

Package ‘GOexpress’

October 16, 2018

Title Visualise microarray and RNAseq data using gene ontology annotations

Version 1.14.0

Date 2017-04-04

Description The package contains methods to visualise the expression profile of genes from a microarray or RNA-seq experiment, and offers a supervised clustering approach to identify GO terms containing genes with expression levels that best classify two or more predefined groups of samples. Annotations for the genes present in the expression dataset may be obtained from Ensembl through the biomaRt package, if not provided by the user. The default random forest framework is used to evaluate the capacity of each gene to cluster samples according to the factor of interest. Finally, GO terms are scored by averaging the rank (alternatively, score) of their respective gene sets to cluster the samples. P-values may be computed to assess the significance of GO term ranking. Visualisation function include gene expression profile, gene ontology-based heatmaps, and hierarchical clustering of experimental samples using gene expression data.

Depends R (>= 3.4), grid, stats, graphics, Biobase (>= 2.22.0)

Imports biomaRt (>= 2.18.0), stringr (>= 0.6.2), ggplot2 (>= 0.9.0), RColorBrewer (>= 1.0), gplots (>= 2.13.0), randomForest (>= 4.6), RCurl (>= 1.95)

Suggests BiocStyle

License GPL (>= 3)

biocViews Software, GeneExpression, Transcription, DifferentialExpression, GeneSetEnrichment, DataRepresentation, Clustering, TimeCourse, Microarray, Sequencing, RNASeq, Annotation, MultipleComparison, Pathways, GO, Visualization

URL <https://github.com/kevinrue/GOexpress>

LazyData true

git_url <https://git.bioconductor.org/packages/GOexpress>

git_branch RELEASE_3_7

git_last_commit 713a65e

git_last_commit_date 2018-04-30

Date/Publication 2018-10-15

Author Kevin Rue-Albrecht [aut, cre],
 Tharvesh M.L. Ali [ctb],
 Paul A. McGettigan [ctb],
 Belinda Hernandez [ctb],
 David A. Magee [ctb],
 Nicolas C. Nalpas [ctb],
 Andrew Parnell [ctb],
 Stephen V. Gordon [ths],
 David E. MacHugh [ths]

Maintainer Kevin Rue-Albrecht <kevinrue67@gmail.com>

R topics documented:

GOexpress-package	2
AlvMac	4
AlvMac_allgenes	6
AlvMac_allGO	7
AlvMac_GOgenes	8
AlvMac_results	9
AlvMac_results.pVal	10
cluster_GO	11
expression_plot	12
expression_plot_symbol	15
expression_profiles	17
expression_profiles_symbol	19
GO_analyse	22
heatmap_GO	26
hist_scores	28
list_genes	29
microarray2dataset	30
plot_design	31
prefix2dataset	32
pValue_GO	33
quantiles_scores	34
rerank	35
subEset	36
subset_scores	37
table_genes	39

Index **41**

Description

Integrates gene expression data with gene ontology annotations to score and visualise genes and gene ontologies best clustering groups of experimental samples. Supports custom annotations, or alternatively provides an interface to the Ensembl annotations using the biomaRt package. The default scoring approach is based on the random forest framework, while a one-way ANOVA is available as an alternative. GO terms are scored and ranked according to the average rank (alternatively, average power) of all associated genes. P-values can be generated to assess the significance of GO term ranking. The ranked list of GO terms is returned, with tools allowing to visualise the statistics on a gene- and ontology-basis.

Details

Package: GOexpress
Type: Package
Version: 1.7.1
Date: 2016-02-06
License: GPL (>= 3)

This package requires only two input variables

1. An `ExpressionSet` containing `assayData` and `phenoData`. The former should be a gene-by-sample matrix providing gene expression values for each gene in each sample. The latter should be an `AnnotatedDataFrame` from the `Biobase` package providing phenotypic information and grouping factors with two or more levels.
2. The name of the grouping factor to investigate, which must be a valid column name in the `phenoData` slot of the above `ExpressionSet`.

Following analysis, visualisation methods include:

- Histogram and quantiles representations of the scores of GO terms
- Permutation-based P-values to assess the significance of GO term ranking
- Filtering of results on various criteria (e.g. number of genes annotated to GO term)
- Re-ordering of GO terms and gene result tables based on score or rank metric
- Table of statistics for genes annotated to a given GO term
- Hierarchical clustering of samples based on the expression level of genes annotated to a given GO term
- Heatmap of samples and genes based on the expression level of genes annotated to a given GO term
- Expression profile of a gene against one given factor (e.g. Time) while grouping samples on another given factor (e.g. Treatment)
- Univariate analysis of the expression level of a gene in the different groups of each experimental factor.
- Venn diagram of the counts of genes shared between a list of GO terms.

Author(s)

Maintainer: Kevin Rue-Albrecht <kevinrue67@gmail.com>

See Also

Main method for an example usage: [GO_analyse](#).

Packages [Biobase](#), [randomForest](#), [RColorBrewer](#), [VennDiagram](#).

Methods [biomaRt::getBM](#), [ggplot2::ggplot](#), [gplots::heatmap.2](#), [gplots::bluered](#), [gplots::greenred](#), [grid::grid.newpage](#), [grid::grid.layout](#), [stringr::str_extract](#).

Examples

```
# Sample input data available with package:
data(AlvMac)

# Sample output data available with package:
data(AlvMac_results)

# Supported species and microarrays:
data(microarray2dataset)
data(prefix2dataset)
```

AlvMac

Sample data from a RNAseq experiment.

Description

An example ExpressionSet including expression data and phenotypic information about the samples.

The expression data is saved in the assayData slot of the ExpressionSet. It is a gene-by-sample matrix, containing a subset of data from an *in vitro* stimulation of bovine macrophages with different mycobacterial strains. Column names are sample names, and row names are Ensembl gene identifiers of the *Bos taurus* species. Each cell contains the log₂-transformed normalised expression level of each gene in each sample.

The phenotypic information is saved in the phenoData slot of the ExpressionSet. Row names are sample names and columns contain descriptive information about each sample, including experimental factors (e.g. Treatment, Timepoint, Animal).

Usage

```
data(AlvMac)
```

Details

Gene expression was measured in poly-A purified strand-specific RNA libraries using the RNA-Sequencing Illumina(R) HiSeq(R) 2000 platform as paired-end 2 x 90 nucleotide reads. Raw reads from pooled RNA libraries were first deconvoluted according to sample-specific nucleotide barcodes. Read pairs containing adapter sequence in either read mate were discarded, and similarly read pairs of low overall quality in either mate were also discarded. Paired-end reads from each filtered individual library were aligned to the *Bos taurus* reference genome (*B. taurus* UMD3.1.71 genome release) using the STAR aligner software. For each library, raw counts for each gene based on sense strand data were obtained using the featureCounts software from the Subread package. The featureCounts parameters were set to unambiguously assign uniquely aligned paired-end reads in a stranded manner to the exons of genes within the *Bos taurus* reference genome annotation (*B.*

taurus UMD3.1.71 genome annotation). The gene count outputs were further processed using the edgeR Bioconductor package.

The gene expression quantitation pipeline within the edgeR package was customised to: (1) filter out all bovine rRNA genes; (2) filter out genes displaying expression levels below the minimally-set threshold of one count per million [CPM] in at least ten individual libraries (number of biological replicates); (3) calculate normalisation factors for each library using the trimmed mean of M-values method; (4) log₂-transform CPM values based on the normalised library size.

To generate this test data subset, we extracted 100 genes from the original dataset of 12,121 genes. All 7 genes associated with the GO term "GO:0034142" (i.e. "toll-like receptor 4 signaling pathway") present in the original full-size filtered-normalised dataset were kept, all 3 Ensembl gene identifiers annotated to the gene symbol 'RPL36A', and finally another random 90 random genes, making a total of 100 genes measured in 117 samples. Samples include all 10 biological replicates collected at four different time-points, see `data(targets)`. The TLR4 pathway was found in the full dataset as the top-ranking biological pathway discriminating the different mycobacterial infections (unpublished observations).

Value

`assayData` is a matrix of expression levels for 100 genes (rows) measured in 117 samples (columns).

- `rownames` are Ensembl gene identifiers of the *Bos taurus* species.
- `colnames` are samples identifiers.

`phenoData` is a data frame with 117 samples and 7 descriptive fields (e.g. experimental factors) in the columns listed below:

- `rownames` are unique identifiers. Here, sample names.
- `File` contains local filenames where the RNAseq counts were obtained from.
- `Sample` contains individual sample name.
- `Animal` contains the unique identifier of the animal corresponding to the biological replicate, stored as a factor.
- `Treatment` contains the infection status of the sample, stored as a factor (CN: Control, MB: *M. bovis*, TB: *M. tuberculosis*)
- `Time` contains the time of measurement in hours post-infection, stored as a factor.
- `Group` contains a combination of the Treatment and Time factors above, stored as a factor itself.
- `Timepoint` contains the time of measurement, stored as a numeric value. This field is useful to use on the X-axis of expression plots. See function `expression_plot()`.

Source

Publication in review process.

