

Package ‘EnrichmentBrowser’

April 14, 2017

Version 2.4.6

Date 2017-03-02

Title Seamless navigation through combined results of set-based and network-based enrichment analysis

Author Ludwig Geistlinger, Gergely Csaba, Ralf Zimmer

Maintainer Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

Depends R(>= 3.0.0), Biobase, GSEABase, pathview

Imports AnnotationDbi, ComplexHeatmap, DESeq2, EDASeq, GO.db, KEGGREST, KEGGgraph, MASS, ReportingTools, Rgraphviz, S4Vectors, SPIA, SummarizedExperiment, biocGraph, edgeR, geneplotter, graph, hwriter, limma, methods, safe, topGO

Suggests ALL, BiocStyle, airway, hgu95av2.db

Description The EnrichmentBrowser package implements essential functionality for the enrichment analysis of gene expression data. The analysis combines the advantages of set-based and network-based enrichment analysis in order to derive high-confidence gene sets and biological pathways that are differentially regulated in the expression data under investigation. Besides, the package facilitates the visualization and exploration of such sets and pathways.

License Artistic-2.0

biocViews Microarray, RNASeq, GeneExpression, DifferentialExpression, Pathways, GraphAndNetwork, Network, GeneSetEnrichment, NetworkEnrichment, Visualization, ReportWriting

NeedsCompilation no

R topics documented:

comb.ea.results	2
compile.grn.from.kegg	4
config.ebrowser	5
de.ana	6
download.kegg.pathways	8
ea.browse	8
ebrowser	10
get.go.genesets	12
get.kegg.genesets	14

ggea.graph	15
make.example.data	16
map.ids	18
nbea	19
normalize	22
plots	24
probe.2.gene.eset	25
read.eset	26
sbea	27

Index	31
--------------	-----------

comb.ea.results	<i>Combining enrichment analysis results</i>
-----------------	--

Description

Different enrichment analysis methods usually result in different gene set rankings for the same dataset. This function allows to combine results from the different set-based and network-based enrichment analysis methods. This includes the computation of average gene set ranks across methods.

Usage

```
comb.ea.results( res.list,
  rank.col=config.ebrowser("GSP.COL"),
  decreasing=FALSE,
  rank.fun = c("comp.ranks", "rel.ranks", "abs.ranks"),
  comb.fun = c("mean", "median", "min", "max", "sum") )
```

Arguments

res.list	A list of enrichment analysis result lists (as returned by the functions sbea and nbea).
rank.col	Rank column. Column name of the enrichment analysis result table that should be used to rank the gene sets. Defaults to the gene set p-value column, i.e. gene sets are ranked according to gene set significance.
decreasing	Logical. Should smaller (decreasing=FALSE, default) or larger (decreasing=TRUE) values in rank.col be ranked better? In case of gene set p-values the smaller the better, in case of gene set scores the larger the better.
rank.fun	Ranking function. Used to rank gene sets according to the result table of individual enrichment methods (as returned from the gs.ranking function). This is typically done according to gene set p-values, but can also take into account gene set scores/statistics, especially in case of gene sets with equal p-value. Can be either one of the predefined functions ('comp.ranks', 'rel.ranks', 'abs.ranks') or a user-defined function. Defaults to 'comp.ranks', i.e. competitive (percentile) ranks are computed by calculating for each gene set the percentage of gene sets with a p-value as small or smaller. Alternatively, 'rel.ranks', i.e. relative ranks are computed in 2 steps:

1. Ranks are assigned according to distinct gene set p-value **categories**, i.e. gene sets with equal p-value obtain the **same** rank. Thus, the gene sets with lowest p-value obtain rank 1, and so on.
2. As opposed to absolute ranks (`rank.fun = 'abs.ranks'`), which are returned from step 1, relative ranks are then computed by dividing the absolute rank by number of distinct p-value categories and multiplying with 100 (= percentile rank).

`comb.fun` Rank combination function. Used to combine gene set ranks across methods. Can be either one of the predefined functions (mean, median, max, min, sum) or a user-defined function. Defaults to 'sum', i.e. the rank sum across methods is computed.

Value

An enrichment analysis result list that can be detailedly explored by calling [ea.browse](#) and from which a flat gene set ranking can be extracted by calling [gs.ranking](#).

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[sbea](#), [nbea](#), [ea.browse](#)

Examples

```
# (1) expression set:
# simulated expression values of 100 genes
# in two sample groups of 6 samples each
eset <- make.example.data(what="eset")
eset <- de.ana(eset)

# (2) gene sets:
# draw 10 gene sets with 15-25 genes
gs <- make.example.data(what="gs", gnames=featureNames(eset))

# (3) make artificial enrichment analysis results:
# 2 ea methods with 5 significantly enriched gene sets each
ora.res <- make.example.data(what="ea.res", method="ora", eset=eset, gs=gs)
gsea.res <- make.example.data(what="ea.res", method="gsea", eset=eset, gs=gs)

# (4) combining the results
res.list <- list(ora.res, gsea.res)
comb.res <- comb.ea.results(res.list)

# (5) result visualization and exploration
gs.ranking(comb.res)

# user-defined ranking and combination functions
# (a) dummy ranking, give 1:nrow(res.tbl)
dummy.rank <- function(res.tbl) seq_len(nrow(res.tbl))

# (b) weighted average for combining ranks
```

```
wavg <- function(r) mean(c(1,2) * r)

comb.res <- comb.ea.results(res.list, rank.fun=dummy.rank, comb.fun=wavg)
```

compile.grn.from.kegg *Compilation of a gene regulatory network from KEGG pathways*

Description

To perform network-based enrichment analysis a gene regulatory network (GRN) is required. There are well-studied processes and organisms for which comprehensive and well-annotated regulatory networks are available, e.g. the RegulonDB for *E. coli* and Yeasttract for *S. cerevisiae*. However, in many cases such a network is missing. A first simple workaround is to compile a network from regulations in the KEGG database.

Usage

```
compile.grn.from.kegg( pwys, out.file = NULL )
```

Arguments

pwys	Either a list of KEGGPathway objects or an absolute file path of a zip compressed archive of pathway xml files in KGML format. Alternatively, you can specify an organism in KEGG three letter code, e.g. 'hsa' for 'Homo sapiens', and the pathways will be downloaded automatically.
out.file	Optional output file the gene regulatory network will be written to.

