for differential expression between X and Y: $Y = \text{gamma} * X + \text{alpha} * C + \text{epsilon}$,
epsilon~i.i.d. $N(0, \text{sigma}^{**2})$.

Test statistic: conditional R^{**2} (or proportion of variance explained) between Y and X.

Null hypothesis: gamma=0.

Parameters

- **dx** (*numpy.ndarray (shape=(n_x, n_cell))*) – Normalized matrix X.
- **dy** (*numpy.ndarray (shape=(n_y, n_cell), dtype=float) or None*) – Normalized matrix Y. If None, indicates dy=dx, i.e. self-association between pairs within X.
- **dc** (*numpy.ndarray (shape=(n_cov, n_cell), dtype=float)*) – Normalized covariate matrix C.
- **bsx** (*int*) – Batch size, i.e. number of Xs in each computing batch. Defaults to 500.
- **bsy** (*int*) – Batch size, i.e. number of Xs in each computing batch. Defaults to 500.
- **nth** (*int*) – Number of parallel processes. Set to 0 for using automatically detected CPU counts.
- **ka** (*dict*) – Keyword arguments passed to normalisr.association.association_test_1. See below.

Returns

- **P-values** (*numpy.ndarray(shape=(n_x, n_y))*) – Differential expression P-value matrix.
- **dot** (*numpy.ndarray(shape=(n_x, n_y))*) – Inner product between X and Y pairs, after removing covariates.

- **alpha** (*numpy.ndarray(shape=(n_x,n_y,n_cov)) or None.*) – Maximum likelihood estimators of alpha, separately tested for each grouping. None if lowmem=True.
- **varx** (*numpy.ndarray(shape=(n_x))*) – Variance of grouping after removing covariates.
- **vary** (*numpy.ndarray(shape=(n_y))*) – Variance of gene expression after removing covariates. It can depend on the grouping being tested depending on parameter single.

Keyword Arguments **dimreduce** (*numpy.ndarray(shape=(n_y,), dtype=int) or int*) – If dy doesn't have full rank, such as due to prior covariate removal (although the recommended method is to leave covariates in dc), this parameter allows to specify the loss of ranks/degrees of freedom to allow for accurate P-value computation. Default is 0, indicating no rank loss.

`normalisr.association.inv_rank(m, tol=1e-08, method='auto', logger=None, mpc=0, qr=0, **ka)`

Computes matrix (pseudo-)inverse and rank with SVD.

Eigenvalues smaller than tol*largest eigenvalue are set to 0. Rank of inverted matrix is also returned. Provides to limit the number of eigenvalues to speed up computation. Broadcasts to the last 2 dimensions of the matrix.

Parameters

- **m** (*numpy.ndarray(shape=(...,n,n), dtype=float)*) – 2-D or higher matrix to be inverted
- **tol** (*float*) – Eigenvalues < tol*maximum eigenvalue are treated as zero.
- **method** (*str*) – Method to compute eigenvalues:
 - auto: Uses scipy for n<mpc or mpc==0 and sklearn otherwise
 - scipy: Uses scipy.linalg.svd
 - scipys: NOT IMPLEMENTED. Uses scipy.sparse.linalg.svds
 - sklearn: Uses sklearn.decomposition.TruncatedSVD
- **logger** (*object*) – Logger to output warning. Defaults (None) to logging module
- **mpc** (*int*) – Maximum rank or number of eigenvalues/eigenvectors to consider. Defaults to 0 to disable limit. For very large input matrix, use a small value (e.g. 500) to save time at the cost of accuracy.
- **qr** (*int*) – Whether to use QR decomposition for matrix inverse. Only effective when method=sklearn, or =auto that defaults to sklearn. * 0: No * 1: Yes with default settings * 2+: Yes with n_iter=qr for sklearn.utils.extmath.randomized_svd
- **ka** (*Keyword args passed to method*) –

Returns

- **mi** (*numpy.ndarray(shape=(...,n,n), dtype=float)*) – Pseudo-inverse matrices
- **r** (*numpy.ndarray(shape=(...), dtype=int) or int*) – Matrix ranks

`normalisr.association.prod1(vx, vy, dx, dy)`

Pickleable function for matrix product that keeps information

Parameters

- **vx** (*any*) – Information passed
- **vy** (*any*) – Information passed
- **dx** (*numpy.ndarray(shape=(...,n))*) – Matrix for multiplication

- **dy** (*numpy.ndarray* (*shape*=(*..*,*n*))) – Matrix for multiplication

Returns

- **vx** (*any*) – vx
- **vy** (*any*) – vy
- **product** (*numpy.ndarray*(*shape*=(*...*))) – dx@dy.T

3.2 normalisr.binnet

`normalisr.binnet.bh` (*pv*, *weight=None*)

Converts P-values to Q-values using Benjamini–Hochberg procedure.