Examples

```
# Load the data
data(AlvMac)

# Structure of the data
str(AlvMac)

# Dimensions (rows, columns) of the data
```

```
dim(AlvMac)

# Subset of first 5 features and 5 samples
AlvMac[1:5, 1:5]

# Phenotypic information
pData(AlvMac)

# Phenotypic information about factor "Group"
AlvMac$Group

# Conversion of a factor to a character vector
as.character(AlvMac$Group)

# Number of samples (rows) and annotations (columns)
dim(pData(AlvMac))
```

AlvMac_allgenes	<i>Example of custom gene feature annotations.</i>
-----------------	--

Description

An example data.frame providing the identifier, name and description corresponding to Ensembl gene identifiers present in the AlvMac example ExpressionSet.

Usage

```
data("AlvMac_allgenes")
```

Details

This data-frame includes only Ensembl gene identifiers present in the AlvMac example ExpressionSet. See the help page of the GO_analyse function for an example usage.

Value

A data frame detailing information about the 100 gene features present in the AlvMac example ExpressionSet:

- gene_id are Ensembl gene identifiers.
- external_gene_name contains the corresponding gene name.
- description contains the corresponding description.

Source

These annotations were obtained from the Ensembl BioMart server using the biomaRt package to access the Ensembl release 75 <http://feb2014.archive.ensembl.org>.

Examples

```
# Load the data
data(AlvMac_allgenes)

# Structure of the data
str(AlvMac_allgenes)

# First few rows
head(AlvMac_allgenes)
```

AlvMac_allGO

Example of custom gene ontology annotations.

Description

An example data.frame providing the identifier, name and namespace corresponding to gene ontology identifiers, compatible with the AlvMac example ExpressionSet.

Usage

```
data("AlvMac_allGO")
```

Details

This data-frame includes all gene ontologies present in the btaurus_gene_ensembl dataset of the Ensembl BioMart server, including those associated with no gene identifier in the AlvMac example ExpressionSet.

See the help page of the GO_analyse function for an example usage.

Value

A data frame detailing information about 13,302 gene ontologies:

- go_id are gene ontology identifiers.
- name_1006 contains the corresponding gene ontology name.
- namespace_1003 contains the corresponding gene ontology namespace (i.e. "biological_process", "molecular_function", or "cellular_component")

Source

These annotations were obtained from the Ensembl BioMart server using the biomaRt package to access the Ensembl release 75 <http://feb2014.archive.ensembl.org>.

Examples

```
# Load the data
data(AlvMac_allGO)

# Structure of the data
str(AlvMac_allGO)

# First few rows
head(AlvMac_allGO)
```

AlvMac_GOgenes	<i>Example of custom mapping between gene ontology identifiers and gene features.</i>
----------------	---

Description

An example `data.frame` associating Ensembl gene identifiers to gene ontology identifiers, compatible with the AlvMac example `ExpressionSet`.

Usage

```
data("AlvMac_GOgenes")
```

Details

This data-frame includes all annotations between Ensembl gene identifiers and gene ontology identifiers present in the `btaurus_gene_ensembl` dataset of the Ensembl BioMart server, including gene identifiers absent from the AlvMac example `ExpressionSet`.

Importantly, gene identifiers present in these annotations but absent from the `ExpressionSet` will be given a score of 0 (minimum valid score; indicates no power to discriminate the predefined groups of sample) and a rank equal to the number of genes present in the entire dataset plus one (worst rank, while preserving discrete continuity of the ranking). This is helpful where features considered uninformative were filtered out of the `ExpressionSet`.

See the help page of the `GO_analyse` function for an example usage.

Value

A data frame with 191,614 associations between the following 2 variables.

- `gene_id` are Ensembl gene identifiers of the *Bos taurus* species.
- `go_id` are gene ontology identifiers.

Source

These annotations were obtained from the Ensembl BioMart server using the `biomaRt` package to access the Ensembl release 75 <http://feb2014.archive.ensembl.org>.

Examples

```
# Load the data
data(AlvMac_GOgenes)

# Structure of the data
str(AlvMac_GOgenes)

# First few rows
head(AlvMac_GOgenes)
```

AlvMac_results	<i>Sample output from the GO_analyse() function on an RNAseq experiment.</i>
----------------	--

Description

This variable may be used to test the filtering and visualisation methods implemented in the package. It contains the output of the command `AlvMac_results = GO_analyse(eSet=AlvMac, f="Treatment")` applied to the toy input data AlvMac.

Usage

```
data(AlvMac_results)
```

Value

A list of 9 slots summarising the input and results of the analysis:

- `GO` contains a table ranking all GO terms related to genes in the expression dataset based on the average ability of their related genes to cluster the samples according to the predefined grouping factor.
- `mapping` contains the table mapping genes present in the dataset to GO terms.
- `genes` contains a table ranking all genes present in the expression dataset based on their ability to cluster the samples according to the predefined grouping factor (see 'factor' below).
- `factor` contains the grouping factor analysed.
- `method` contains the statistical framework used.
- `subset` contains the filters used to select a subset of samples from the original ExpressionSet for analysis.
- `rank.by` contains the metric used to rank the scoring tables.
- `n.tree` contains number of trees built during the randomForest analysis.
- `m.try` contains the number of features randomly sampled as candidates at each split in each tree built during the randomForest analysis.

Warning

Running the above command again, you might obtain slightly different scores and ranks due to the stochastic process of sampling used by the random forest algorithm. However, the ranking metric was found to be robust and stable across run, given adequate number of trees and predictor variables sampled.

To produce reproducible results, use the `set.seed()` function prior to running any randomising or sampling function.

Source

Source data are part of a publication in review.

Examples

```
data(AlvMac_results)
str(AlvMac_results)
head(AlvMac_results$GO, n=20)
head(AlvMac_results$GO$genes, n=20)
```

AlvMac_results.pVal *Sample output from the pValue_GO() function on an RNAseq experiment.*

Description

This variable may be used to test the filtering and visualisation methods implemented in the package.

It contains the output of successively applying the commands `AlvMac_results = GO_analyse(eSet=AlvMac, f="Trea` and `AlvMac_results.pVal = pValue_GO(AlvMac_results, N=100, ranked.by=result$rank.by, rank.by='P')` to the toy input data `AlvMac`.

Usage

```
data("AlvMac_results.pVal")
```

Value

A list of 9 slots summarising the input and results of the analysis:

- `GO` contains a table ranking all GO terms related to genes in the expression dataset based on the average ability of their related genes to cluster the samples according to the predefined grouping factor.
- `mapping` contains the table mapping genes present in the dataset to GO terms.
- `genes` contains a table ranking all genes present in the expression dataset based on their ability to cluster the samples according to the predefined grouping factor (see 'factor' below).
- `factor` contains the grouping factor analysed.
- `method` contains the statistical framework used.
- `subset` contains the filters used to select a subset of samples from the original `ExpressionSet` for analysis.
- `rank.by` contains the metric used to rank the scoring tables.
- `n.tree` contains number of trees built during the `randomForest` analysis.
- `m.try` contains the number of features randomly sampled as candidates at each split in each tree built during the `randomForest` analysis.
- `p.iterations` contains the number of permutations performed to compute the P-value in the GO slot.

Warning

Running the above command again, you might obtain slightly different scores and ranks due to the stochastic process of sampling used by the random forest algorithm. However, the ranking metric was found to be robust and stable across run, given adequate number of trees and predictor variables sampled.

To produce reproducible results, use the `set.seed()` function prior to running any randomising or sampling function.

Source

Source data are part of a publication in review.

Examples

```
data(AlvMac_results.pVal)
str(AlvMac_results.pVal)
head(AlvMac_results.pVal, n=20)
```

cluster_GO	<i>Generates a hierarchical clustering of the samples</i>
------------	---

Description

Clusters the samples using only the expression levels of genes associated with a given go_id.

Usage

```
cluster_GO(
  go_id, result, eSet, f=result$factor, subset=NULL,
  method_dist="euclidean", method_hclust="average", cex=0.8,
  main=paste(go_id, result$GO[result$GO$go_id == go_id, "name_1006"]),
  xlab="Distance", cex.main=1, main.Lsplit=NULL, ...)
```

Arguments

go_id	A Gene Ontology (GO) identifier.
result	The output of GO_analyse() or a subset of it obtained from subset_scores().
eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the AssayData slot, and a phenotypic information data-frame in the phenodata slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
f	The grouping factor in phenodata to label the samples by.
subset	A named list to subset eSet. Names must be column names existing in colnames(pData(eSet)). Values must be vectors of values existing in the corresponding column of pData(eSet).
method_dist	The method used to calculate distance between samples. See the dist() method from package stats.
method_hclust	The method used to cluster samples. See the hclust() method from package stats.
cex	A numeric value defining the character expansion of text in the plot.
main	A character string for the main title of the plot.
xlab	A label for the x axis, defaults to "Distance".
cex.main	Scaling factor of the main title font size. Default is 1. We suggest to use it in combination with the argument main.Lsplit for GO terms with long names.

`main.Lsplit` Number of characters after which a new-line character will be inserted in the main title. If this would occur within a word, the new-line character will be inserted before this word. Default is NULL, leaving the title on a single line.