Value

if(is.null(out.file)): the gene regulatory network; else: none, as the gene regulatory network is written to file

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[KEGGPathway-class](#), [parseKGML](#), [download.kegg.pathways](#)

Examples

```
# (1) download human pathways
# pwys <- download.kegg.pathways("hsa")
# (2) compile gene regulatory network
# grn <- compile.grn.from.kegg(pwys)

pwys <- system.file("extdata/hsa_kegg_pwys.zip", package="EnrichmentBrowser")
hsa.grn <- compile.grn.from.kegg(pwys)
```

Description

Function to get and set configuration parameters determining the default behavior of the EnrichmentBrowser

Usage

```
config.ebrowser( key, value = NULL )
```

Arguments

key	Configuration parameter.
value	Value to overwrite the current value of key.

Details

Important pData, fData, and result column names:

- SMPL.COL: pData column storing the sample IDs (default: "SAMPLE")
- GRP.COL: pData column storing binary group assignment (default: "GROUP")
- BLK.COL: pData column defining paired samples or sample blocks (default: "BLOCK")
- PRB.COL: fData column storing probe/feature IDs ("PROBEID", read-only)
- EZ.COL: fData column storing gene ENTREZ IDs ("ENTREZID", read-only)
- SYM.COL: fData column storing gene symbols ("SYMBOL", read-only)
- GN.COL: fData column storing gene names ("GENENAME", read-only)
- FC.COL: fData column storing (log2) fold changes of differential expression between sample groups (default: "FC")
- ADJP.COL: fData column storing adjusted (corrected for multiple testing) p-values of differential expression between sample groups (default: "ADJ.PVAL")
- GS.COL: result table column storing gene set IDs (default: "GENE.SET")
- GSP.COL: result table column storing gene set significance (default: "P.VALUE")
- PMID.COL: gene table column storing PUBMED IDs ("PUBMED", read-only)

Important URLs (all read-only):

- NCBI.URL: <http://www.ncbi.nlm.nih.gov/>
- PUBMED.URL: <http://www.ncbi.nlm.nih.gov/pubmed/>
- GENE.URL: <http://www.ncbi.nlm.nih.gov/gene/>
- KEGG.URL: <http://www.genome.jp/dbget-bin/>
- KEGG.GENE.URL: http://www.genome.jp/dbget-bin/www_bget?
- KEGG.SHOW.URL: http://www.genome.jp/dbget-bin/show_pathway?
- GO.SHOW.URL: <http://amigo.geneontology.org/amigo/term/>

Default output directory:

- EBROWSER.HOME: `system.file(package="EnrichmentBrowser")`
- OUTDIR.DEFAULT: `file.path(EBROWSER.HOME, "results")`

Gene set size:

- GS.MIN.SIZE: minimum number of genes per gene set (default: 5)
- GS.MAX.SIZE: maximum number of genes per gene set (default: 500)

Result appearance:

- RESULT.TITLE: (default: "Table of Results")
- NR.SHOW: maximum number of entries to show (default: 20)

Value

If `is.null(value)` this returns the value of the selected configuration parameter. Otherwise, it updates the selected parameter with the given value.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

Examples

```
# getting config information
config.ebrowser("GS.MIN.SIZE")

# setting config information
# WARNING: this is for advanced users only!
# inappropriate settings will impair EnrichmentBrowser's functionality
config.ebrowser(key="GS.MIN.SIZE", value=3)
```

de.ana

Differential expression analysis between two sample groups

Description

The function carries out a differential expression analysis between two sample groups. Resulting fold changes and derived p-values are returned. Raw p-values are corrected for multiple testing.

Usage

```
de.ana( expr, grp = NULL, blk = NULL,
        de.method = c("limma", "edgeR", "DESeq"), padj.method = "BH", stat.only=FALSE, min.cpm=2 )
```

Arguments

expr	Expression data. A numeric matrix. Rows correspond to genes, columns to samples. Alternatively, this can also be an object of class <code>ExpressionSet</code> (in case of microarray data) or an object of class <code>SeqExpressionSet</code> (in case of RNA-seq data). See the man page of <code>read.eset</code> for prerequisites for the expression data.
grp	<i>*BINARY*</i> group assignment for the samples. Use '0' and '1' for unaffected (controls) and affected (cases) samples, respectively. If NULL, this is assumed to be defined via a column named 'GROUP' in the pData slot if 'expr' is a (Seq)ExpressionSet.
blk	Optional. For paired samples or sample blocks. This can also be defined via a column named 'BLOCK' in the pData slot if 'expr' is a (Seq)ExpressionSet.
de.method	Differential expression method. Use 'limma' for microarray and RNA-seq data. Alternatively, differential expression for RNA-seq data can be also calculated using edgeR ('edgeR') or DESeq2 ('DESeq'). Defaults to 'limma'.
padj.method	Method for adjusting p-values to multiple testing. For available methods see the man of page the of the stats function <code>p.adjust</code> . Defaults to 'BH'.
stat.only	Logical. Should only the test statistic be returned? This is mainly for internal use, in order to carry out permutation tests on the DE statistic for each gene. Defaults to FALSE.
min.cpm	In case of RNA-seq data: should genes not satisfying a minimum counts-per-million (cpm) threshold be excluded from the analysis? This is typically recommended. See the edgeR vignette for details. The default filter is to exclude genes with $\text{cpm} < 2$ in more than half of the samples.

Value

A DE-table with measures of differential expression for each gene/row, i.e. a two-column matrix with log₂ fold changes in the 1st column and derived p-values in the 2nd column. If 'expr' is a (Seq)ExpressionSet, the DE-table will be automatically appended to the fData slot.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

`read.eset` describes prerequisites for the expression data, `normalize` for normalization of expression data, `voom` for preprocessing of RNA-seq data, `p.adjust` for multiple testing correction, `eBayes` for DE analysis with limma, `glmFit` for DE analysis with edgeR, and `DESeq` for DE analysis with DESeq.

Examples

```
# (1) microarray data: intensity measurements
ma.eset <- make.example.data(what="eset", type="ma")
ma.eset <- de.ana(ma.eset)
head(fData(ma.eset))

# (2) RNA-seq data: read counts
rseq.eset <- make.example.data(what="eset", type="rseq")
rseq.eset <- de.ana(rseq.eset, de.method="DESeq")
head(fData(rseq.eset))
```

download.kegg.pathways

Download of KEGG pathways for a particular organism

Description

The function downloads all metabolic and non-metabolic pathways in KEGG XML format for a specified organism.

Usage

```
download.kegg.pathways( org, out.dir = NULL, zip = FALSE )
```

Arguments

org	Organism in KEGG three letter code, e.g. 'hsa' for 'homo sapiens'.
out.dir	Output directory. If not null, pathways are written to files in the specified directory.
zip	Logical. In case pathways are written to file ('out.dir' is not null): should output files be zipped?