Parameters

- **pv** (*numpy.ndarray* (*shape*=(*n*,))) – P-values.
- **weight** (*numpy.ndarray* (*shape*=(*n*,)) or *None*) – Weight of each P-value. Defaults (*None*) to equal.

Returns Q-values.

Return type *numpy.ndarray*(*shape*=(*n*,))

References

Controlling the false discovery rate: a practical and powerful approach to multiple testing, Benjamini and Hochberg. 1995

`normalisr.binnet.binnet` (*net*, *qcut*)

Binarizes P-value co-expression network to thresholded Q-value network.

Q-values are computed separately per row to account for differences in the number of genes co-expressed, especially by master regulators, using Benjamini–Hochberg procedure. Co-expression Q-value matrix is thresholded for return.

Parameters

- **net** (*numpy.ndarray* (*shape*=(*n_gene*,*n_gene*), *dtype=float*)) – Symmetric co-expression P-value matrix.
- **qcut** (*float*) – Cutoff for Q-value network.

Returns Binarized, asymmetric co-expression matrix.

Return type *numpy.ndarray*(*shape*=(*n_gene*,*n_gene*),*dtype=bool*)

`normalisr.binnet.nodiag` (*d*, *split=False*)

Removes diagonals from 2D matrix.

Parameters

- **d** (*numpy.ndarray* (*shape*=[*n1*,*n2*])) – Input matrix
- **split** (*bool*) – Whether to output list of 1D arrays for each row instead of a single concatenated 1D array.

Returns

- **If split** (*List of numpy.ndarray*(*shape*=[*n2* or *n2-1*]) for each row.)

- **Else** (*numpy.ndarray(shape=(n1*n2-min(n1,n2),))*) – Reduced 1D array

`normalisr.binnet.rediag(d, fill=0, shape=None)`

Converts a ‘nodiag’ed vector back to matrix of original shape.

Parameters

- **d** (*numpy.ndarray(shape=(n,))*) – Output array from `nodiag(split=False)`
- **fill** (*any*) – Values to fill in the diagonal entries
- **shape** (*tuple or None*) – Shape of original matrix. If omitted (*None*), assume original is a square matrix.

Returns Recovered original matrix.

Return type `numpy.ndarray(shape=shape)`

3.3 normalisr.coex

`normalisr.coex.coex(dt, dc, **ka)`

Performs co-expression analyses for all gene pairs.

Performs parallel computation with multiple processes on the same machine.

Model for co-expression between genes i & j: $X_i = \gamma X_j + \alpha C + \epsilon$,

$\epsilon \sim \text{i.i.d. } N(0, \sigma^2)$.

Test statistic: conditional R^2 (or proportion of variance explained) between X_i and X_j .

Null hypothesis: $\gamma = 0$.

Parameters

- **dt** (*numpy.ndarray(shape=(n_gene, n_cell), dtype=float)*) – Normalized expression matrix X .
- **dc** (*numpy.ndarray(shape=(n_cov, n_cell), dtype=float)*) – Normalized covariate matrix C .
- **ka** (*dict*) – Keyword arguments passed to `normalisr.association.association_test_1`. See below.

Returns

- **P-values** (*numpy.ndarray(shape=(n_gene, n_gene))*) – Co-expression P-value matrix.
- **dot** (*numpy.ndarray(shape=(n_gene, n_gene))*) – Inner product of expression between gene pairs, after removing covariates.
- **var** (*numpy.ndarray(shape=(n_gene,))*) – Variance of gene expression after removing covariates. Pearson $R = (((\text{dot}/\text{numpy.sqrt(var)}).T)/\text{numpy.sqrt(var)}).T$.