... Additional parameters passed on to `dist()`, `hclust()` and `plot()`.

Value

Returns the output of the `plot()` function.

Author(s)

Kevin Rue-Albrecht

See Also

Method [GO_analyse](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# Hierarchical clustering of samples based on the same GO term
cluster_GO(
  go_id="GO:0034142", result=AlvMac_results, eSet=AlvMac, cex=0.7
)

# Re-label sample by another factor
cluster_GO(
  go_id="GO:0034142", result=AlvMac_results, eSet=AlvMac, cex=0.7, f="Group"
)
```

expression_plot

Plots the expression profile of a gene by levels of a factor

Description

This function will plot the expression profile of a gene across a valid X-axis variable in the phenodata while representing the mean and confidence interval of groups of samples defined by levels of another valid grouping factor in the phenodata.

Usage

```
expression_plot(
  gene_id, result, eSet, x_var, f=result$factor, subset=NULL,
  xlab=x_var, ylab="log2(cpm)", ylim=range(exprs(eSet)),
  col.palette="Accent",
  col=brewer.pal(n=length(levels(pData(eSet)[,f])), name=col.palette),
  level=0.95, title=NULL, title.size=2, axis.title.size=20,
  axis.text.size=15, axis.text.angle=0,
  legend.title.size=20, legend.text.size=15, legend.key.size=30)
```

Arguments

gene_id	An gene or probeset identifier present in rownames(expr_data).
result	An output of the GO_analyse() or subset_scores() function.
eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the AssayData slot, and a phenotypic information data-frame in the phenodata slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
x_var	A column name in phenodata to plot on the X-axis. If representing time on the X-axis, users should store the time-points as numeric values in the AnnotatedDataFrame to adequately space the time-points.
f	A column name in phenodata to group the samples when representing mean and confidence interval. The factor specified in the initial GO_analyse() call is used by default. Unexpected grouping factors of samples can reveal interesting trends (e.g. "Animal", "Tissue", "CellType", ...).
subset	A named list to subset eSet. Names must be column names existing in colnames(pData(eSet)). Values must be vectors of values existing in the corresponding column of pData(eSet).
xlab	Title of the X-axis. Default is the value of x_var.
ylab	Title of the Y-axis. Default is "log2(cpm)".
ylim	Numeric vector of length 2 specifying the lower and upper bounds of the Y axis. Default is scaled to the full range of expression values in the expression dataset, to ease comparison of different genes. If set to NULL, the axis will be scaled to fit the plotted data only.
col.palette	A valid RColorBrewer palette name to fetch the colormap from. Default is palette "Accent".
col	A vector of color names or codes. The number of colors provided must match the number of levels of the grouping factor. If specified, overrides argument col.palette.
level	The confidence interval level to visualise around the mean of each group. Default is 0.95.
title	Changes the plot title. Default is a combination of the gene id and the associated gene.
title.size	Changes the font size of the title. Default is 2.
axis.title.size	Changes the font size of the axes title. Default is 20.
axis.text.size	Changes the font size of the axes text labels. Default is 15.
axis.text.angle	Changes the angle of the X axis text labels. Default is 0 (horizontal).
legend.title.size	Changes the font size of the legend title. Default is 20.
legend.text.size	Changes the font size of the legend text labels. Default is 15.
legend.key.size	Changes the size of the legend keys (in points). Default is 30.

Value

The ggplot object.

Warning

Common issues:

- It may not be possible to produce plots where the combination of X-axis variable and grouping factor leaves too few replicates to compute a confidence interval for each X value. This is a limitation imposed by the `ggplot2` package to produce proper statistics and confidence intervals. In such cases, it may be preferable to use the `expression_profiles()` method.

Author(s)

Kevin Rue-Albrecht

References

- [ggplot2](#) package.

See Also

Package [Biobase](#), methods [expression_plot_symbol](#), [GO_analyse](#) and [ggplot](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# Expression by gene identifier (TNIP3)
expression_plot(
  gene_id="ENSBTAG00000047107",
  result=AlvMac_results, eSet=AlvMac, x_var="Timepoint"
)

# Same gene, plotted by animal and grouped by treatment (merging time points)
expression_plot(
  gene_id="ENSBTAG00000047107",
  result=AlvMac_results, eSet=AlvMac, x_var="Animal",
  f="Treatment")

# Same gene, plotted by animal and grouped by time-point (merging treatments)
expression_plot(
  gene_id="ENSBTAG00000047107",
  result=AlvMac_results, eSet=AlvMac, x_var="Animal",
  f="Time")
```

 expression_plot_symbol

Plots the expression profile of a gene by levels of a factor

Description

This function will plot the expression profile of a gene across a valid X-axis variable from the AnnotatedDataFrame while representing the mean and confidence interval of groups of samples defined by levels of a valid grouping factor from the AnnotatedDataFrame.

In the case of a gene name represented by multiple gene or probeset identifiers in the dataset, a lattice of plots will be produced. Each of the plots in this lattice can subsequently be plotted separately using its associated index.

Usage

```
expression_plot_symbol(
  gene_symbol, result, eSet, x_var, f=result$factor, subset=NULL,
  index=0, xlab=x_var, ylab="log2(cpm)", ylim=range(exprs(eSet)),
  col.palette="Accent",
  col=brewer.pal(n=length(levels(pData(eSet)[,f])), name=col.palette),
  level=0.95, titles=c(), title.size=2, axis.title.size=20,
  axis.text.size=15, axis.text.angle=0,
  legend.title.size=20, legend.text.size=20, legend.key.size=30)
```

Arguments

gene_symbol	A gene name present in rownames(AlvMac_results\$genes).
result	An output of the GO_analyse() or subset_scores() function.
eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the AssayData slot, and a phenotypic information data-frame in the phenodata slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
x_var	A column name in phenodata to plot on the X-axis. If representing time on the X-axis, users should store the time-points as numeric values in the AnnotatedDataFrame to adequately space the time-points.
f	A column name in phenodata to group the samples when representing mean and confidence interval. The factor specified in the initial GO_analyse() call is used by default. Unexpected grouping factors of samples can reveal interesting trends (e.g. "Animal", "Tissue", "CellType", ...).
subset	A named list to subset eSet. Names must be column names existing in colnames(pData(eSet)). Values must be vectors of values existing in the corresponding column of pData(eSet).
index	In the case where multiple gene or probeset identifiers are associated with the gene name, index=2 will plot the expression profile of the second feature identifier alone, for instance.
xlab	Title of the X-axis. Default is the value of x_var.

<code>ylab</code>	Title of the Y-axis. Default is "log2(cpm)".
<code>ylim</code>	Numeric vector of length 2 specifying the lower and upper bounds of the Y axis. Default is scaled to the full range of expression values in the expression dataset, to ease comparison of different genes. If set to NULL, the axis will be scaled to fit the plotted data only.
<code>col.palette</code>	A valid RColorBrewer palette name to fetch the colormap from.
<code>col</code>	A vector of color names or codes. The number of colors provided must match the number of levels of the grouping factor. Default to a palette with an adequate set of colors. If specified, overrides argument <code>col.palette</code> .
<code>level</code>	The confidence interval level to visualise around the mean of each group. Default is 0.95.
<code>titles</code>	Character vector providing as many titles as there are plots to replace the default titles. Default is a combination of the gene id and the associated gene.
<code>title.size</code>	Changes the font size of the title. Default is 2.
<code>axis.title.size</code>	Changes the font size of the axes title. Default is 20.
<code>axis.text.size</code>	Changes the font size of the axes text labels. Default is 15.
<code>axis.text.angle</code>	Changes the angle of the X axis text labels. Default is 0 (horizontal).
<code>legend.title.size</code>	Changes the font size of the legend title. Default is 20.
<code>legend.text.size</code>	Changes the font size of the legend text labels. Default is 15.
<code>legend.key.size</code>	Changes the size of the legend keys (in points). Default is 30.

Value

The `ggplot` object, or the vector of feature identifiers if multiple features are annotated to the gene symbol.

Warning

Common issues:

- It may not be possible to produce plots where the combination of X-axis variable and grouping factor leaves too few replicates to compute a confidence interval for each X value. This is a limitation imposed by the `ggplot2` package to produce proper statistics and confidence intervals. In such cases, it may be preferable to use the `expression_profiles_symbol()` method.

Author(s)

Kevin Rue-Albrecht

References

- [ggplot2](#) package.

