

OrderedList

April 19, 2010

`compareLists`

Compare Ordered Lists with Weighted Overlap Score

Description

The two orderings received as parameters are compared using the weighted overlap score and compared with a random distribution of that score (yielding an empirical p-value).

Usage

```
compareLists(ID.List1, ID.List2, mapping = NULL,  
             two.sided=TRUE, B = 1000, alphas = NULL,  
             invar.q = 0.5, min.weight = 1e-5,  
             no.reverse=FALSE)
```

Arguments

<code>ID.List1</code>	first ordered list of identifiers to be compared.
<code>ID.List2</code>	second ordered list to be compared, must have the same length as <code>ID.List1</code> .
<code>mapping</code>	maps identifiers between the two lists. This is a matrix with two columns. All items in <code>ID.List1</code> must match to exactly one entry of column 1 of the mapping, each element in <code>ID.List2</code> must match exactly one element in column 2 of the mapping. If mapping is <code>NULL</code> , the two lists are expected to contain the same identifiers and there must be a one-to-one relationship between the two.
<code>two.sided</code>	whether the score is to be computed considering both ends of the list, or just the top members.
<code>B</code>	the number of permutations used to estimate empirical p-values.
<code>alphas</code>	a set of alpha candidates to be evaluated. If set to <code>NULL</code> , alphas are determined such that reasonable maximal ranks to be considered result.
<code>invar.q</code>	quantile of genes expected to be invariant. These are not used during shuffling, since they are expected to stay away from the ends of the lists, even when the data is perturbed to generate the <code>NULL</code> distribution. The default of 0.5 is reasonable for whole-genome gene expression analysis, but must be reconsidered when the compared lists are deduced from other sources.
<code>min.weight</code>	the minimal weight to be considered.
<code>no.reverse</code>	skip computing scores for reversed second list.

Details

The two lists received as arguments are matched against each other according to the given mapping. The comparison is performed from both ends by default. Permutations of lists are used to generate random scores and compute empirical p-values. The evaluation is also performed for the case the lists should be reversed. From the resulting output, the set of overlapping list identifiers can be extracted using function `getOverlap`.

Value

An object of class `listComparison` is returned. It contains the following list elements:

<code>n</code>	the length of the lists
<code>call</code>	the input parameters
<code>nn</code>	the maximal number of genes corresponding to the alphas and the minimal weight
<code>scores</code>	scores for the straight list comparisons
<code>revScores</code>	scores for the reversed list comparison
<code>pvalues</code>	p-values for the straight list comparison
<code>revPvalues</code>	p-values for the reversed list comparison
<code>overlap</code>	number of overlapping identifiers per rank in straight comparison
<code>revOverlap</code>	number of overlapping identifiers per rank in reversed comparison
<code>randomScores</code>	random scores per weighting parameter
<code>ID.List1</code>	same as input <code>ID.List1</code>
<code>ID.List2</code>	same as input <code>ID.List2</code>

There are print and plot methods for `listComparison` objects. The plot method takes a parameter which to specify whether "overlap" or "density" is to be drawn.

Author(s)

Claudio Lottaz, Stefanie Scheid

References

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

See Also

[OrderedList](#), [getOverlap](#)

Examples

```
### Compare two artificial lists with some overlap
data(OL.data)
list1 <- as.character(OL.data$map$prostate)
list2 <- c(sample(list1[1:500]), sample(list1[501:1000]))
x <- compareLists(list1, list2)
x
getOverlap(x)
```

getOverlap

Extracting the Intersecting IDs From a listComparison Object

Description

This function extracts the intersecting set of list identifiers from an object of class `listComparison` as output of function `compareLists`. The user has to specify the maximum rank to be considered to receive the intersecting set up to this rank.