Keyword Arguments

- **bs** (*int*) – Batch size, i.e. number of genes in each computing batch. Use 0 for default: Data transfer limited to 1GB, capped at `bs=500`.
- **nth** (*int*) – Number of parallel processes. Set to 0 for using automatically detected CPU counts.

- **dimreduce** (*numpy.ndarray* (*shape*=(*n_gene*,), *dtype=int*) or *int*) – If *dt* doesn't have full rank, such as due to prior covariate removal (although the recommended method is to leave covariates in *dc*), this parameter allows to specify the loss of ranks/degrees of freedom to allow for accurate P-value computation. Default is 0, indicating no rank loss.

3.4 normalisr.de

`normalisr.de.de` (*dg*, *dt*, *dc*, *bs*=0, *nth*=1, *single*=0, ***ka*)

Performs differential expression analyses for all genes against all groupings.

Allows multiple options to treat other groupings when testing on one grouping.

Performs parallel computation with multiple processes on the same machine.

Model for differential expression between gene Y and grouping X: $Y = \gamma X + \alpha C + \epsilon$, $\epsilon \sim \text{i.i.d. } N(0, \sigma^2)$.

Test statistic: conditional R^2 (or proportion of variance explained) between Y and X.

Null hypothesis: $\gamma = 0$.

Parameters

- **dg** (*numpy.ndarray* (*shape*=(*n_group*, *n_cell*))) – Grouping matrix for a list of X to be tested, e.g. grouping by gene knock-out.
- **dt** (*numpy.ndarray* (*shape*=(*n_gene*, *n_cell*), *dtype=float*)) – Normalized expression matrix Y.
- **dc** (*numpy.ndarray* (*shape*=(*n_cov*, *n_cell*), *dtype=float*)) – Normalized covariate matrix C.
- **bs** (*int*) – Batch size, i.e. number of groupings and genes in each computing batch. For *single*=0,1, splits groupings & genes. Defaults to 500. For *single*=4, only allows *bs*=0 (default automatic setting).
- **nth** (*int*) – Number of parallel processes. Set to 0 for using automatically detected CPU counts.
- **single** (*int*) – Option to deal with other groupings when testing one groupings v.s. gene expression.
 - 0: Ignores other groupings.
 - 1: Excludes all cells belonging to any other grouping (*value*==1), assuming *dg*=0,1 only.
 - 4: Treats other groupings as covariates for mean expression.
- **ka** (*dict*) – Keyword arguments passed to `normalisr.association.association_test_*`. See below.

Returns

- **P-values** (*numpy.ndarray* (*shape*=(*n_group*, *n_gene*))) – Differential expression P-value matrix.
- **gamma** (*numpy.ndarray* (*shape*=(*n_group*, *n_gene*))) – Differential expression log fold change matrix.
- **alpha** (*numpy.ndarray* (*shape*=(*n_group*, *n_gene*, *n_cov*))) – Maximum likelihood estimators of alpha, separately tested for each grouping.

- **varg** (*numpy.ndarray(shape=(n_group))*) – Variance of grouping after removing covariates.
- **var** (*numpy.ndarray(shape=(n_group,n_gene))*) – Variance of gene expression after removing covariates. It can depend on the grouping being tested depending on parameter **single**.

Keyword Arguments

- **dimreduce** (*numpy.ndarray(shape=(n_gene,), dtype=int) or int*) – If **dt** doesn't have full rank, such as due to prior covariate removal (although the recommended method is to leave covariates in **dc**), this parameter allows to specify the loss of ranks/degrees of freedom to allow for accurate P-value computation. Default is 0, indicating no rank loss.
- **method** (*str, only for single=4*) – Method to compute eigenvalues in SVD-based matrix inverse (for removal of covariates):
 - **auto**: Uses **scipy** for **n_matrix**<**mpc** or **mpc**==0 and **sklearn** otherwise. Default.
 - **scipy**: Uses **scipy.linalg.svd**.
 - **scipys**: NOT IMPLEMENTED. Uses **scipy.sparse.linalg.svds**.
 - **sklearn**: Uses **sklearn.decomposition.TruncatedSVD**.
- **mpc** (*int, only for single=4*) – Uses only the top **mpc** singular values as non-zero in SVD-based matrix inverse. Here effectively reduces covariates to their top principal components. This reduction is performed after including other groupings as additional covariates. Defaults to 0 to disable dimension reduction. For very large grouping matrix, use a small value (e.g. 100) to save time at the cost of accuracy.
- **qr** (*int, only for single=4*) – Whether to use QR decomposition method for SVD in matrix inverse. Only effective when **method**=**sklearn**, or **=auto** and defaults to **sklearn**. Takes the following values:
 - 0: No (default).
 - 1: Yes with default settings.
 - 2+: Yes with **n_iter**=**qr** for **sklearn.utils.extmath.randomized_svd**.
- **tol** (*float, only for single=4*) – Eigenvalues < **tol***(maximum eigenvalue) are treated as zero in SVD-based matrix inverse. Default is 1E-8.