Value

if(is.null(out.dir)): a list of KEGGPathway objects else: none, as pathways are written to file

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[keggList](#), [keggGet](#), [KEGGPathway-class](#), [parseKGML](#)

Examples

```
pwys <- download.kegg.pathways("hsa")
```

ea.browse

Exploration of enrichment analysis results

Description

Functions to extract a flat gene set ranking from an enrichment analysis result object and to detailedly explore it.

Usage

```
ea.browse( res, nr.show = -1, graph.view = NULL, html.only = FALSE )

gs.ranking( res, signif.only = TRUE )
```

Arguments

res	Enrichment analysis result list (as returned by the functions sbea and nbea).
nr.show	Number of gene sets to show. As default all statistically significant gene sets are displayed.
graph.view	Optional. Should a graph-based summary (reports and visualizes consistency of regulations) be created for the result? If specified, it needs to be a gene regulatory network, i.e. either an absolute file path to a tabular file or a character matrix with exactly <i>*THREE*</i> cols; 1st col = IDs of regulating genes; 2nd col = corresponding regulated genes; 3rd col = regulation effect; Use '+' and '-' for activation/inhibition.
html.only	Logical. Should the html file only be written (without opening the browser to view the result page)? Defaults to FALSE.
signif.only	Logical. Display only those gene sets in the ranking, which satisfy the significance level? Defaults to TRUE.

Value

gs.ranking: DataFrame with gene sets ranked by the corresponding p-value;
ea.browse: none, opens the browser to explore results.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[sbea](#), [nbea](#), [comb.ea.results](#)

Examples

```
# real data
# (1) reading the expression data from file
exprs.file <- system.file("extdata/exprs.tab", package="EnrichmentBrowser")
pdat.file <- system.file("extdata/pData.tab", package="EnrichmentBrowser")
fdat.file <- system.file("extdata/fData.tab", package="EnrichmentBrowser")
probe.eset <- read.eset(exprs.file, pdat.file, fdat.file)
gene.eset <- probe.2.gene.eset(probe.eset)
gene.eset <- de.ana(gene.eset)
annotation(gene.eset) <- "hsa"

# artificial enrichment analysis results
gs <- make.example.data(what="gs", gnames=featureNames(gene.eset))
ea.res <- make.example.data(what="ea.res", method="ora", eset=gene.eset, gs=gs)

# (5) result visualization and exploration
gs.ranking(ea.res)
ea.browse(ea.res)
```

Description

This is the all-in-one wrapper function to perform the standard enrichment analysis pipeline implemented in the EnrichmentBrowser package.

Given flat gene expression data, the data is read in and subsequently subjected to chosen enrichment analysis methods.

The results from different methods can be combined and investigated in detail in the default browser.

Usage

```
ebrowser( meth, exprs, pdat, fdat, org, data.type = c(NA, "ma", "rseq"),
          norm.method = "quantile", de.method = "limma",
          gs, grn = NULL, perm = 1000, alpha = 0.05, beta = 1,
          comb = FALSE, browse = TRUE, nr.show = -1 )
```

Arguments

meth	Enrichment analysis method. Currently, the following enrichment analysis methods are supported: 'ora', 'safe', 'gsea', 'samgs', 'ggea', 'spia', 'nea', and 'pathnet'. See sbea and nbea for details.
exprs	Expression matrix. A tab separated text file containing <i>*normalized*</i> expression values on a <i>*log*</i> scale. Columns = samples/subjects; rows = features/probes/genes; NO headers, row or column names. Supported data types are log2 counts (microarray single-channel), log2 ratios (microarray two-color), and log2-counts per million (RNA-seq logCPMs). See limma's user guide for definition and normalization of the different data types. Alternatively, this can be an object of ExpressionSet-class , assuming the expression matrix in the 'exprs' slot.
pdat	Phenotype data. A tab separated text file containing annotation information for the samples in either <i>*two or three*</i> columns. NO headers, row or column names. The number of rows/samples in this file should match the number of columns/samples of the expression matrix. The 1st column is reserved for the sample IDs; The 2nd column is reserved for a <i>*BINARY*</i> group assignment. Use '0' and '1' for unaffected (controls) and affected (cases) sample class, respectively. For paired samples or sample blocks a third column is expected that defines the blocks. If 'exprs' is an object of ExpressionSet-class , the 'pdat' argument can be left unspecified, which then expects group and optional block assignments in respectively named columns 'GROUP' (mandatory) and 'BLOCK' (optional) in the 'pData' slot of the ExpressionSet.
fdat	Feature data. A tab separated text file containing annotation information for the features. Exactly <i>*TWO*</i> columns; 1st col = feature IDs; 2nd col = corresponding KEGG gene ID for each feature ID in 1st col; NO headers, row or column names. The number of rows/features in this file should match the number of rows/features of the expression matrix. If 'exprs' is an object of ExpressionSet-class , the 'fdat' argument can be left unspecified, which then expects feature and gene IDs in respectively named columns 'PROBE' and 'GENE' in the 'fData' slot of the ExpressionSet.

org	Organism under investigation in KEGG three letter code, e.g. 'hsa' for 'Homo sapiens'. See also kegg.species.code to convert your organism of choice to KEGG three letter code.
data.type	Expression data type. Use 'ma' for microarray and 'rseq' for RNA-seq data. If NA, data.type is automatically guessed. If the expression values in 'eset' are decimal numbers they are assumed to be microarray intensities. Whole numbers are assumed to be RNA-seq read counts. Defaults to NA.
norm.method	Determines whether and how the expression data should be normalized. For available microarray normalization methods see the man page of the limma function normalizeBetweenArrays . For available RNA-seq normalization methods see the man page of the EDASeq function betweenLaneNormalization . Defaults to 'quantile', i.e. normalization is carried out so that quantiles between arrays/lanes/samples are equal. Use 'none' to indicate that the data is already normalized and should not be normalized by ebrowser. See the man page of normalize for details.
de.method	Determines which method is used for per-gene differential expression analysis. See the man page of de.ana for details. Defaults to 'limma', i.e. differential expression is calculated based on the typical limma lmFit procedure.
gs	Gene sets. Either a list of gene sets (vectors of KEGG gene IDs) or a text file in GMT format storing all gene sets under investigation.
grn	Gene regulatory network. Either an absolute file path to a tabular file or a character matrix with exactly *THREE* cols; 1st col = IDs of regulating genes; 2nd col = corresponding regulated genes; 3rd col = regulation effect; Use '+' and '-' for activation/inhibition.
perm	Number of permutations of the expression matrix to estimate the null distribution. Defaults to 1000. Can also be an integer vector matching the length of 'meth' to assign different numbers of permutations for different methods.
alpha	Statistical significance level. Defaults to 0.05.
beta	Log2 fold change significance level. Defaults to 1 (2-fold).
comb	Logical. Should results be combined if more than one enrichment method is selected? Defaults to FALSE.
browse	Logical. Should results be displayed in the browser for interactive exploration? Defaults to TRUE.
nr.show	Number of gene sets to show. As default all statistical significant gene sets are displayed.