Usage

```
getOverlap(x, max.rank = NULL, percent = 0.95)
## S3 method for class 'listComparisonOverlap':
plot(x, which="overlap", no.title=FALSE, no.legend=FALSE,
     list.name1="List 1", list.name2="List 2", ...)
```

Arguments

<code>x</code>	An object of class <code>listComparison</code> .
<code>max.rank</code>	The maximum rank to be considered.
<code>percent</code>	The final list of overlapping genes consists of those probes that contribute a certain percentage to the overall similarity score. Default is <code>percent=0.95</code> . To get the full list of genes, set <code>percent=1</code> .
<code>which</code>	select what to draw, either 'overlap' or 'scores'.
<code>no.title</code>	whether to generate a title automatically.
<code>no.legend</code>	whether to generate a legend automatically.
<code>list.name1</code>	A name for the first list provided to <code>compareLists</code> .
<code>list.name2</code>	A name for the second list provided to <code>compareLists</code> .
<code>...</code>	Further arguments passed on to generic <code>plot</code> .

Details

Function `compareLists` returns a list comparison for several choices of `alpha`. The number of genes to be taken into account differs dependent on `alpha`. One might now want to fix the number of genes and receive the resulting set of intersecting list identifiers. To this end, the user chooses a maximum rank to be considered from the values in column 'Genes' of the `listComparison` object. The direction (original or reversed) will internally be set to the direction yielding the higher similarity score.

If `two.sided` was `TRUE`, the first `max.rank` IDs on top of the lists and the first `max.rank` identifiers at the end of the lists are considered. If `two.sided` was `FALSE`, only the `max.rank` top identifiers are evaluated for overlap.

Value

An object of class `listComparisonOverlap` is returned. It contains the following list elements:

<code>n</code>	the length of the lists.
<code>call</code>	the parameters of the input object.

nn	the input max.rank.
score	the observed similarity score.
pvalue	p-values for the observed score.
overlaps	number of overlapping identifiers per rank.
randomScores	random scores for given parameters.
direction	numerical value. Returns '1' if the similarity score is higher for the originally ordered lists and '-1' if the score is higher for the comparison of one original to one reversed list.
intersect	Vector with the sorted overlapping list identifiers, which contribute percent to the overall similarity score.

There are print and plot methods for `listComparisonOverlap` objects. The plot method takes a parameter `which` to specify whether "overlap" or "scores" is to be drawn.

Author(s)

Claudio Lottaz, Stefanie Scheid

References

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

See Also

[OrderedList](#), [compareLists](#)

Examples

```
### Compare two artificial lists with some overlap
data(OL.data)
list1 <- as.character(OL.data$map$prostate)
list2 <- c(sample(list1[1:500]), sample(list1[501:1000]))
x <- compareLists(list1, list2)
x
getOverlap(x)
```

OL.data

Gene Expression and Clinical Information of Two Cancer Studies

Description

The data contains a list with three elements: `breast`, `prostate` and `map`. The first two are expression sets of class `ExpressionSet` taken from the breast cancer study of Huang et al. (2003) and the prostate cancer study of Singh et al. (2002). Both data sets were preprocessed as described in Yang et al. (2006). The data sets serve as illustration for function `prepareData`. Hence the sets contain only a random subsample of the original probes. We further removed unneeded samples from both studies.

The labels of the `breast` expression set were extended with 'B' to create two data sets where the probe IDs differ but can be mapped onto each other. The mapping is stored in the data frame `map`, which consists of the two probe ID vectors.

Usage

```
data(OL.data)
```

References

Huang E, Cheng S, Dressman H, Pittman J, Tsou M, Horng C, Bild A, Iversen E, Liao M, Chen C, West M, Nevins J, and Huang A (2003): Gene expression predictors of breast cancer outcomes, *Lancet* **361**, 1590–1596.

Singh D, Febbo PG, Ross K, Jackson DG, Manola J, Ladd C, Tamayo P, Renshaw AA, D'Amico AV, Richie JP, Lander E, Loda M, Kantoff PW, Golub TR, and Sellers WR (2002): Gene expression correlates of clinical prostate cancer behavior, *Cancer Cell* **1**, 203–209.