See Also

Package [Biobase](#), methods [expression_plot](#), [GO_analyse](#) and [ggplot](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# Expression by gene identifier (TNIP3)
expression_plot_symbol(
  gene_symbol="PIK3AP1",
  result=AlvMac_results, eSet=AlvMac, x_var="Timepoint"
)

# Same gene, plotted by animal and grouped by treatment (merging time points)
expression_plot_symbol(
  gene_symbol="PIK3AP1",
  result=AlvMac_results, eSet=AlvMac, x_var="Animal",
  f="Treatment"
)

# Same gene, plotted by animal and grouped by time-point (merging treatments)
expression_plot_symbol(
  gene_symbol="PIK3AP1",
  result=AlvMac_results, eSet=AlvMac, x_var="Animal",
  f="Time")
```

expression_profiles *Plots the individual expression profile of a gene in samples series*

Description

This function will plot the expression profile of a gene in individual samples series across a valid X-axis variable in the phenodata, while colouring sample groups according to another variable in the phenodata, and using different line types according to yet another (or the same) variable in the phenodata.

Usage

```
expression_profiles(
  gene_id, result, eSet, x_var, seriesF, subset=NULL,
  colourF=result$factor, linetypeF=colourF, line.size=1.5,
  xlab=x_var, ylab="log2(cpm)", ylim=range(exprs(eSet)),
  col.palette="Accent",
  col=brewer.pal(n=length(levels(pData(eSet)[,colourF])),
                 name=col.palette),
  lty=1:length(levels(pData(eSet)[,linetypeF])),
  title=NULL, title.size=2, axis.title.size=20,
  axis.text.size=15, axis.text.angle=0,
  legend.title.size=20, legend.text.size=15,
  legend.key.size=30)
```

Arguments

gene_id	An gene or probeset identifier present in rownames(expr_data).
result	An output of the GO_analyse() or subset_scores() function.

eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the AssayData slot, and a phenotypic information data-frame in the phenodata slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
x_var	A column name in phenodata to plot on the X-axis. If representing time on the X-axis, users should store the time-points as numeric values in the AnnotatedDataFrame to adequately space the time-points.
seriesF	A column name in phenodata to track samples of an identical series across the X-axis. The combination of seriesF and x_var should uniquely identify each sample in eSet.
subset	A named list to subset eSet. Names must be column names existing in colnames(pData(eSet)). Values must be vectors of values existing in the corresponding column of pData(eSet).
colourF	A column name in phenodata to colour the corresponding groups of samples. The factor specified in the initial GO_analyse() call is used by default.
linetypeF	A column name in phenodata to assign a different line type to the corresponding groups of samples. Default mirrors colourF.
line.size	The width of the plotted lines. Default is 1.5.
xlab	Title of the X-axis. Default is the value of x_var.
ylab	Title of the Y-axis. Default is "log2(cpm)".
ylim	Numeric vector of length 2 specifying the lower and upper bounds of the Y axis. Default is scaled to the full range of expression values in the expression dataset, to ease comparison of different genes. If set to NULL, the axis will be scaled to fit the plotted data only.
col.palette	A valid RColorBrewer palette name to fetch the colormap from. Default is palette "Accent".
col	A vector of color names or codes. The number of colors provided must match the number of levels of the factor colourF. If specified, overrides argument col.palette. Default colours are obtained from the col.palette.
lty	A vector of numeric values corresponding to line types. The number of line types provided must match the number of levels of the factor linetypeF. Default line types are 1:length(levels(pData(eSet)[, linetypeF])).
title	Changes the plot title. Default is a combination of the gene id and the associated gene.
title.size	Changes the font size of the title. Default is 2.
axis.title.size	Changes the font size of the axes title. Default is 20.
axis.text.size	Changes the font size of the axes text labels. Default is 15.
axis.text.angle	Changes the angle of the X axis text labels. Default is 0 (horizontal).
legend.title.size	Changes the font size of the legend title. Default is 20.
legend.text.size	Changes the font size of the legend text labels. Default is 15.
legend.key.size	Changes the size of the legend keys (in points). Default is 30.

Details

In order to track and visualise individual sample series, each sample in the ExpressionSet should be associated with a unique combination of seriesF and x_var. This may require the generation of a new factor in the phenodata, combining all experimental factors except that plotted on the X-axis. See below for an example on the training dataset.

Value

The ggplot object.

Author(s)

Kevin Rue-Albrecht

References

- [ggplot2](#) package.

See Also

Package [Biobase](#), methods [expression_plot_symbol](#), [expression_plot_symbol](#) and [ggplot](#).

Examples

```
# load the sample output data
data(AlvMac_results)

AlvMac$Series <- paste(AlvMac$Animal, AlvMac$Treatment, sep="_")

expression_profiles(
  gene_id = "ENSBTAG00000047107", result = AlvMac_results,
  eSet=AlvMac, x_var = "Timepoint", seriesF="Series",
  legend.title.size=10, legend.text.size=10, legend.key.size=15)

expression_profiles(
  gene_id = "ENSBTAG00000047107", result = AlvMac_results,
  eSet=AlvMac, x_var = "Timepoint", seriesF="Series",
  linetypeF="Animal",
  legend.title.size=10, legend.text.size=10, legend.key.size=15)
```

expression_profiles_symbol

Plots the individual expression profile of a gene in samples series

Description

This function will plot the expression profile of a gene in individual samples series across a valid X-axis variable in the phenodata, while colouring sample groups according to another variable in the phenodata, and using different line types according to yet another (or the same) variable in the phenodata.

Usage

```
expression_profiles_symbol(
  gene_symbol, result, eSet, x_var, seriesF, subset=NULL,
  colourF=result$factor, linetypeF=colourF, line.size=1.5,
  index=0, xlab=x_var, ylab="log2(cpm)", ylim=range(exprs(eSet)),
  col.palette="Accent",
  col=brewer.pal(n=length(levels(pData(eSet)[,colourF])),
    name=col.palette),
  lty=1:length(levels(pData(eSet)[,linetypeF])),
  titles=c(), title.size=2, axis.title.size=20,
  axis.text.size=15, axis.text.angle=0,
  legend.title.size=20, legend.text.size=15,
  legend.key.size=30)
```

Arguments

gene_symbol	A gene name present in rownames(AlvMac_results\$genes).
result	An output of the GO_analyse() or subset_scores() function.
eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the AssayData slot, and a phenotypic information data-frame in the phenodata slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
x_var	A column name in phenodata to plot on the X-axis. If representing time on the X-axis, users should store the time-points as numeric values in the AnnotatedDataFrame to adequately space the time-points.
seriesF	A column name in phenodata to track samples of an identical series across the X-axis. The combination of seriesF and x_var should uniquely identify each sample in eSet.
subset	A named list to subset eSet. Names must be column names existing in colnames(pData(eSet)). Values must be vectors of values existing in the corresponding column of pData(eSet).
colourF	A column name in phenodata to colour the corresponding groups of samples. The factor specified in the initial GO_analyse() call is used by default.
linetypeF	A column name in phenodata to assign a different line type to the corresponding groups of samples. Default mirrors colourF.
line.size	The width of the plotted lines. Default is 1.5.
index	In the case where multiple gene or probeset identifiers are associated with the gene name, index=2 will plot the expression profile of the second feature identifier alone, for instance.
xlab	Title of the X-axis. Default is the value of x_var.
ylab	Title of the Y-axis. Default is "log2(cpm)".
ylim	Numeric vector of length 2 specifying the lower and upper bounds of the Y axis. Default is scaled to the full range of expression values in the expression dataset, to ease comparison of different genes. If set to NULL, the axis will be scaled to fit the plotted data only.
col.palette	A valid RColorBrewer palette name to fetch the colormap from. Default is palette "Accent".

<code>col</code>	A vector of color names or codes. The number of colors provided must match the number of levels of the factor <code>colourF</code> . If specified, overrides argument <code>col.palette</code> . Default colours are obtained from the <code>col.palette</code> .
<code>lty</code>	A vector of numeric values corresponding to line types. The number of line types provided must match the number of levels of the factor <code>linetypeF</code> . Default line types are <code>1:length(levels(pData(eSet)[, linetypeF]))</code> .
<code>titles</code>	Character vector providing as many titles as there are plots to replace the default titles. Default is a combination of the gene id and the associated gene.
<code>title.size</code>	Changes the font size of the title. Default is 2.
<code>axis.title.size</code>	Changes the font size of the axes title. Default is 20.
<code>axis.text.size</code>	Changes the font size of the axes text labels. Default is 15.
<code>axis.text.angle</code>	Changes the angle of the X axis text labels. Default is 0 (horizontal).
<code>legend.title.size</code>	Changes the font size of the legend title. Default is 20.
<code>legend.text.size</code>	Changes the font size of the legend text labels. Default is 15.
<code>legend.key.size</code>	Changes the size of the legend keys (in points). Default is 30.

Details

In order to track and visualise individual sample series, each sample in the `ExpressionSet` should be associated with a unique combination of `seriesF` and `x_var`. This may require the generation of a new factor in the phenodata, combining all experimental factors except that plotted on the X-axis. See below for an example on the training dataset.

Value

The `ggplot` object, or the vector of feature identifiers if multiple features are annotated to the gene symbol.

Author(s)

Kevin Rue-Albrecht

References

- [ggplot2](#) package.

See Also

Package [Biobase](#), methods [expression_plot_symbol](#), [expression_plot_symbol](#) and [ggplot](#).

Examples

```
# load the sample output data
data(AlvMac_results)