3.5 normalisr.gocovt

```
normalisr.gocovt.goe(genelist, go_file, goa_file, bg=None, nmin=5, conversion=None, evidence_set={'EXP', 'HDA', 'HGI', 'HMP', 'HTP', 'IBA', 'IBD', 'IDA', 'IGI', 'IKR', 'IMP', 'IPI', 'IRD', 'ISA', 'ISM', 'ISO', 'ISS'})
```

Finds GO enrichment with **goatools** (0.7.11 tested).

WARNING: This method is inexact for multi-maps in gene name conversion. However, it has a negligible effect in top GO component removal in single-cell co-expression.

Parameters

- **genelist** (*list of str*) – Genes to search for enrichment.
- **go_file** (*str*) – File path for GO DAG (downloadable at <http://geneontology.org/docs/download-ontology/>).
- **goa_file** (*str*) – File path for GO associations. See parameter **conversion**.

- **bg** (*list of str*) – Background genes.
- **nmin** (*int*) – Minimum number of principal genes required in GO.
- **conversion** (*tuple*) – Conversion of **gene ID system** from gene list to the GO annotation.
 - name_from: Gene naming system of genelist. For gene names, use ‘symbol,alias’.
 - name_to: Gene naming system of goa_file. Examples:
 - * Human: use ‘uniprot.Swiss-Prot’ (for GO annotations downloaded from http://geneontology.org/gene-associations/goa_human.gaf.gz).
 - * Mouse: use ‘MGI’ (for GO annotations downloaded from <http://current.geneontology.org/annotations/mgi.gaf.gz>).
 - species: Species for gene name conversion. Examples: ‘human’, ‘mouse’.
- **evidence_set** (*set of str*) – **GO evidences** to include. Defaults to non-expression based results to avoid circular reasoning bias.

Returns

- **goe** (*pandas.DataFrame*) – GO enrichment.
- **gotop** (*str*) – Top enriched GO ID
- **genes** (*list of str or None*) – Intersection list of genes in gotop and also bg. None if bg is None.

`normalisr.gocovt.gotop(net, namet, go_file, goa_file, n=100, **ka)`

Finds the top variable GO enrichment of top principal genes in the binary co-expression network.

Principal genes are those with most co-expressed genes. They reflect the most variable pathways in the dataset. When the variable pathways are housekeeping related, they may conceal cell-type-specific co-expression patterns from being observed and understood. This function identifies the most variable pathway with gene ontology enrichment study of the top principal genes. Background genes are all genes provided.

Parameters

- **net** (*numpy.ndarray (shape=(n_gene, n_gene), dtype=bool)*) – Binary co-expression network matrix.
- **namet** (*list of str*) – Gene names matching the rows and columns of net.
- **go_file** (*str*) – File path for GO DAG (downloadable at <http://geneontology.org/docs/download-ontology/>).
- **goa_file** (*str*) – File path of GO annotation. See parameter **conversion** in `normalisr.gocovt.goe`.
- **n** (*int*) – Number of top principal genes to include for GO enrichment. Default is 100, giving good performance in general.
- **ka** (*dict*) – **IMPORTANT**: Keyword arguments passed to `normalisr.gocovt.goe` to determine how to perform GO enrichment study. If you see no gene mapped, check your gene name conversion rule in **conversion** parameter of `normalisr.gocovt.goe`. GO annotation have a specific gene ID system.

Returns

- **principals** (*list of str*) – List of principal genes.
- **goe** (*pandas.DataFrame*) – GO enrichment results.

- **gotop** (*str*) – Top enriched GO ID.
- **genes** (*list of str*) – List of genes in the gotop GO ID.