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

See Also

[OL.result](#)

OL.result

Three Examples of Class 'OrderedList'

Description

The data set consists of an `OrderedList` object derived by applying function `OrderedList` on the expression sets in `OL.data`. The function calls are given in the example section below.

Usage

```
data(OL.result)
```

References

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

See Also

[OL.data](#), [OrderedList](#)

Examples

```
## Not run:
a <- prepareData(
  list(data=OL.data$breast, name="breast", var="Risk", out=c("high", "low"), pa
  list(data=OL.data$prostate, name="prostate", var="outcome", out=c("Rec", "NF
  mapping=OL.data$map
)
OL.result <- OrderedList(a)

## End(Not run)
```

Description

Function `OrderedList` aims for the *comparison of comparisons*: given two expression studies with one ranked (ordered) list of genes each, we might observe considerable overlap among the top-scoring genes. `OrderedList` quantifies this overlap by computing a weighted similarity score, where the top-ranking genes contribute more to the score than the genes further down the list. The final list of overlapping genes consists of those probes that contribute a certain percentage to the overall similarity score.

Usage

```
OrderedList(eset, B = 1000, test = "z", beta = 1, percent = 0.95,
           verbose = TRUE, alpha=NULL, min.weight=1e-5, empirical=FALSE)
```

Arguments

<code>eset</code>	Expression set containing the two studies of interest. Use <code>prepareData</code> to generate <code>eset</code> .
<code>B</code>	Number of internal sub-samples needed to optimize <code>alpha</code> .
<code>test</code>	String, one of 'fc' (log ratio = log fold change), 't' (t-test with equal variances) or 'z' (t-test with regularized variances). The z-statistic is implemented as described in Efron et al. (2001).
<code>beta</code>	Either 1 or 0.5. In a comparison where the class labels of the studies match, we set <code>beta=1</code> . For example, in each single study the first class relates to bad prognosis while the second class relates to good prognosis. If a matching is not possible, we set <code>beta=0.5</code> . For example, we compare a study with good/bad prognosis classes to a study, in which the classes are two types of cancer tissues.
<code>percent</code>	The final list of overlapping genes consists of those probes that contribute a certain percentage to the overall similarity score. Default is <code>percent=0.95</code> . To get the full list of genes, set <code>percent=1</code> .
<code>verbose</code>	Logical value for message printing.
<code>alpha</code>	A vector of weighting parameters. If set to <code>NULL</code> (the default), parameters are computed such that top 100 to the top 2500 ranks receive weights above <code>min.weight</code> .
<code>min.weight</code>	The minimal weight to be taken into account while computing scores.
<code>empirical</code>	If <code>TRUE</code> , empirical confidence intervals will be computed by randomly permuting the class labels of each study. Otherwise, a hypergeometric distribution is used. Confidence intervals appear when using <code>plot.OrderedList</code> .

Details

In short, the similarity measure is computed as follows: Based on two-sample test statistics like the t-test, genes within each study are ranked from most up-regulated down to most down-regulated. Thus we have one ordered list per study. Now for each rank going both from top (up-regulated

end) and from bottom (down-regulated end) we count the number of overlapping genes. The total overlap A_n for rank n is defined as:

$$A_n = O_n(G_1, G_2) + O_n(f(G_1), f(G_2))$$

where G_1 and G_2 are the two ordered list, $f(G_1)$ and $f(G_2)$ are the two flipped lists with the down-regulated genes on top and O_n is the size of the overlap of its two arguments. A preliminary version of the weighted overlap over all ranks n is then given as:

$$T_\alpha(G_1, G_2) = \sum_n \exp -\alpha n A_n.$$

The final similarity score includes the case that we cannot match the classes in each study exactly and thus do not know whether up-regulation in one list corresponds to up- or down-regulation in the other list. Here parameter β comes into play:

$$S_\alpha(G_1, G_2) = \max \beta T_\alpha(G_1, G_2), (1 - \beta) T_\alpha(G_1, f(G_2)).$$

Parameter β is set by the user but parameter α has to be tuned in a simulation using sub-samples and permutations of the original class labels.