AlvMac$Series <- paste(AlvMac$Animal, AlvMac$Treatment, sep="_")

expression_profiles_symbol(
  gene_symbol="TNIP3", result = AlvMac_results,
  eSet=AlvMac, seriesF="Series", x_var = "Timepoint", linetypeF="Animal",
  line.size=1.5, title.size=1.5, legend.title.size=10, legend.text.size=10,
  legend.key.size=30)

expression_profiles_symbol(
  gene_symbol="TNIP3", result = AlvMac_results,
  eSet=AlvMac, seriesF="Series", x_var = "Timepoint", linetypeF="Animal",
  line.size=1.5, lty=rep(1,10), title.size=1.5, legend.title.size=10,
  legend.text.size=10, legend.key.size=30)
```

GO_analyse

Identifies gene ontologies clustering samples according to predefined factor.

Description

Combines gene expression data with Gene Ontology (GO) annotations to rank and visualise genes and GO terms enriched for genes best clustering predefined groups of samples based on gene expression levels.

This methods semi-automatically retrieves the latest information from Ensembl using the `biomaRt` package, except if custom GO annotations are provided. Custom GO annotations have two main benefits: firstly they allow the analysis of species not supported in the Ensembl BioMart server, and secondly they save time skipping calls to the Ensembl BioMart server for species that are supported. The latter also presents the possibility of using an older release of the Ensembl annotations, for example.

Using default settings, a random forest analysis is performed to evaluate the ability of each gene to cluster samples according to a predefined grouping factor (one-way ANOVA available as an alternative). Each GO term is scored and ranked according to the average rank (alternatively, average power) of all associated genes to cluster the samples according to the factor. The ranked list of GO terms is returned, with tools allowing to visualise the statistics on a gene- and ontology-basis.

Usage

```
GO_analyse(
  eSet, f, subset=NULL, biomaRt_name = "ENSEMBL_MART_ENSEMBL",
  biomaRt_dataset="", microarray="",
  method="randomForest", rank.by="rank", do.trace=100, ntree=1000,
  mtry=ceiling(2*sqrt(nrow(eSet))), GO_genes=NULL, all_GO=NULL,
  all_genes=NULL, FUN.GO=mean, ...)
```

Arguments

eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the AssayData slot, and a phenotypic information data-frame in the phenodata slot. In the expression matrix, row names are identifiers of expressed features, and column names are identifiers of the individual samples. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for statistical tests and visualisation methods.
f	A column name in phenodata used as the grouping factor for the analysis.
subset	A named list to subset eSet for the analysis. Names must be column names existing in colnames(pData(eSet)). Values must be vectors of values existing in the corresponding column of pData(eSet). The original ExpressionSet will be left unchanged.
biomart_name	The Ensembl BioMart database to which connect to.
biomart_dataset	The Ensembl BioMart dataset identifier corresponding to the species studied. If not specified and no custom annotations were provided, the method will attempt to automatically identify the adequate dataset from the first feature identifier in the dataset. Use data(prefix2dataset) to access a table listing valid choices.
microarray	The identifier in the Ensembl BioMart corresponding to the microarray platform used. If not specified and no custom annotations were provided, the method will attempt to automatically identify the platform used from the first feature identifier in the dataset. Use data(microarray2dataset) to access a table listing valid choices.
method	The statistical framework to score genes and gene ontologies. Either "randomForest" or "rf" to use the random forest algorithm, or alternatively either of "anova" or "a" to use the one-way ANOVA model. Default is "randomForest".
rank.by	Either of "rank" or "score" to chose the metric used to order the gene and GO term result tables. Default to 'rank'.
do.trace	Only used if method="randomForest". If set to TRUE, gives a more verbose output as randomForest is run. If set to some integer, then running output is printed for every do.trace trees. Default is 100.
ntree	Only used if method="randomForest". Number of trees to grow. This should be set to a number large enough to ensure that every input row gets predicted at least a few times
mtry	Only used if method="randomForest". Number of features randomly sampled as candidates at each split. Default value is 2*sqrt(gene_count) which is approximately 220 genes for a dataset of 12,000 genes.
GO_genes	Custom annotations associating features present in the expression dataset to gene ontology identifiers. This must be provided as a data-frame of two columns, named gene_id and go_id. If provided, no call to the Ensembl BioMart server will be done, and arguments all_GO and all_genes should be provided as well, to enable all downstream features of GOexpress. An example is provided in AlvMac_GOgenes.
all_GO	Custom annotations used to annotate each GO identifier present in GO_genes with the ontology name (e.g. "apoptotic process") and namespace (i.e. "biological_process", "molecular_function", or "cellular_component"). This must be provided as a data-frame containing at least one column named go_id, and

preferably two more columns named `name_1006` and `namespace_1003` for consistency with the Ensembl BioMart. Supported alternative column headers are `name` and `namespace`. Respectively, `name` should be used to provide the description of the GO term, and `namespace` should contain one of "biological_process", "molecular_function" and "cellular_component". `name` is used to generate the title of ontology-based figured, and `namespace` is important to enable subsequent filtering of results by their corresponding value. An example is provided in `data(AlvMac_allGO)`.

<code>all_genes</code>	Custom annotations used to annotate each feature identifier in the expression dataset with the gene name or symbol (e.g. "TNF"), and an optional description. This must be provided as a data-frame containing at least a column named <code>gene_id</code> and preferably two more columns named <code>external_gene_name</code> and <code>description</code> for consistency with the Ensembl BioMart. A supported alternative header is <code>name</code> . While <code>external_gene_name</code> is important to enable subsequent visualisation of results by gene symbol, <code>description</code> is only displayed for readability of result tables. An example is provided in <code>data(AlvMac_allgenes)</code> .
<code>FUN.GO</code>	Function to summarise the score and rank of all feature associated with each gene ontology. Default is mean function. If using "lambda-like" (anonymous) functions, these must take a list of numeric values as an input, and return a single numeric value as an output.
<code>...</code>	Additional arguments passed on to the <code>randomForest</code> and <code>useMart</code> methods, if applicable.

Details

The default scoring functions strongly favor GO terms associated with fewer genes at the top of the ranking. This bias may actually be seen as a valuable feature which enables the user to browse through GO terms of increasing "granularity", i.e. associated with increasingly large sets of genes, although consequently being increasingly vague and blurry (e.g. "protein binding" molecular function associated with over 6,000 genes).

It is suggested to use the `subset_scores()` function to subsequently filter out GO terms with fewer than 5+ genes associated with it. Indeed, those GO terms are more sensitive to outlier genes as they were scored on the average of a handful of genes.

Additionally, the `pValue_GO` function may be used to generate a permutation-based P-value indicating the chance of seeing each GO term reaching an equal or higher rank – or score – by chance.

Value

A list containing the results of the analysis. Some elements are specific to the output of each analysis method.

Core elements:

<code>GO</code>	A table ranking all GO terms related to genes in the expression dataset based on the average ability of their related genes to cluster the samples according to the predefined grouping factor.
<code>mapping</code>	The table mapping genes present in the dataset to GO terms.
<code>genes</code>	A table ranking all genes present in the expression dataset based on their ability to cluster the samples according to the predefined grouping factor.
<code>factor</code>	The predefined grouping factor.
<code>method</code>	The statistical framework used.

subset	The filters used to run the analysis only on a subset of the samples. NULL if no filter was applied.
rank.by	The metric used to rank order the genes and gene ontologies.
FUN.GO	The function used to summarise the score and rank of all gene features associated with each gene ontology.

Random Forest additional elements:

n tree	Number of trees grown.
m try	Number of variables randomly sampled as candidates at each split.

One-way ANOVA does not have additional arguments.

Warning

Make sure that the factor *f* is an actual factor in the R language meaning. This is important for the underlying statistical framework to identify the groups of samples defined by their level of this factor.

If the column defining the factor (e.g. "Treatment") in `phenodata` is not an R factor, use `pData(targets)$Treatment =` to convert the character values into an actual R factor with appropriate levels.

Author(s)

Kevin Rue-Albrecht

See Also

Methods [subset_scores](#), [pValue_GO](#), [getBM](#), [randomForest](#), and [oneway.test](#).

Examples

```
# Load example data subset
data(AlvMac)
# Load a local copy of annotations obtained from the Ensembl Biomart server
data(AlvMac_GOgenes)
data(AlvMac_allgenes)
data(AlvMac_allGO)

# Run the analysis on factor "Treatment",
# considering only treatments "MB" and "TB" at time-point "48H"
# using a local copy of annotations obtained from the Ensembl BioMart server
AlvMac_results <- GO_analyse(
  eSet=AlvMac, f="Treatment",
  subset=list(Time=c("48H"), Treatment=c("MB", "TB")),
  GO_genes=AlvMac_GOgenes, all_genes=AlvMac_allgenes, all_GO=AlvMac_allGO
)

# Valid Ensembl BioMart datasets are listed in the following variable
data(prefix2dataset)

# Valid microarray= values are listed in the following variable
data(microarray2dataset)

## Not run:
# Other valid but time-consuming examples:
```

```

# Run the analysis on factor "Treatment" including all samples
GO_analyse(eSet=AlvMac, f="Treatment")

# Run the analysis on factor "Treatment" using ANOVA method
GO_analyse(eSet=AlvMac, f="Treatment", method="anova")

# Use alternative GO scoring/summarisation functions (Default is: average)

# Named functions
GO_analyse(eSet=AlvMac, f="Treatment", FUN.GO = median)

# Anonymous functions (simple example without scientific value)
GO_analyse(eSet=AlvMac, f="Treatment", FUN.GO = function(x){median(x)/100})