`normalisr.gocovt.pc1(d)`

Computes top principal component.

Parameters `d` (`numpy.ndarray (shape=(n1, n2))`) –

Returns

Return type `numpy.ndarray(shape=(n2,))`

`normalisr.gocovt.pccovt(dt, dc, namet, genes, condcov=True)`

Introduces an extra covariate from the top principal component of given genes.

The extra covariate is the top principal component of normalized expressions of the selected genes. Adding a covariate from housekeeping pathway can reveal cell-type-specific activities in co-expression networks.

Parameters

- **dt** (`numpy.ndarray (shape=(n_gene, n_cell))`) – Normalized expression matrix.
- **dc** (`numpy.ndarray (shape=(n_cov, n_cell))`) – Existing normalized covariate matrix.
- **namet** (*list of str*) – List of gene names for rows in dt.
- **genes** (*list of str*) – List of gene names to include in finding top PC of their expression as an extra covariate.
- **condcov** (*bool*) – Whether to condition on existing covariates before computing top PC. Default: True.

Returns New normalized covariate matrix.

Return type `numpy.ndarray(shape=(n_cov+1,n_cell))`

3.6 normalisr.lcpm

`normalisr.lcpm.lcpm(reads, normalize=True, nth=0, ntot=None, varscale=0, seed=None, lowmem=True)`

Computes Bayesian log CPM from raw read counts.

The technical sampling process is modelled as a Binomial distribution. The logCPM given read counts is a Bayesian inference problem and follows (shifted) Beta distribution. We use the expectation of posterior logCPM as the estimated expression levels. Resampling function is also provided to account for variances in the posterior distribution.

Warning: Modifying keyword arguments other than `nth` or `seed` is neither recommended nor supported for function ‘lcpm’. Do so at your own risk.

Parameters

- **reads** (`numpy.ndarray (shape=(n_gene, n_cell), dtype='uint')`) – Read count matrix.
- **normalize** (*bool*) – Whether to normalize output to logCPM per cell. Default: True.
- **nth** (*int*) – Number of threads to use. Defaults to 0 to use all cores automatically detected.

- **ntot** (*int*) – Manually sets value of total number of reads in binomial distribution. Since the posterior distribution stabilizes quickly as ntot increases, a large number, e.g. 1E9 is good for general use. Defaults to None to disable manual value.
- **varscale** (*float*) – Resamples estimated expression using the posterior Beta distribution. varscale sets the scale of variance than its actual value from the posterior distribution. Defaults to 0, to compute expectation with no variance.
- **seed** (*int*) – Initial random seed if set.
- **lowmem** (*bool*) – Low memory mode disable mean and var in Returns and therefore saves memory.

Returns

- **lcpm** (*numpy.ndarray(shape=(n_gene,n_cell))*) – Estimated expression as logCPM from read counts.
- **mean** (*numpy.ndarray(shape=(n_gene,n_cell)) or None*) – Mean/Expectation of lcpm's every entry's posterior distribution. None if lowmem=True.
- **var** (*numpy.ndarray(shape=(n_gene,n_cell)) or None*) – Variance of lcpm's every entry's posterior distribution. None if lowmem=True.
- **cov** (*numpy.ndarray(shape=(3,n_cell))*) – Cellular summary covariates computed from read count matrix that may confound lcpm. Contains:
 - cov[0]: Log total read count per cell
 - cov[1]: Number of 0-read genes per cell
 - cov[2]: cov[0]**2

`normalisr.lcpm.scaling_factor(dt, varname='nt0mean', v0=0, v1='max')`

Computes scaling factor of variance normalization for every gene.

Lowly expressed genes need full variance normalization because of technical confounding from sequencing depth. Highly expressed genes do not need variance normalization because they are already accurately measured. The scaling factor operates as an exponential factor on the variance normalization scale for each gene. It should be maximum/minimum for genes with lowest/highest expression.

Warning: Modifying keyword arguments is neither recommended nor supported for function 'scaling_factor'. Do so at your own risk.