Value

Returns an object of class `OrderedList`, which consists of a list with entries:

<code>n</code>	Total number of genes.
<code>label</code>	The concatenated study labels as provided by <code>eset</code> .
<code>p</code>	The p-value specifying the significance of the similarity.
<code>intersect</code>	Vector with sorted probe IDs of the overlapping genes, which contribute <code>percent</code> to the overall similarity score.
<code>alpha</code>	The optimal regularization parameter <code>alpha</code> .
<code>direction</code>	Numerical value. Returns '1' if the similarity score is higher for the originally ordered lists and '-1' if the score is higher for the comparison of one original to one flipped list. Of special interest if <code>beta=0.5</code> .
<code>scores</code>	Matrix of observed test scores with genes in rows and studies in columns.
<code>sim.scores</code>	List with four elements with output of the resampling with optimal <code>alpha</code> . <code>SIM.observed</code> : The observed similarity score. <code>SIM.alternative</code> : Vector of observed similarity scores simulated using sub-sampling within the distinct classes of each study. <code>SIM.random</code> : Vector of random similarity scores simulated by randomly permuting the class labels of each study. <code>subSample</code> : TRUE to indicate that sub-sampling was used.
<code>pauc</code>	Vector with pAUC-scores for each candidate of the regularization parameter <code>alpha</code> . The maximal pAUC-score defines the optimal <code>alpha</code> . See also <code>plot.OrderedList</code> .
<code>call</code>	List with some of the input parameters.
<code>empirical</code>	List with confidence interval values. Is NULL if <code>empirical=FALSE</code> .

Author(s)

Xinan Yang, Claudio Lottaz, Stefanie Scheid

References

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

Efron B, Tibshirani R, Storey JD, and Tusher V (2001): Empirical Bayes analysis of a microarray experiment, *Journal of the American Statistical Society* **96**, 1151–1160.

See Also

[prepareData](#), [OL.data](#), [OL.result](#), [plot.OrderedList](#), [print.OrderedList](#), [compareLists](#)

Examples

```
### Let's compare the two example studies.
### The first entries of 'out' both relate to bad prognosis.
### Hence the class labels match between the two studies
### and we can use 'OrderedList' with default 'beta=1'.
data(OL.data)
a <- prepareData(
  list(data=OL.data$breast, name="breast", var="Risk", out=c("high", "low"), pa
  list(data=OL.data$prostate, name="prostate", var="outcome", out=c("Rec", "NF
  mapping=OL.data$map
)
## Not run:
OL.result <- OrderedList(a)

## End(Not run)

### The same comparison was done beforehand.
data(OL.result)
OL.result
plot(OL.result)
```

plot.OrderedList *Plotting Function for OrderedList Objects*

Description

The function generates three different plots, which can be selected via argument `which`. With default `which=NULL`, all three figures are plotted into one graphics device.

Usage

```
## S3 method for class 'OrderedList':
plot(x, which = NULL, no.title=FALSE, ...)
```

Arguments

<code>x</code>	Object of class <code>OrderedList</code> .
<code>which</code>	Select one of the three figures described in the details section below.
<code>no.title</code>	logical, whether to skip plotting a title.
<code>...</code>	Additional graphical arguments.

Details

which is one of 'pauc', 'scores' or 'overlap'. If `NULL`, all figures are produced in a row.

Option 'pauc' selects the plot of pAUC-scores, based on which the optimal α is chosen. The pAUC-score measure the separability between the two distributions of observed and expected similarity scores. The similarity scores depend on α and thus α is chosen where the pAUC-scores are maximal. The optimal α is marked by a vertical line.

Figure 'scores' shows kernel density estimates of the two score distributions underlying the pAUC-score for optimal α . The red curve correspondence to simulated observed scores and the black curve to simulated expected scores. The vertical red line denotes the actually observed similarity score. The bottom rugs mark the simulated values. The two distributions got the highest pAUC-score of separability and thus provide the best signal-to-noise separation.