# Syntax examples without actual data:

# To force the use of the Ensembl BioMart for the human species, use:
GO_analyse(eSet, f, biomart_dataset = "hsapiens_gene_ensembl")

# To force use of the bovine affy_bovine microarray annotations use:
GO_analyse(eSet, f, microarray = "affy_bovine")

## End(Not run)

```

heatmap_GO	<i>Generates a heatmap and hierarchical clustering of the samples and the genes</i>
------------	---

Description

Clusters the samples and the genes associated with a GO term using the expression levels of genes related to a given ontology. Represents expression levels of those genes in a heatmap.

Usage

```

heatmap_GO(
  go_id, result, eSet, f=result$factor, subset=NULL, gene_names=TRUE,
  NA.names=FALSE, margins=c(7, 5),
  scale="none", cexCol=1.2, cexRow=0.5,
  labRow=NULL,
  cex.main=1, trace="none", expr.col=bluered(75),
  row.col.palette="Accent",
  row.col=c(),
  main=paste(
    go_id, result$GO[result$GO$go_id == go_id,"name_1006"]
  ),
  main.Lsplit=NULL,
  ...)

```

Arguments

<code>go_id</code>	A Gene Ontology (GO) identifier.
<code>result</code>	The output of <code>GO_analyse()</code> or a subset of it obtained from <code>subset_scores()</code> .
<code>eSet</code>	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the <code>assayData</code> slot, and a phenotypic information data-frame in the <code>phenodata</code> slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
<code>f</code>	The grouping factor in <code>phenodata</code> to label the samples by.
<code>subset</code>	A named list to subset <code>eSet</code> . Names must be column names existing in <code>colnames(pData(eSet))</code> . Values must be vectors of values existing in the corresponding column of <code>pData(eSet)</code> .
<code>gene_names</code>	A boolean value. Default is TRUE, to label genes by their associated gene name. If FALSE, labels the genes by their feature identifier in the expression dataset (i.e. Ensembl gene identifier or microarray probeset).
<code>NA.names</code>	A boolean value. Default is TRUE. If labelling genes by their associated gene name (see argument <code>gene_names</code>), whether to display the gene feature identifier for gene features without associated gene name.
<code>margins</code>	A numeric vector of length 2 specifying the number of lines of margins to apply for the bottom and right margins, respectively. Defaults are 7 (bottom) and 5 (right). For Ensembl gene identifiers, we suggest setting manually a bottom margin of 13. See <code>heatmap.2()</code> .
<code>scale</code>	Character indicating if the values should be centered and scaled in either the row direction or the column direction, or none. Default is "none". See <code>heatmap.2()</code> .
<code>cexCol, cexRow</code>	Positive numbers, used as <code>cex.axis</code> in for the row or column axis labeling. Defaults are 1.2 and 1, respectively. See <code>heatmap.2()</code> .
<code>labRow</code>	A character vector of names to re-label the rows (samples). If vector of length 1, it is assumed to be a column name in the <code>phenoData</code> slot. Default are the values of the factor <code>f</code> . See <code>heatmap.2()</code> .
<code>cex.main</code>	Scaling factor of the main title font size. Default is 1. We suggest to use it in combination with the argument <code>main.Lsplit</code> for GO terms with long names.
<code>trace</code>	Character string indicating whether a solid "trace" line should be drawn across each 'row' or down each 'column', both' or 'none'. The distance of the line from the center of each color-cell is proportional to the size of the measurement. Defaults to 'none'.
<code>expr.col</code>	Character vector indicating the colors to represent the different levels of gene expression. Defaults to a colormap of 75 shades ranging from blue (low) to red (high) and centered around white. If using differential expression data, you should probably use <code>greenred(75)</code> instead.
<code>row.col.palette</code>	A valid RColorBrewer palette name to fetch the colormap from, to color-code the groups of samples.
<code>row.col</code>	A vector of color names or codes. The number of colors provided must match the number of levels of the grouping factor. Default to an palette of up to 9 colors marking the different levels of the predefined grouping factor on the left side of the heatmap.

```

main          Main title of the figure. Default is paste(go_id, go_name).
main.Lsplit   Number of characters after which a new-line character will be inserted in the
              main title. If this would occur within a word, the new-line character will be
              inserted before this word. Default is NULL, leaving the title on a single line.
...          Additional arguments passed on to heatmap.2().

```

Value

Returns the output of the heatmap.2() function.

Author(s)

Kevin Rue-Albrecht

See Also

Method [heatmap.2](#), [GO_analyse](#), and [brewer.pal.info](#).

Examples

```

# load the sample output data
data(AlvMac_results)

# Heatmap the top-ranked GO term (toll-like receptor 4 signaling pathway) as
# example
heatmap_GO(go_id="GO:0034142", result=AlvMac_results, eSet=AlvMac)

# Same with larger sample labels on the right hand side.
heatmap_GO(go_id="GO:0034142", result=AlvMac_results, eSet=AlvMac, cexRow=1)

# Change the color-coding to green-black-red gradient (more appropriate for
# differential expression values)
library(gplots)
heatmap_GO(
  go_id="GO:0034142", result=AlvMac_results, eSet=AlvMac,
  expr.col=greenred(75)
)

```

hist_scores

Plots the distribution of scores following an GOexpress analysis.

Description

Plots the an histogram representing the frequencies of scores in the output variable of the GO_analyse() function.

This function can also be used on the output of subset_scores() function as it returns a value formatted identically to the output of the GO_analyse() function.

Usage

```

hist_scores(result,
            main=paste("Distribution of average scores in",
                      deparse(substitute(result))),
            xlab="Average score", ...)

```

Arguments

result The output of the GO_analyse() function.
 main, xlab These arguments to title have useful defaults here.
 ... Additional arguments passed on to hist().

Value

Returns the output of the hist() function.

Author(s)

Kevin Rue-Albrecht

See Also

Method [hist](#), and [GO_analyse](#).

Examples

```
# load the sample output data with p.values computed
data(AlvMac_results.pVal)

# Histogram of scores (labelled with counts)
hist_scores(result=AlvMac_results, breaks=20, labels=TRUE)

# filter for Biological Processes associated with 5+ genes and <=0.05 P-value
filtered_results <- subset_scores(
  result=AlvMac_results.pVal, total_count=5, p.val=0.05,
  namespace="BP")

# Histogram of scores (labelled with counts)
hist_scores(result=filtered_results, breaks=20, labels=TRUE)
```

list_genes

Returns the genes associated with a Gene Ontology

Description

Given a Gene Ontology (GO) identifier represented in the dataset, returns a character vector listing the feature identifiers annotated to it.

Usage

```
list_genes(go_id, result, data.only=TRUE)
```

Arguments

go_id A Gene Ontology (GO) identifier represented in the dataset.
 result The output of GO_analyse() or a subset of it obtained from subset_scores().
 data.only Whether to return only the feature identifiers present in the given dataset or alternatively returns all feature identifiers associated with the GO term in the Ensembl BioMart. Default is TRUE.

Value

A character vector listing the feature identifiers of the genes associated with the GO term.

Author(s)

Kevin Rue-Albrecht

See Also

Method [GO_analyse](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# List of genes associated with the GO term "toll-like receptor 4 signaling
# pathway"
list_genes(result=AlvMac_results, go_id="GO:0034142")
```

microarray2dataset	<i>Table mapping probeset identifier prefixes to datasets in the Ensembl BioMart.</i>
--------------------	---

Description

The species corresponding to an probeset identifier can often be identified from the prefix of the identifier (e.g. Bt.457.1.S1_at corresponds to *Bos taurus* which is only associated with a single microarray platform in the Ensembl BioMart). This table maps some manually curated unique patterns to the corresponding species and microarray platform.

Usage

```
data(microarray2dataset)
```

Details

All Agilent microarray share the same prefix pattern, making it very difficult to differentiate. Many Affymetrix microarrays also share the same prefix "AFFX" for several probesets.

The dataset and microarray arguments of the GO_analyse method are the best way to specify which BioMart information should be used to annotate the features in your expression dataset.

Value

A data frame with 158 rows and the 6 columns. Rows are sorted alphabetically by species name and each refer to a unique combination of dataset and microarray identifier in the Ensembl BioMart. The columns are described below:

- dataset contains species-specific biomaRt dataset names.
- microarray contains the microarray identifier in the Ensembl BioMart dataset above.
- sample contains a sample probeset from this combination of Ensembl BioMart dataset and microarray.

- `species` contains an example probeset sampled from the corresponding Ensembl BioMart dataset and microarray.
- `pattern` contains the corresponding manually curated pattern representative of probesets for this microarray platform. This pattern may or may not be unique to the microarray. Therefore we encourage users to use the `'microarray'` argument of the `GO_analyse` method to specify the Ensembl BioMart dataset to use.
- `unique` contains a boolean value stating whether the pattern was found unique to the microarray or not.