Parameters

- **dt** (*numpy.ndarray(shape=(n_gene,n_cell), dtype='uint')*) – Read count matrix.
- **varname** (*str*) – Variable used to compute scaling factor for each gene. Can be:
 - logtpropmean: $\log(\text{dt.mean}(\text{axis}=1)/\text{dt.mean}(\text{axis}=1).\text{sum}())$
 - logtmeanprop: $\log((\text{dt}/\text{dt.sum}(\text{axis}=0)).\text{mean}(\text{axis}=1))$
 - nt0mean: $(\text{dt}==0).\text{mean}(\text{axis}=1)$
 - lognt0mean: $\log((\text{dt}==0).\text{mean}(\text{axis}=1))$
 - log1-nt0mean: $\log(1-(\text{dt}==0).\text{mean}(\text{axis}=1))$
 Defaults to nt0mean.
- **v0** (*float*) – Variable values to set scaling factor to 0 (for v0) and 1 (for v1). Linear assignment is applied for values inbetween. Can be:

- max: max
- min: min
- any float: that float
- **v1** (*float*) – Variable values to set scaling factor to 0 (for v0) and 1 (for v1). Linear assignment is applied for values inbetween. Can be:
 - max: max
 - min: min
 - any float: that float

Returns Scaling factor of variance normalization for each gene

Return type `numpy.ndarray(shape=(n_gene,))`

`normalisr.lcpm.trigamma(x)`

Tri-gamma function .

Parameters **x** (*float or numpy.ndarray*) – Input value(s)

Returns

Return type `float or numpy.ndarray`

3.7 normalisr.norm

`normalisr.norm.compute_var(dt, dc, stepmax=1, eps=1e-06)`

Computes variance normalization scale for each cell.

Performs a log-linear fit of the variance of each cell with covariates. Optionally use Expectation-Maximization(EM)-like method to iteratively fit mean and variance. For EM-like method, early-stopping is suggested because of overfitting issues.

Parameters

- **dt** (*numpy.ndarray(shape=(n_gene, n_cell))*) – Bayesian LogCPM expression level matrix.
- **dc** (*numpy.ndarray(shape=(n_cov, n_cell))*) – Covariate matrix.
- **stepmax** (*int*) – Maximum number of EM-like iterations of mean and variance normalization. Stop of iteration is also possible when relative accuracy target is reached. Defaults to 1, indicating no iterative normalization.
- **eps** (*float*) – Relative accuracy target for early stopping. Constrains the maximum relative difference of fitted variance across cells compared to the last step. Defaults to 1E-6.

Returns Inverse sqrt of fitted variance for each cell, i.e. the multiplier for variance normalization. For iterative normalization, the optimal step will be returned, defined as having the minimal max relative change across cells.

Return type `numpy.ndarray(shape=(n_cell,))`

`normalisr.norm.normcov(dc, c=True)`

Normalizes each continuous covariate to 0 mean and unit variance.

Optionally introduces constant 1 covariate as intercept. Categorical covariates should be in binary/one-hot form, and will be left unchanged.

Parameters

- **dc** (*numpy.ndarray* (*shape*=(*n_cov*, *n_cell*))) – Current covariate matrix. Use empty matrix with *n_cov*=0 if no covariate.
- **c** (*bool*) – Whether to introduce a constant 1 covariate.

Returns Processed covariate matrix.

Return type *numpy.ndarray*(*shape*=(*n_cov*+1 if *c* else *n_cov*, *n_cell*))

`normalisr.norm.normvar` (*dt*, *dc*, *w*, *wt*, *dextra*=None, *cat*=1, *nth*=1, *bs*=500, *keepvar*=True)

Performs mean and variance normalizations.

Expression levels are normalized at mean and then at variance levels. Effectively each gene *x* is multiplied by *w****wt*[*x*] before removing covariates as *dc**(*w****wt*[*x*]). Continuous covariates are normalized at variance levels. Effectively covariates are transformed to *dc***w*. Therefore, variance normalization for expression are scaled differently for each gene.