Finally, 'overlap' displays the numbers of overlapping genes in the two gene lists. The overlap size is drawn as a step function over the respective ranks. Top ranks correspond to up-regulated and bottom ranks to down-regulated genes. In addition, the expected overlap and 95% confidence intervals derived from a hypergeometric distribution are plotted. If `empirical=TRUE` in `OrderedList` the confidence intervals were derived empirically from shuffling the data and computing the overlap under the null hypothesis.

Value

No value is returned.

Author(s)

Xinan Yang, Stefanie Scheid

References

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

See Also

[OrderedList](#)

Examples

```
data(OL.result)
plot(OL.result)
```

prepareData

Combining Two Studies into an Expression Set

Description

The function prepares a collection of two expression sets (`ExpressionSet`) and/or Affy batches (`AffyBatch`) to be passed on to the main function `OrderedList`. For each data set, one has to specify the variable in the corresponding phenodata from which the grouping into two distinct classes is done. The data sets are then merged into one `ExpressionSet` together with the re-arranged phenodata. If the studies were done on different platforms but a subset of genes can be mapped from one chip to the other, this information can be provided via the `mapping` argument.

Please note that both data sets have to be *pre-processed* beforehand, either together or independent of each other. In addition, the gene expression values have to be on an *additive scale*, that is logarithmic or log-like scale.

Usage

```
prepareData(eset1, eset2, mapping = NULL)
```

Arguments

eset1	The main inputs are the distinct studies. Each study is stored in a named list, which has five elements: <code>data</code> , <code>name</code> , <code>var</code> , <code>out</code> and <code>paired</code> , see details below.
eset2	Same as <code>eset1</code> for the second data set.
mapping	Data frame containing one named vector for each study. The vectors are comprised of probe IDs that fit to the rownames of the corresponding expression set. For each study, the IDs are ordered identically. For example, the k th row of <code>mapping</code> provides the label of the k th gene in each single study. If all studies were done on the same chip, no mapping is needed (default).

Details

Each study has to be stored in a list with five elements:

<code>data</code>	Object of class <code>ExpressionSet</code> or <code>AffyBatch</code> .
<code>name</code>	Character string with comparison label.
<code>var</code>	Character string with phenodata variable. Based on this variable, the samples for the two-sample testing will be selected.
<code>out</code>	Vector of two character strings with the levels of <code>var</code> that define the two clinical classes. The order of the two levels is important.
<code>paired</code>	Logical - TRUE if samples are paired (e.g. two measurements per patients) or FALSE if all samples are independent.

Value

An object of class `ExpressionSet` containing the joint data sets with appropriate phenodata.

Author(s)

Stefanie Scheid

References

Yang X, Bentink S, Scheid S, and Spang R (2006): Similarities of ordered gene lists, to appear in *Journal of Bioinformatics and Computational Biology*.

See Also

[OL.data](#), [OrderedList](#)

Examples

```
data(OL.data)

### 'map' contains the appropriate mapping between 'breast' and 'prostate' IDs.
### Let's first concatenate two studies.
A <- prepareData(
```


Index

*Topic **datagen**

compareLists, 1
getOverlap, 3
prepareData, 9

*Topic **datasets**

OL.data, 4
OL.result, 5

*Topic **hplot**

plot.OrderedList, 8

*Topic **htest**

OrderedList, 6

*Topic **print**

print.OrderedList, 11

compareLists, 1, 4, 8

getOverlap, 2, 3

OL.data, 4, 5, 8, 10

OL.result, 5, 5, 8

OrderedList, 2, 4, 5, 6, 9–11

plot.listComparison
(compareLists), 1

plot.listComparisonOverlap
(getOverlap), 3

plot.OrderedList, 6, 7, 8, 8

prepareData, 4, 6, 8, 9

print.listComparison
(compareLists), 1

print.listComparisonOverlap
(getOverlap), 3

print.OrderedList, 8, 11