Value

None, opens the browser to explore results.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

References

Limma User's guide: <http://www.bioconductor.org/packages/limma>

See Also

[read.eset](#) to read expression data from file; [probe.2.gene.eset](#) to transform probe to gene level expression; [kegg.species.code](#) maps species name to KEGG code. [get.kegg.genesets](#) to retrieve gene set definitions from KEGG; [compile.grn.from.kegg](#) to construct a GRN from KEGG pathways; [sbea](#) to perform set-based enrichment analysis; [nbea](#) to perform network-based enrichment analysis; [comb.ea.results](#) to combine results from different methods; [ea.browse](#) for exploration of resulting gene sets

Examples

```
# expression data from file
exprs.file <- system.file("extdata/exprs.tab", package="EnrichmentBrowser")
pdat.file <- system.file("extdata/pData.tab", package="EnrichmentBrowser")
fdat.file <- system.file("extdata/fData.tab", package="EnrichmentBrowser")

# getting all human KEGG gene sets
# hsa.gs <- get.kegg.genesets("hsa")
gs.file <- system.file("extdata/hsa_kegg_gs.gmt", package="EnrichmentBrowser")
hsa.gs <- parse.genesets.from.GMT(gs.file)

# set-based enrichment analysis
ebrowser(  meth="ora",
           exprs=exprs.file, pdat=pdat.file, fdat=fdat.file,
           gs=hsa.gs, org="hsa", nr.show=3)

# compile a gene regulatory network from KEGG pathways
# hsa.grn <- compile.grn.from.kegg("hsa")
pwys <- system.file("extdata/hsa_kegg_pwys.zip", package="EnrichmentBrowser")
hsa.grn <- compile.grn.from.kegg(pwys)

# network-based enrichment analysis
ebrowser(  meth="ggea",
           exprs=exprs.file, pdat=pdat.file, fdat=fdat.file,
           gs=hsa.gs, grn=hsa.grn, org="hsa", nr.show=3 )

# combining results
ebrowser(  meth=c("ora", "ggea"), comb=TRUE,
           exprs=exprs.file, pdat=pdat.file, fdat=fdat.file,
           gs=hsa.gs, grn=hsa.grn, org="hsa", nr.show=3 )
```

```
get.go.genesets
```

```
Definition of gene sets according to the Gene Ontology (GO)
```

Description

This function retrieves GO gene sets for an organism under investigation either via download from BioMart or based on BioC annotation packages.

Usage

```
get.go.genesets( org, onto = c("BP", "MF", "CC"), mode = c("GO.db","biomart") )
```

Arguments

org	An organism in (KEGG) three letter code, e.g. 'hsa' for 'Homo sapiens'.
onto	Character. Specifies one of the three GO ontologies: 'BP' (biological process), 'MF' (molecular function), 'CC' (cellular component). Defaults to 'BP'.
mode	Character. Determines in which way the gene sets are retrieved. This can be either 'GO.db' or 'biomart'. The 'GO.db' mode creates the gene sets based on BioC annotation packages - which is fast, but represents not necessarily the most up-to-date mapping. In addition, this option is only available for the currently supported model organisms in BioC. The 'biomart' mode downloads the mapping from BioMart - which can be time consuming, but allows to select from a larger range of organisms and contains the latest mappings. Defaults to 'GO.db'.

Value

A list of gene sets (vectors of gene IDs).

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

References

<http://geneontology.org/>

See Also

[annFUN](#) for general GO2gene mapping used in the 'GO.db' mode, the [biomaRt](#) package for general queries to BioMart, [get.kegg.genesets](#) for defining gene sets according to KEGG, [parse.genesets.from.GMT](#) to parse user-def. gene sets from file.

Examples

```
# Typical usage for gene set enrichment analysis:
# Biological process terms based on BioC annotation (for human)
gs <- get.go.genesets("hsa")

# eq.:
# gs <- get.go.genesets(org="hsa", onto="BP", mode="GO.db")

# Alternatively:
# downloading from BioMart
# this may take a few minutes ...

gs <- get.go.genesets(org="hsa", mode="biomart")
```

get.kegg.genesets	<i>Definition of gene sets according to KEGG pathways for a specified organism</i>
-------------------	--

Description

To perform a gene set enrichment analysis on KEGG pathways, it is necessary to build up the gene set database in a format that the GSEA method can read. Parsing a list of gene sets from a flat text file in GMT format. This function performs the necessary steps, including the retrieval of the participating gene IDs for each pathway and the conversion to GMT format.

Usage

```
get.kegg.genesets( pwys, gmt.file = NULL )
```

```
parse.genesets.from.GMT( gmt.file )
```

Arguments

pwys	Either a list of KEGGPathway objects or an absolute file path of a zip compressed archive of pathway xml files in KGML format. Alternatively, an organism in KEGG three letter code, e.g. 'hsa' for 'Homo sapiens'.
gmt.file	Gene set file in GMT format. See details.

Details

The GMT (Gene Matrix Transposed) file format is a tab delimited file format that describes gene sets. In the GMT format, each row represents a gene set. Each gene set is described by a name, a description, and the genes in the gene set. See references.

Value

A list of gene sets (vectors of gene IDs).

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

References

GMT file format http://www.broadinstitute.org/cancer/software/gsea/wiki/index.php/Data_formats

KEGG Organism code http://www.genome.jp/kegg/catalog/org_list.html

See Also

[keggList](#), [keggLink](#), [KEGGPathway-class](#), [parseKGML](#)

Examples

```
# WAYS TO DEFINE GENE SETS ACCORDING TO HUMAN KEGG PATHWAYS

# (1) from scratch: via organism ID

gs <- get.kegg.genesets("hsa")

# (2) extract from pathways
# download human pathways via:
# pwys <- download.kegg.pathways("hsa")
pwys <- system.file("extdata/hsa_kegg_pwys.zip", package="EnrichmentBrowser")
gs <- get.kegg.genesets(pwys)

# (3) parsing gene sets from GMT
gmt.file <- system.file("extdata/hsa_kegg_gs.gmt", package="EnrichmentBrowser")
gs <- parse.genesets.from.GMT(gmt.file)
```

ggea.graph

GGEA graphs of consistency between regulation and expression

Description

Gene graph enrichment analysis (GGEA) is a network-based enrichment analysis method implemented in the `EnrichmentBrowser` package. The idea of GGEA is to evaluate the consistency of known regulatory interactions with the observed gene expression data. A GGEA graph for a gene set of interest displays the consistency of each interaction in the network that involves a gene set member. Nodes (genes) are colored according to expression (up-/down-regulated) and edges (interactions) are colored according to consistency, i.e. how well the interaction type (activation/inhibition) is reflected in the correlation of the expression of both interaction partners.