Parameters

- **dt** (*numpy.ndarray* (*shape*=(*n_gene*, *n_cell*))) – Bayesian logCPM matrix.
- **dc** (*numpy.ndarray* (*shape*=(*n_cov*, *n_cell*))) – Covariate matrix.
- **w** (*numpy.ndarray* (*shape*=(*n_cell*,))) – Computed variance normalization multiplier.
- **wt** (*numpy.ndarray* (*shape*=(*n_gene*,))) – Computed scaling factor for each gene.
- **dextra** (*numpy.ndarray* (*shape*=(*n_extra*, *n_cell*))) – Extra data matrix also to be normalized like continuous covariates.
- **cat** (*int*) – Whether to normalize categorical/binary covariates (those with only 0 or 1s). Defaults to 1.
 - 0: No
 - 1: No except constant-1 covariate (intercept)
 - 2: Yes
- **nth** (*int*) – Number of parallel threads.
- **bs** (*int*) – Batch size for each job.
- **keepvar** (*bool*) – Whether to maintain the variance of each gene invariant in mean normalization step. If so, expression variances are scaled back to original after mean normalization and before variance normalization. This function only affects overall variance level and its downstreams (e.g. differential expression log fold change). This function would not affect P-value computation. Default: True.

Returns

- (*dtn*, *dcn*) or (*dtn*, *dcn*, *dextran*) if *dextra* is not None
- **dtn** (*numpy.ndarray*(*shape*=(*n_gene*, *n_cell*))) – Normalized gene expression matrix.
- **dcn** (*numpy.ndarray*(*shape*=(*n_cov*, *n_cell*))) – Normalized covariate matrix.
- **dextran** (*numpy.ndarray*(*shape*=(*n_extra*, *n_cell*))) – Normalized extra data matrix.

`normalisr.norm.normvar1` (*dt*, *dc*, *w2*)

A single-threaded function to normalize a subset of transcriptome.

Should not be invoked directly.

Parameters

- **dt** (*numpy.ndarray (shape=(n_gene, n_cell))*) – Bayesian logCPM matrix subset.
- **dc** (*numpy.ndarray (shape=(n_cov, n_cell))*) – Covariate matrix.
- **w2** (*numpy.ndarray (shape=(n_gene,))*) – Intermediate weight computed in normvar for subset.

Returns Normalized gene expression matrix subset.

Return type *numpy.ndarray(shape=(n_gene, n_cell))*

3.8 normalisr.parallel

`normalisr.parallel.autopooler(n, it, *a, chunksize=1, dummy=False, return_iter=False, unordered=False, **ka)`

Uses multiprocessing.Pool or multiprocessing.dummy.Pool to run iterator in parallel.

Parameters

- **n** (*int*) – Number of parallel processes. Set to 0 to use auto detected CPU count.
- **it** (*iterator of (function, tuple, dict)*) – Each iteration computes **function(*tuple,**dict)**. **function** must be picklable, i.e. a base level function in a module or file.
- **a** (*tuple*) – Arguments passed to Pool.
- **chunksize** (*int*) – Number of iterations passed to each process each time.
- **dummy** (*bool*) – Whether to use multiprocessing.dummy instead
- **return_iter** (*bool*) – Not Implemented. Whether to return iterator of results instead. If not, return list of results.
- **unordered** (*bool*) – Whether the order of output matters.
- **ka** (*dict*) – Keyword arguments passed to Pool

Returns Results returned by function(*tuple,**dict), in same order of the iterator if not unordered.

Return type list (or iterator if return_iter) of any

3.9 normalisr.qc

`normalisr.qc.qc_outlier(dw, pcut=1e-10, outrate=0.02)`

Quality control by removing cell outliers by variance.

Fit normal distribution on the inverse sqrt variance to detect outliers. This is performed by iterative estimation of normal distribution with non-outliers and then determination of outliers with the normal distribution.

Parameters

- **dw** (*numpy.ndarray (shape=(n_cell,))*) – Fitted inverse sqrt cell variance.
- **pcut** (*float*) – Bonferroni P-value cutoff for asserting outliers in a normal distribution of fitted cell variance. Default: 1E-10.
- **outrate** (*float*) – Upper bound of proportion of outliers on either side of variance distribution. Used for initial outlier assignment and final validity check. Default: 0.02.

Returns Whether each cell passed QC

Return type `numpy.ndarray(shape=(n_cell,),dtype=bool)`

`normalisr.qc.qc_reads(reads, n_gene, nc_gene, ncp_gene, n_cell, nt_cell, ntp_cell)`

Quality control by bounding read counts.

Quality control is performed separately on genes based on their cell statistics and on cells based on their gene statistics, iteratively until dataset remains unchanged. A gene or cell is removed if any of the QC criteria is violated at any time in the iteration. All QC parameters can be set to 0 to disable QC filtering for that criterion.