Source

The `microarray2dataset.build` method stored in the `toolkit.R` script was used to query the Ensembl BioMart server and build this table.

Examples

```
data(microarray2dataset)
microarray2dataset
```

plot_design	<i>Plot Univariate Effects for genes associated with a Gene Ontology</i>
-------------	--

Description

Successively plots univariate effects of one or more [factors](#), typically for a designed experiment as analyzed by [aov\(\)](#).

Usage

```
plot_design(
  go_id, result, eSet, subset=NULL,
  factors=colnames(pData(eSet)), main="", main.Lsplit=NULL, ...)
```

Arguments

<code>go_id</code>	A Gene Ontology (GO) identifier represented by at least one gene in the dataset.
<code>result</code>	The output of <code>GO_analyse()</code> or a subset of it obtained from <code>subset_scores()</code> .
<code>eSet</code>	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the <code>AssayData</code> slot, and a phenotypic information data-frame in the <code>phenodata</code> slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
<code>subset</code>	A named list to subset <code>eSet</code> . Names must be column names existing in <code>colnames(pData(eSet))</code> . Values must be vectors of values existing in the corresponding column of <code>pData(eSet)</code> .
<code>factors</code>	A set of column names from <code>phenodata</code> . Each of these values will be represented on the X-axis to investigate its effect on the average expression of a given genes for each level of that factor.
<code>main</code>	Changes the main title of the plots.

`main.Lsplit` Number of characters after which a new-line character will be inserted in the main title. If this would occur within a word, the new-line character will be inserted before this word. Default is NULL, leaving the title on a single line.

`...` Additional arguments which will be passed on to the `plot.design()` function.

Value

The output of the `plot.design()` function.

Author(s)

Kevin Rue-Albrecht

See Also

Method [plot.design](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# Univariate plot
plot_design(go_id="GO:0034142", eSet=AlvMac, result=AlvMac_results)
```

prefix2dataset

Table mapping Ensembl gene identifier prefixes to BioMart datasets.

Description

The species corresponding to an Ensembl gene identifier can typically be identified from the prefix of the identifier (e.g. ENSBTAG corresponds to *Bos taurus*). This table maps each known unique prefix to the corresponding species.

Usage

```
data(prefix2dataset)
```

Details

C. elegans, *D. melanogaster*, and *S. cerevisiae* have atypical identifier pattern and prefixes in their Ensembl gene identifiers. However, the automatically extracted prefix for *C. elegans* and *D. melanogaster* – respectively "WBgene" and "FBgn" – can be used as such to identify datasets from those species. On the other hand, prefixes used for the *S. cerevisiae* include "YHR", "YAL", and many others. Consequently, expression data from *S. cerevisiae* species is identified without referring to the "prefix2dataset" table; instead, such datasets are identified if the first gene identifier in the dataset starts with "Y".

Value

A data frame with 69 rows and 4 columns. Each row refers to one dataset in the Ensembl BioMart. The columns are described below:

- `dataset` contains the biomaRt dataset name.
- `species` contains the corresponding species name.
- `prefix` contains the corresponding unique prefix.
- `sample` contains a sample Ensembl gene identifier from this dataset.

Source

The `prefix2dataset.build` method stored in the `toolkit.R` script was used to query the Ensembl BioMart server and build this table.

Examples

```
data(prefix2dataset)
str(prefix2dataset)
prefix2dataset
```

pValue_GO

Compute p-values for ontologies by randomising gene labels.

Description

Permutes the rank of genes processed by `GO_analyse()` to compute a P-value for each ontology

Usage

```
pValue_GO(
  result, N=1000, ranked.by=result$rank.by, rank.by='P',
  FUN.GO=result$FUN.GO)
```

Arguments

<code>result</code>	The output of <code>GO_analyse()</code> or a subset of it obtained from <code>subset_scores()</code> .
<code>N</code>	The number of permutation desired. Default to 1000.
<code>ranked.by</code>	Either of 'rank' or 'score'. The metric used to compare whether GO terms in the randomised data are relative to their original value. Default to the current ordering method of the <code>result</code> object.
<code>rank.by</code>	Either of 'rank', 'score', or 'p.val'. The metric used to order the GO terms after computing of the P-value. Default to 'p.val'.
<code>FUN.GO</code>	Function to summarise the score and rank of all feature associated with each gene ontology. Logically, default is the function used in the call to <code>GO_analyse</code> . If using "lambda-like" (anonymous) functions, these must take a list of numeric values as an input, and return a single numeric value as an output.

Details

This function is relatively lengthy. Its procedure can be divided in several steps: (1) assemble a copy of all genes used in the original scoring step, then for each of the N permutations, (2) permute the gene labels, (3) aggregate the rank – or score – of all genes belonging to each GO term, (4) decide whether each GO term is ranked – or scored – better or worse than originally, (5) return a p-value based on the number of permutations where each GO term was better ranked – or scored – than originally.

Value

A list formatted identically to the results of the analysis, with an added 'p.val' column in the GO slot, and re-ordered by the chosen metric.

Author(s)

Kevin Rue-Albrecht

See Also

Method [GO_analyse](#).

Examples

```
## Not run:
# Time-consuming examples:
# Limited here to N=1 for time constraint
# Recommended N=1000 (or at least 100).

# Run the analysis on factor "Treatment" including all samples
AlvMac.Pvals <- pValue_GO(result, N=1, ranked.by=result$rank.by, rank.by='P')

## End(Not run)
```

quantiles_scores *Returns the quantiles of scores following an GOexpress analysis.*

Description

Returns a set of quantiles in indicating the scores reached by given proportions of the GO terms.

This function can also be used on the output of `subset_scores()` function as it returns a value formatted identically to the output of the `GO_analyse()` function.

Usage

```
quantiles_scores(result, probs=c(0.9, 0.95, 0.99, 0.999, 0.9999),
  quantiles=FALSE)
```

Arguments

result	The output of <code>GO_analyse()</code> or a subset of it obtained from <code>subset_scores()</code> .
probs	Numeric vector of probabilities with values in [0,1]. (Values up to 2e-14 outside that range are accepted and moved to the nearby endpoint.) See <code>quantile</code>
quantiles	A numeric vector of the percentiles for which the scores are desired.

Value

A named vector of percentiles and corresponding scores.

Author(s)

Kevin Rue-Albrecht

See Also

Method [quantile](#).

Examples

```
# load the sample output data with p.values computed
data(AlvMac_results.pVal)

# filter for Biological Processes associated with 5+ genes and <=0.05 P-value
filtered_results <- subset_scores(
  result=AlvMac_results.pVal, total_count=5, p.val=0.05,
  namespace="BP")

# Quantiles of scores
quantiles_scores(result=filtered_results)
```

 rerank

Reorder the result variable by alternative metrics.

Description

Reorder the ranked tables of GO terms and genes either by increasing (average) rank or decreasing (average) score.

Usage

```
rerank(result, rank.by = 'rank')
```

Arguments

result	The output of <code>GO_analyse()</code> or a subset of it obtained from <code>subset_scores()</code> .
rank.by	Either of 'rank', 'score' or 'p.val'; the metric to rank the GO terms and genes. Note that 'pval' is only applicable on the output of the <code>pValue_GO()</code> function. See details for breaking ties.

Details

Taking an example, to rank GO terms by P-value and break ties by average rank, rerank first by 'rank', and then rerank the resulting object by 'p.val'.

Value

A list formatted identically to the results of the analysis, but ordered by the chosen metric.

Note

The name `reorder()` was not used to avoid conflict with package `stats`.

Author(s)

Kevin Rue-Albrecht

See Also

Method [GO_analyse](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# Re-rank the GO terms and genes based on the actual score instead of the rank
reranked.byScore <- rerank(result=AlvMac_results, rank.by="score")

# load the sample output data with p.values computed
data(AlvMac_results.pVal)

# To rank by P-value, while breaking the ties by ave_rank,
# rank first by ave_rank
reranked.byRank <- rerank(result=AlvMac_results.pVal, rank.by="rank")
# rank the result by P-value
reranked.pVal_rank <- rerank(result=reranked.byRank, rank.by="p.val")
```

subEset

Subset an ExpressionSet to values of phenotypic data columns.

Description

Given a list of column names and corresponding values present in the `phenoData` slot of an `ExpressionSet` object, this method returns the subset of the `ExpressionSet` restricted to samples associated with the given values in the given columns.

Usage

```
subEset(eSet, subset = list())
```

Arguments

eSet	ExpressionSet of the Biobase package including a gene-by-sample expression matrix in the <code>AssayData</code> slot, and a phenotypic information data-frame in the <code>phenoData</code> slot. In the expression matrix, row names are Ensembl gene identifiers or probeset identifiers, and column names are sample identifiers. In the phenotypic data-frame, row names are sample identifiers, column names are grouping factors and phenotypic traits usable for the one-way ANOVA.
subset	A named list to subset eSet. Names must be column names existing in <code>colnames(pData(eSet))</code> . Values must be vectors of values existing in the corresponding column of <code>pData(eSet)</code> .