Usage

```
ggea.graph( gs, grn, eset,
            alpha = 0.05, beta = 1, max.edges = 50, cons.thresh = 0.7 )

ggea.graph.legend()
```

Arguments

gs	Gene set under investigation. This should be a character vector of KEGG gene IDs.
grn	Gene regulatory network. Character matrix with exactly <i>*THREE*</i> cols; 1st col = IDs of regulating genes; 2nd col = corresponding regulated genes; 3rd col = regulation effect; Use '+' and '-' for activation/inhibition.
eset	Expression set. An object of class <code>ExpressionSet</code> containing the gene expression set. See <code>read.eset</code> and <code>probe.2.gene.eset</code> for required annotations in the <code>pData</code> and <code>fData</code> slot.
alpha	Statistical significance level. Defaults to 0.05.
beta	Log2 fold change significance level. Defaults to 1 (2-fold).

max.edges Maximum number of edges that should be displayed. Defaults to 50.
cons.thresh Consistency threshold. Graphical parameter that correspondingly increases line width of edges with a consistency above the chosen threshold (defaults to 0.7).

Value

None, plots to a graphics device.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[nbea](#) to perform network-based enrichment analysis. [ea.browse](#) for exploration of resulting gene sets.

Examples

```
# (1) expression set:  
# simulated expression values of 100 genes  
# in two sample groups of 6 samples each  
eset <- make.example.data(what="eset")  
eset <- de.ana(eset)  
  
# (2) gene sets:  
# draw 10 gene sets with 15-25 genes  
gs <- make.example.data(what="gs", gnames=featureNames(eset))  
  
# (3) compiling artificial regulatory network  
grn <- make.example.data(what="grn", nodes=featureNames(eset))  
  
# (4) plot consistency graph  
ggee.graph(gs=gs[[1]], grn=grn, eset=eset)  
  
# (5) get legend  
ggee.graph.legend()
```

make.example.data *Example data for the EnrichmentBrowser package*

Description

Functionality to construct example data sets for demonstration. This includes expression data, gene sets, gene regulatory networks, and enrichment analysis results.

Usage

```
make.example.data( what = c("eset", "gs", "grn", "ea.res"), ... )
```


Arguments

- what** Kind of example data set to be constructed. This should be one out of:
- eset: Expression set
 - gs: Gene set list
 - grn: Gene regulatory network
 - ea.res: Enrichment analysis result object as returned by the functions [sbea](#) and [nbea](#)
- ...** Additional arguments to fine-tune the specific example data sets.
- For `what='eset'`:
- type: Expression data type. Should be either 'ma' (Microarray intensity measurements) or 'rseq' (RNA-seq read counts).
 - nfeat: Number of features/genes. Defaults to 100.
 - nsmpl: Number of samples. Defaults to 12.
 - blk: Create sample blocks. Defaults to TRUE.
 - norm: Should the expression data be normalized? Defaults to FALSE.
 - de.ana: Should an differential expression analysis be carried out automatically? Defaults to FALSE.
- For `what='gs'`:
- gnames: gene names from which the sets will be sampled. Per default the sets will be drawn from `c(g1, ..., g100)`.
 - n: number of sets. Defaults to 10.
 - min.size: minimal set size. Defaults to 15.
 - max.size: maximal set size. Defaults to 25.
- For `what='grn'`:
- nodes: gene node names for which edges will be drawn. Per default node names will be `c(g1, ..., g100)`.
 - edge.node.ratio: ratio number of edges / number of nodes. Defaults to 3, i.e. creates 3 times more edges than nodes.
- For `what='ea.res'`:
- eset: Expression set. Calls `make.example.data(what="eset")` per default.
 - gs: Gene sets. Calls `make.example.data(what="gs")` per default.
 - method: Enrichment analysis method. Defaults to 'ora'.
 - alpha: Statistical significance level. Defaults to 0.05.

Value

Depends on the 'what' argument.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifl.lmu.de>

Examples

```
eset <- make.example.data(what="eset")
```

map.ids	<i>Mapping between gene ID types for feature names of an expression set</i>
---------	---

Description

Functionality to map feature names of an expression set between common gene ID types such as ENSEMBL and ENTREZ.

Usage

```
id.types( org )
```

```
map.ids( eset, org=NA, from="ENSEMBL", to="ENTREZID" )
```

Arguments

eset	Expression set. An object of ExpressionSet-class . Expects the featureNames to be of gene ID type given in argument 'from'.
org	Organism in KEGG three letter code, e.g. 'hsa' for 'Homo sapiens'. See references.
from	Gene ID type from which should be mapped. Corresponds to the gene ID type of the featureNames of argument 'eset'. Defaults to 'ENSEMBL'.
to	Gene ID type to which should be mapped. Corresponds to the gene ID type the featuresNames of argument 'eset' should be updated with. Defaults to 'ENTREZID'.

Details

The function 'id.types' lists the valid values which the arguments 'from' and 'to' can take. This corresponds to the names of the available gene ID types for the mapping.

Value

id.types: character vector listing the available gene ID types for the mapping;
map.ids: An object of [ExpressionSet](#).

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

References

KEGG Organism code http://www.genome.jp/kegg/catalog/org_list.html

See Also

[ExpressionSet-class](#), [mapIds](#), [keytypes](#)

Examples

```
# create an expression set with 3 genes and 3 samples
eset <- make.example.data("eset", nfeat=3, nsmpl=3)
featureNames(eset) <- paste0("ENSG000000000", c("003", "005", "419"))
eset <- map.ids(eset, org="hsa")
```

nbea

Network-based enrichment analysis (NBEA)

Description

This is the main function for network-based enrichment analysis. It implements and wraps existing implementations of several frequently used state-of-art methods and allows a flexible inspection of resulting gene set rankings.