Parameters

- **reads** (`numpy.ndarray((n_gene, n_cell), dtype='uint')`) – Read count matrix.
- **n_gene** (`int`) – Lower bound on total read counts for gene QC.
- **nc_gene** (`int`) – Lower bound on number of expressed cells for gene QC.
- **ncp_gene** (`float`) – Lower bound on proportion of expressed cells for gene QC.
- **n_cell** (`int`) – Lower bound on total read counts for cell QC.
- **nt_cell** (`int`) – Lower bound on number of expressed genes for cell QC.
- **ntp_cell** (`float`) – Lower bound on proportion of expressed genes for cell QC.

Returns

- **genes_select** (`numpy.ndarray(dtype='uint')`) – Array of indices of genes passed QC.
- **cells_select** (`numpy.ndarray(dtype='uint')`) – Array of indices of cells passed QC.

INDICES AND TABLES

- `genindex`
- `modindex`
- `search`

PYTHON MODULE INDEX

n

- `normalizr.association`, [15](#)
- `normalizr.binnet`, [19](#)
- `normalizr.coex`, [20](#)
- `normalizr.de`, [21](#)
- `normalizr.gocovt`, [22](#)
- `normalizr.lcpm`, [24](#)
- `normalizr.norm`, [26](#)
- `normalizr.normalizr`, [6](#)
- `normalizr.parallel`, [28](#)
- `normalizr.qc`, [28](#)

A

association_test_1() (in module *normalisr.association*), 15
 association_test_2() (in module *normalisr.association*), 16
 association_test_4() (in module *normalisr.association*), 16
 association_tests_1() (in module *normalisr.association*), 17
 autopooler() (in module *normalisr.parallel*), 28

B

bh() (in module *normalisr.binnet*), 19
 binnet() (in module *normalisr.binnet*), 19
 binnet() (in module *normalisr.normalisr*), 6

C

coex() (in module *normalisr.coex*), 20
 coex() (in module *normalisr.normalisr*), 6
 compute_var() (in module *normalisr.norm*), 26
 compute_var() (in module *normalisr.normalisr*), 7

D

de() (in module *normalisr.de*), 21
 de() (in module *normalisr.normalisr*), 7

G

goe() (in module *normalisr.gocovt*), 22
 gotop() (in module *normalisr.gocovt*), 23
 gotop() (in module *normalisr.normalisr*), 8

I

inv_rank() (in module *normalisr.association*), 18

L

lcpm() (in module *normalisr.lcpm*), 24
 lcpm() (in module *normalisr.normalisr*), 9

M

module
 normalisr.association, 15

normalisr.binnet, 19
normalisr.coex, 20
normalisr.de, 21
normalisr.gocovt, 22
normalisr.lcpm, 24
normalisr.norm, 26
normalisr.normalisr, 6
normalisr.parallel, 28
normalisr.qc, 28

N

nodiag() (in module *normalisr.binnet*), 19
normalisr.association
 module, 15
normalisr.binnet
 module, 19
normalisr.coex
 module, 20
normalisr.de
 module, 21
normalisr.gocovt
 module, 22
normalisr.lcpm
 module, 24
normalisr.norm
 module, 26
normalisr.normalisr
 module, 6
normalisr.parallel
 module, 28
normalisr.qc
 module, 28
normcov() (in module *normalisr.norm*), 26
normcov() (in module *normalisr.normalisr*), 10
normvar() (in module *normalisr.norm*), 27
normvar() (in module *normalisr.normalisr*), 10
normvar1() (in module *normalisr.norm*), 27

P

pc1() (in module *normalisr.gocovt*), 24
 pccovt() (in module *normalisr.gocovt*), 24
 pccovt() (in module *normalisr.normalisr*), 11

`prod1()` (*in module normalisr.association*), 18

Q

`qc_outlier()` (*in module normalisr.normalisr*), 11

`qc_outlier()` (*in module normalisr.qc*), 28

`qc_reads()` (*in module normalisr.normalisr*), 12

`qc_reads()` (*in module normalisr.qc*), 29

R

`rediaq()` (*in module normalisr.binnet*), 20

S

`scaling_factor()` (*in module normalisr.lcpm*), 25

`scaling_factor()` (*in module normalisr.normalisr*),
12

T

`trigamma()` (*in module normalisr.lcpm*), 26