Value

A subset of the given ExpressionSet restricted to samples associated with the given values in the given columns

Author(s)

Kevin Rue-Albrecht

See Also

Class [ExpressionSet](#).

Examples

```
# Load example data set
data(AlvMac)

# Subset it to only samples of "CN" and "MB" treatments, and also only "2H",
# "6H", and "24H" time-points
sub.AlvMac <- subEset(
  eSet=AlvMac,
  subset=list(
    Treatment=c("CN", "MB"),
    Time=c("2H", "6H")
  )
)
```

subset_scores	<i>Returns a filtered list from GO_analyse results.</i>
---------------	---

Description

Given an output variable from a GO_analyse analysis and a set of valid filters and thresholds, returns an identically formatted list keeping only the rows of the score table passing all the filters.

Usage

```
subset_scores(result, ...)
```

Suggested use:
 # on the output of GO_analyse()
 # subset_scores(result, total=5, namespace="BP")
 # or after application of pValue_GO()
 # subset_scores(result, total=5, namespace="BP", p.val=0.05)

Arguments

result	The output of the GO_analyse() function.
...	Additional pairs of filter and threshold values in the format "filter=threshold". Filters must be valid names from colnames(result\$GO).

Details

It is highly recommended to filter out GO terms with very few genes (e.g. less than 5 genes), as the scoring function is biased for those GO terms (see UsersGuide).

Suggested filters:

total_count, total:	Filter keeping only GO terms associated with at least the given count of genes in the annotation
p.val, p, P:	Filter keeping only the GO terms with P-value lower or equal to the given cutoff (This filter is not applicable to the 'biological_process' namespace)
namespace, namespace_1003:	Filter keeping only the GO terms of a given type. Valid values are "biological_process", "molecular_function", "cellular_component", "ontology", "ontology_1003"

Additional filters:

data_count, data:	Filter keeping only GO terms associated with at least the given count of genes, present in the ExpressionSet
ave_rank, rank:	Filter keeping only GO terms with an ave_rank value equal or lower than the given cutoff (average of ranks)
ave_score, score:	Filter keeping only GO terms with an ave_score value equal or higher than the given cutoff (average of scores)

Value

A list formatted identically to the `result` object provided, with an added or updated `filters.GO` slot stating the filters and cutoffs applied, and restricted to:

- the gene ontologies passing the given filters and cutoff values.
- the genes mapped to those remaining ontologies (i.e. all genes not associated with an ontology are discarded in the output object).
- gene-mapping annotations related to the remaining gene ontologies.

Note

It is possible to further filter a filter result object. Some warnings may appear if the new filter cutoff conflict with the ones previously applied (stored in the `filters.GO` slot of the filtered object). Example of conflicting filter values are:

- Filtering for 'biological_process' a result object that was previously filtered for 'molecular_function' (no more 'BP' terms are present in the first filtered object).
- filtering for a `total_count` of 5 or more a result object previously filtered for a `total_count` of 10 or more.

Author(s)

Kevin Rue-Albrecht

See Also

Method [GO_analyse](#).

Examples

```
# load the sample output data
data(AlvMac_results.pVal)

# have an overview of the result variable
```

```

str(AlvMac_results.pVal)

# filter for Biological Processes associated with 5+ genes and <=0.05 P-value
filtered_results <- subset_scores(
  result=AlvMac_results.pVal, total_count=5, p.val=0.05,
  namespace="BP")

# have an overview of the filtered result variable
str(filtered_results)

```

table_genes	<i>Returns a table listing the genes associated with a given Gene Ontology</i>
-------------	--

Description

Given a Gene Ontology (GO) identifier represented in the dataset and the output variable of a `GO_analyse()` function, `table_genes()` returns a table listing the genes associated with that `go_id`, their associated gene name, and description.

Usage

```
table_genes(go_id, result, data.only=FALSE, order.by='rank')
```

Arguments

<code>go_id</code>	A Gene Ontology (GO) identifier.
<code>result</code>	The output of the <code>GO_analyse()</code> function.
<code>data.only</code>	Whether to return only the feature identifiers present in the given dataset or alternatively returns all feature identifiers associated with the GO term in the Ensembl BioMart.
<code>order.by</code>	The metric or name to order the output table. Default is increasing 'rank', which corresponds to decreasing score. Valid values are 'rank' (increasing) equivalent to 'score' (decreasing), 'gene_id', 'name', equivalent to 'external_gene_name' and 'external_gene_id'. Text ranking in alphabetical order.

Value

A data frame listing the statistics and annotations for the genes present in the expression dataset and associated with the GO term.

Author(s)

Kevin Rue-Albrecht

See Also

Method [GO_analyse](#).

Examples

```
# load the sample output data
data(AlvMac_results)

# Table of result for genes associated with the GO term
# "toll-like receptor 4 signaling pathway"
table_genes(result=AlvMac_results, go_id="GO:0034142")

# ordered by gene name
table_genes(result=AlvMac_results, go_id="GO:0034142", order.by='name')
```


Index

- *Topic **ExpressionSet**
 - subEset, 36
- *Topic **GOexpress**
 - AlvMac, 4
 - AlvMac_allgenes, 6
 - AlvMac_allGO, 7
 - AlvMac_GOgenes, 8
 - AlvMac_results, 9
 - AlvMac_results.pVal, 10
 - cluster_GO, 11
 - expression_plot, 12
 - expression_plot_symbol, 15
 - expression_profiles, 17
 - expression_profiles_symbol, 19
 - GO_analyse, 22
 - GOexpress-package, 2
 - heatmap_GO, 26
 - hist_scores, 28
 - list_genes, 29
 - microarray2dataset, 30
 - plot_design, 31
 - prefix2dataset, 32
 - pValue_GO, 33
 - quantiles_scores, 34
 - rerank, 35
 - subEset, 36
 - subset_scores, 37
 - table_genes, 39
- *Topic **anova**
 - GO_analyse, 22
 - GOexpress-package, 2
- *Topic **clustering**
 - cluster_GO, 11
 - GO_analyse, 22
 - GOexpress-package, 2
 - heatmap_GO, 26
 - plot_design, 31
- *Topic **datasets**
 - AlvMac, 4
 - AlvMac_allgenes, 6
 - AlvMac_allGO, 7
 - AlvMac_GOgenes, 8
 - AlvMac_results, 9
 - AlvMac_results.pVal, 10
 - microarray2dataset, 30
 - prefix2dataset, 32
- *Topic **expression**
 - cluster_GO, 11
 - expression_plot, 12
 - expression_plot_symbol, 15
 - expression_profiles, 17
 - expression_profiles_symbol, 19
 - GO_analyse, 22
 - GOexpress-package, 2
 - heatmap_GO, 26
 - plot_design, 31
- *Topic **gene**
 - cluster_GO, 11
 - expression_plot, 12
 - expression_plot_symbol, 15
 - expression_profiles, 17
 - expression_profiles_symbol, 19
 - GO_analyse, 22
 - GOexpress-package, 2
 - heatmap_GO, 26
 - list_genes, 29
 - plot_design, 31
 - table_genes, 39
- *Topic **ontology**
 - GO_analyse, 22
 - GOexpress-package, 2
 - heatmap_GO, 26
 - table_genes, 39
- *Topic **package**
 - GOexpress-package, 2
- *Topic **randomForest**
 - GO_analyse, 22
 - GOexpress-package, 2
- AlvMac, 4
- AlvMac_allgenes, 6
- AlvMac_allGO, 7
- AlvMac_GOgenes, 8
- AlvMac_results, 9
- AlvMac_results.pVal, 10
- aov, 31

Biobase, [4](#), [14](#), [16](#), [19](#), [21](#)
biomaRt::getBM, [4](#)
brewer.pal.info, [28](#)

cluster_GO, [11](#)

expression_plot, [12](#), [16](#)
expression_plot_symbol, [14](#), [15](#)
expression_profiles, [17](#)
expression_profiles_symbol, [19](#)
ExpressionSet, [37](#)

factors, [31](#)

getBM, [25](#)
ggplot, [14](#), [16](#), [19](#), [21](#)
ggplot2::ggplot, [4](#)
GO_analyse, [4](#), [12](#), [14](#), [16](#), [22](#), [28–30](#), [34](#), [36](#),
[38](#), [39](#)
GOexpress (GOexpress-package), [2](#)
GOexpress-package, [2](#)
gplots:bluered, [4](#)
gplots:greenred, [4](#)
gplots:heatmap.2, [4](#)
grid:grid.layout, [4](#)
grid:grid.newpage, [4](#)

heatmap.2, [28](#)
heatmap_GO, [26](#)
hist, [29](#)
hist_scores, [28](#)

list_genes, [29](#)

microarray2dataset, [30](#)

oneway.test, [25](#)

plot.design, [32](#)
plot_design, [31](#)
prefix2dataset, [32](#)
pValue_GO, [25](#), [33](#)

quantile, [35](#)
quantiles_scores, [34](#)

randomForest, [4](#), [25](#)
RColorBrewer, [4](#)
rerank, [35](#)

stringr:str_extract, [4](#)
subEset, [36](#)
subset_scores, [25](#), [37](#)

table_genes, [39](#)

VennDiagram, [4](#)