Usage

```
nbea( method = EnrichmentBrowser::nbea.methods(), eset, gs, grn,
      prune.grn=TRUE, alpha = 0.05, perm = 1000, padj.method = "none",
      out.file = NULL, browse = FALSE, ... )
```

```
nbea.methods()
```

Arguments

method	Network-based enrichment analysis method. Currently, the following network-based enrichment analysis methods are supported: 'ggea', 'spia', 'pathnet', 'degraph', 'topologygsa', 'ganpa', 'cepa', 'netgsa', and 'nea'. Default is 'ggea'. This can also be the name of a user-defined function implementing network-based enrichment. See Details.
eset	Expression set. An object of class ExpressionSet . See read.eset and probe.2.gene.eset for required annotations in the pData and fData slot.
gs	Gene sets. Either a list of gene sets (vectors of KEGG gene IDs) or a text file in GMT format storing all gene sets under investigation.
grn	Gene regulatory network. Either an absolute file path to a tabular file or a character matrix with exactly <i>*THREE*</i> cols; 1st col = IDs of regulating genes; 2nd col = corresponding regulated genes; 3rd col = regulation effect; Use '+' and '-' for activation/inhibition.
prune.grn	Logical. Should the GRN be pruned? This removes duplicated, self, and reversed edges. Defaults to TRUE.
alpha	Statistical significance level. Defaults to 0.05.
perm	Number of permutations of the expression matrix to estimate the null distribution. Defaults to 1000. If using method='ggea', it is possible to set 'perm=0' to use a fast approximation of gene set significance to avoid permutation testing. See Details.
padj.method	Method for adjusting nominal gene set p-values to multiple testing. For available methods see the man of page the of the stats function p.adjust . Defaults to 'none', i.e. leaves the nominal gene set p-values unadjusted.

<code>out.file</code>	Optional output file the gene set ranking will be written to.
<code>browse</code>	Logical. Should results be displayed in the browser for interactive exploration? Defaults to FALSE.
<code>...</code>	Additional arguments passed to individual nbea methods. This includes currently: <ul style="list-style-type: none"> • <code>beta</code>: Log2 fold change significance level. Defaults to 1 (2-fold). For SPIA and NEA: <ul style="list-style-type: none"> • <code>sig.stat</code>: decides which statistic is used for determining significant DE genes. Options are: <ul style="list-style-type: none"> – <code>'p'</code> (Default): genes with p-value below alpha. – <code>'fc'</code>: genes with $\text{abs}(\log_2(\text{fold change}))$ above beta – <code>'&'</code>: p & fc (logical AND) – <code>' '</code>: p fc (logical OR) For GGEA: <ul style="list-style-type: none"> • <code>cons.thresh</code>: edge consistency threshold between -1 and 1. Defaults to 0.2, i.e. only edges of the GRN with consistency ≥ 0.2 are included in the analysis. Evaluation on real datasets has shown that this works best to distinguish relevant gene sets. Use consistency of -1 to include all edges. • <code>gs.edges</code>: decides which edges of the grn are considered for a gene set under investigation. Should be one out of c(<code>'&'</code>, <code>' '</code>), denoting logical AND and OR. respectively. Accordingly, this either includes edges for which regulator AND / OR target gene are members of the investigated gene set.

Details

`'ggea'`: gene graph enrichment analysis, scores gene sets according to consistency within the given gene regulatory network, i.e. checks activating regulations for positive correlation and repressing regulations for negative correlation of regulator and target gene expression (Geistlinger et al., 2011). When using `'ggea'` it is possible to estimate the statistical significance of the consistency score of each gene set in two different ways: (1) based on sample permutation as described in the original publication (Geistlinger et al., 2011) or (2) using an approximation in the spirit of Bioconductor's npGSEA package that is much faster.

`'spia'`: signaling pathway impact analysis, combines ORA with the probability that expression changes are propagated across the pathway topology; implemented in Bioconductor's SPIA package.

`'pathnet'`: pathway analysis using network information, applies ORA on combined evidence for the observed signal for gene nodes and the signal implied by connected neighbors in the network; implemented in Bioconductor's PathNet package.

`'degraph'`: differential expression testing for gene graphs, multivariate testing of differences in mean incorporating underlying graph structure; implemented in Bioconductor's DEGraph package

`'topologygsa'`: topology-based gene set analysis, uses Gaussian graphical models to incorporate the dependence structure among genes as implied by pathway topology; implemented in CRAN's topologyGSA package.

`'ganpa'`: gene association network-based pathway analysis, incorporates network-derived gene weights in the enrichment analysis; implemented in CRAN's GANPA package.

`'cepa'`: centrality-based pathway enrichment, incorporates network centralities as node weights mapped from differentially expressed genes in pathways; implemented in CRAN's CePa package.

'netgsa': network-based gene set analysis, incorporates external information about interactions among genes as well as novel interactions learned from data; implemented in CRAN's NetGSA package.

'nea': network enrichment analysis, applies ORA on interactions instead of genes; implemented in Bioconductor's neaGUI package.

It is also possible to use additional network-based enrichment methods. This requires to implement a function that takes 'eset', 'gs', 'grn', 'alpha', and 'perm' as arguments and returns a numeric matrix 'res.tbl' with a mandatory column named 'P.VALUE' storing the resulting p-value for each gene set in 'gs'. The rows of this matrix must be named accordingly (i.e. `rownames(res.tbl) == names(gs)`). See examples.

Value

`nbea.methods`: a character vector of currently supported methods;

`nbea`: `if(is.null(out.file))`: an enrichment analysis result object that can be detailedly explored by calling `ea.browse` and from which a flat gene set ranking can be extracted by calling `gs.ranking`. If 'out.file' is given, the ranking is written to the specified file.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

References

Geistlinger et al. (2011) From sets to graphs: towards a realistic enrichment analysis of transcriptomic systems. *Bioinformatics*, 27(13), i366–73.

See Also

Input: `read.eset`, `probe.2.gene.eset`, `get.kegg.genesets` to retrieve gene set definitions from KEGG. `compile.grn.from.kegg` to construct a GRN from KEGG pathways.

Output: `gs.ranking` to rank the list of gene sets. `ea.browse` for exploration of resulting gene sets.

Other: `sbea` to perform set-based enrichment analysis. `comb.ea.results` to combine results from different methods. the SPIA package for more information on signaling pathway impact analysis. the neaGUI package for more information on network enrichment analysis. the PathNet package for more information on pathway analysis using network information.

Examples

```
# currently supported methods
nbea.methods()

# (1) expression data:
# simulated expression values of 100 genes
# in two sample groups of 6 samples each
eset <- make.example.data(what="eset")
eset <- de.ana(eset)

# (2) gene sets:
# draw 10 gene sets with 15-25 genes
gs <- make.example.data(what="gs", gnames=featureNames(eset))

# (3) make 2 artificially enriched sets:
```

```

sig.genes <- featureNames(eset)[fData(eset)$ADJ.PVAL < 0.1]
gs[[1]] <- sample(sig.genes, length(gs[[1]]))
gs[[2]] <- sample(sig.genes, length(gs[[2]]))

# (4) gene regulatory network
grn <- make.example.data(what="grn", nodes=featureNames(eset))

# (5) performing the enrichment analysis
ea.res <- nbea(method="ggea", eset=eset, gs=gs, grn=grn)

# (6) result visualization and exploration
gs.ranking(ea.res, signif.only=FALSE)

# using your own tailored function as enrichment method
dummy.nbea <- function(eset, gs, grn, alpha, perm)
{
  sig.ps <- sample(seq(0,0.05, length=1000),5)
  insig.ps <- sample(seq(0.1,1, length=1000), length(gs)-5)
  ps <- sample(c(sig.ps, insig.ps), length(gs))
  score <- sample(1:100, length(gs), replace=TRUE)
  res.tbl <- cbind(score, ps)
  colnames(res.tbl) <- c("SCORE", "P.VALUE")
  rownames(res.tbl) <- names(gs)
  return(res.tbl[order(ps),])
}

ea.res2 <- nbea(method=dummy.nbea, eset=eset, gs=gs, grn=grn)
gs.ranking(ea.res2)

```

normalize

Normalization of microarray and RNA-seq expression data

Description

This function wraps commonly used functionality from limma for microarray normalization and from EDASeq for RNA-seq normalization.

Usage

```
normalize( eset,
          norm.method = "quantile", within = FALSE, data.type = c("ma", "rseq") )
```

Arguments

- | | |
|-------------|---|
| eset | Expression set. An object of ExpressionSet-class . See the man page of read.eset for prerequisites for the expression data. |
| norm.method | Determines how the expression data should be normalized. For available microarray normalization methods see the man page of the limma function normalizeBetweenArrays . For available RNA-seq normalization methods see the man page of the EDASeq function betweenLaneNormalization . Defaults to 'quantile', i.e. normalization is carried out so that quantiles between arrays/lanes/samples are equal. See details. |

within	Logical. Is only taken into account if <code>data.type='rseq'</code> . Determine whether GC content normalization should be carried out (as implemented in the EDASeq function withinLaneNormalization). Defaults to FALSE. See details.
data.type	Expression data type. Use 'ma' for microarray and 'rseq' for RNA-seq data. If NA, <code>data.type</code> is automatically guessed. If the expression values in 'eset' are decimal numbers they are assumed to be microarray intensities. Whole numbers are assumed to be RNA-seq read counts. Defaults to NA.

Details

Normalization of high-throughput expression data is essential to make results within and between experiments comparable. Microarray (intensity measurements) and RNA-seq (read counts) data exhibit typically distinct features that need to be normalized for. For specific needs that deviate from these standard normalizations, the user should always refer to more specific functions/packages.

Microarray data is expected to be single-channel. For two-color arrays, it is expected here that normalization within arrays has been already carried out, e.g. using [normalizeWithinArrays](#) from `limma`.

RNA-seq data is expected to be raw read counts. Please note that normalization for downstream DE analysis, e.g. with `edgeR` and `DESeq`, is not ultimately necessary (and in some cases even discouraged) as many of these tools implement specific normalization approaches. See the vignette of `EDASeq`, `edgeR`, and `DESeq` for details.

Value

An object of [ExpressionSet-class](#). For RNA-seq data, an object of [SeqExpressionSet-class](#) to conform with downstream DE analysis.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[read.eset](#) describes prerequisites for the expression data;
[normalizeWithinArrays](#) and [normalizeBetweenArrays](#) for normalization of microarray data;
[withinLaneNormalization](#) and [betweenLaneNormalization](#) for normalization of RNA-seq data.

Examples

```
#
# (1) simulating expression data: 100 genes, 12 samples
#

# (a) microarray data: intensity measurements
ma.eset <- make.example.data(what="eset", type="ma")

# (b) RNA-seq data: read counts
rseq.eset <- make.example.data(what="eset", type="rseq")

#
# (2) Normalization
#
```

```
# (a) microarray ...
norm.eset <- normalize(ma.eset)

# (b) RNA-seq ...
norm.eset <- normalize(rseq.eset)

# ... normalize also for GC content
gc.content <- rnorm(100, 0.5, sd=0.1)
fData(rseq.eset)$gc <- gc.content

norm.eset <- normalize(rseq.eset, within=TRUE)
```

plots

Visualization of gene expression

Description

Visualization of differential gene expression via heatmap, p-value histogram and volcano plot (fold change vs. p-value).

Usage

```
pdistr( p )
volcano( fc, p )
exprs.heatmap( expr, grp )
```

Arguments

p	Numeric vector of p-values for each gene.
fc	Numeric vector of fold changes (typically on log2 scale).
expr	Expression matrix. Rows correspond to genes, columns to samples.
grp	<i>*BINARY*</i> group assignment for the samples. Use '0' and '1' for unaffected (controls) and affected (cases) samples, respectively.

Value

None, plots to a graphics device.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifl.lmu.de>

See Also

[de.ana](#) for differential expression analysis, [heatmap](#) and [truehist](#) for generic plotting.

Examples

```
# (1) simulating expression data: 100 genes, 12 samples
eset <- make.example.data(what="eset")

# plot heatmap
exprs.heatmap(expr=exprs(eset), grp=as.factor(pData(eset)$GROUP))

# (2) DE analysis
eset <- de.ana(eset)
pdistr(fData(eset)$ADJ.PVAL)
volcano(fc=fData(eset)$FC, p=fData(eset)$ADJ.PVAL)
```

probe.2.gene.eset *Transformation of probe level expression to gene level expression*

Description

Reads expression data at probe level and summarizes gene expression behavior by averaging over all probes that are annotated to a particular gene.

Usage

```
probe.2.gene.eset( probe.eset, use.mean = TRUE )
```

Arguments

probe.eset	Probe expression set of class ExpressionSet . The fData slot of the expression set must contain a 'GENE' column that lists for each probe the corresponding KEGG gene ID.
use.mean	Logical. Determining, in case of multiple probes for one gene, whether a mean value is computed (use.mean=TRUE), or the probe that discriminate the most between the two sample group is kept (use.mean=FALSE). Defaults to TRUE.

Value

An [ExpressionSet](#) on gene level.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[ExpressionSet-class](#), [read.eset](#) for reading expression data from file, [de.ana](#) for differential expression analysis.

Examples

```
# (1) reading the expression data from file
exprs.file <- system.file("extdata/exprs.tab", package="EnrichmentBrowser")
pdat.file <- system.file("extdata/pData.tab", package="EnrichmentBrowser")
fdat.file <- system.file("extdata/fData.tab", package="EnrichmentBrowser")
probe.eset <- read.eset(exprs.file, pdat.file, fdat.file)
gene.eset <- probe.2.gene.eset(probe.eset)
```

read.eset

*Reading gene expression data from file into an expression set***Description**

The function reads in plain expression data from file with minimum annotation requirements for the pData and fData slots.

Usage

```
read.eset( exprs.file, pdat.file, fdat.file,
           data.type = c(NA, "ma", "rseq"), NA.method = c("mean", "rm", "keep") )
```

Arguments

- | | |
|------------|---|
| exprs.file | Expression matrix. A tab separated text file containing expression values. Columns = samples/subjects; rows = features/probes/genes; NO headers, row or column names. See details. |
| pdat.file | Phenotype data. A tab separated text file containing annotation information for the samples in either <i>*two or three*</i> columns. NO headers, row or column names. The number of rows/samples in this file should match the number of columns/samples of the expression matrix. The 1st column is reserved for the sample IDs; The 2nd column is reserved for a <i>*BINARY*</i> group assignment. Use '0' and '1' for unaffected (controls) and affected (cases) sample class, respectively. For paired samples or sample blocks a third column is expected that defines the blocks. |
| fdat.file | Feature data. A tab separated text file containing annotation information for the features. In case of probe level data: exactly <i>*TWO*</i> columns; 1st col = probe/feature IDs; 2nd col = corresponding gene ID for each feature ID in 1st col; In case of gene level data: The list of gene IDs newline-separated (i.e. just one column). It is recommended to use <i>*ENTREZ*</i> gene IDs (to benefit from downstream visualization and exploration functionality of the enrichment analysis). NO headers, row or column names. The number of rows (features/probes/genes) in this file should match the number of rows/features of the expression matrix. Alternatively, this can also be the ID of a recognized platform such as 'hgu95av2' (Affymetrix Human Genome U95 chip) or 'ecoli2' (Affymetrix E. coli Genome 2.0 Array). See details. |
| data.type | Expression data type. Use 'ma' for microarray and 'rseq' for RNA-seq data. If NA, data.type is automatically guessed. If the expression values in 'read.eset' are decimal numbers they are assumed to be microarray intensities. Whole numbers are assumed to be RNA-seq read counts. Defaults to NA. |
| NA.method | Determines how to deal with NA's (missing values). This can be one out of: <ul style="list-style-type: none"> • mean: replace NA's by the row means for a feature over all samples. • rm: rows (features) that contain NA's are removed. • keep: do nothing. Missing values are kept (which, however, can then cause several issues in the downstream analysis) |

Defaults to 'mean'.

Details

See the limma's user guide <http://www.bioconductor.org/packages/limma> for definition and normalization of the different expression data types.

In case of microarray data the feature IDs typically correspond to probe IDs. Thus, the fdat.file should define a mapping from probe ID (1st column) to corresponding KEGG gene ID (2nd column). The mapping can be defined automatically by providing the ID of a recognized platform such as 'hgu95av2' (Affymetrix Human Genome U95 chip). This requires that a corresponding '.db' package exists (see http://www.bioconductor.org/packages/release/BiocViews.html#___ChipName for all available chips/packages) and that you have it installed. *However, this option should be used with care*. Existing mappings might be outdated and sometimes the KEGG gene ID does not correspond to the Entrez ID (e.g. for E. coli and S. cerevisae). In these cases probe identifiers are mapped twice (probe ID -> Entrez ID -> KEGG ID), which almost always results in loss of information. Thus, mapping quality should always be checked and in case properly defined with a 2-column fdat.file.

Value

An object of [ExpressionSet](#).

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

See Also

[ExpressionSet-class](#)

Examples

```
# reading the expression data from file
exprs.file <- system.file("extdata/exprs.tab", package="EnrichmentBrowser")
pdat.file <- system.file("extdata/pData.tab", package="EnrichmentBrowser")
fdat.file <- system.file("extdata/fData.tab", package="EnrichmentBrowser")
eset <- read.eset(exprs.file, pdat.file, fdat.file)
```

sbea

Set-based enrichment analysis (SBEA)

Description

This is the main function for the enrichment analysis of gene sets. It implements and wraps existing implementations of several frequently used state-of-art methods and allows a flexible inspection of resulting gene set rankings.

Usage

```
sbea( method = EnrichmentBrowser::sbea.methods(), eset, gs, alpha = 0.05,
      perm = 1000, padj.method = "none", out.file = NULL, browse = FALSE, ... )

sbea.methods()
```

Arguments

method	Set-based enrichment analysis method. Currently, the following set-based enrichment analysis methods are supported: 'ora', 'safe', 'gsea', 'padog', 'roast', 'camera', 'gsa', 'gsva', 'globaltest', 'samgs', 'ebm', and 'mgsa'. For basic ora also set 'perm=0'. Default is 'ora'. This can also be the name of a user-defined function implementing set-based enrichment. See Details.
eset	Expression set. An object of class <code>ExpressionSet</code> . See <code>read.eset</code> and <code>probe.2.gene.eset</code> for required annotations in the <code>pData</code> and <code>fData</code> slots.
gs	Gene sets. Either a list of gene sets (vectors of KEGG gene IDs) or a text file in GMT format storing all gene sets under investigation.
alpha	Statistical significance level. Defaults to 0.05.
perm	Number of permutations of the expression matrix to estimate the null distribution. Defaults to 1000. For basic ora set 'perm=0'. Using <code>method="gsea"</code> and 'perm=0' invokes the permutation approximation from the <code>npGSEA</code> package.
padj.method	Method for adjusting nominal gene set p-values to multiple testing. For available methods see the man page of the stats function <code>p.adjust</code> . Defaults to 'none', i.e. leaves the nominal gene set p-values unadjusted.
out.file	Optional output file the gene set ranking will be written to.
browse	Logical. Should results be displayed in the browser for interactive exploration? Defaults to FALSE.
...	Additional arguments passed to individual sbea methods. This includes currently for ORA and MGSA: <ul style="list-style-type: none"> • beta: Log2 fold change significance level. Defaults to 1 (2-fold). • sig.stat: decides which statistic is used for determining significant DE genes. Options are: <ul style="list-style-type: none"> – 'p' (Default): genes with p-value below alpha. – 'fc': genes with $\text{abs}(\log_2(\text{fold change}))$ above beta – '&': p & fc (logical AND) – ' ': p fc (logical OR)

Details

'ora': overrepresentation analysis, simple and frequently used test based on the hypergeometric distribution (see Goeman and Buhlmann, 2007, for a critical review).

'safe': significance analysis of function and expression, generalization of ORA, includes other test statistics, e.g. Wilcoxon's rank sum, and allows to estimate the significance of gene sets by sample permutation; implemented in the `safe` package (Barry et al., 2005).

'gsea': gene set enrichment analysis, frequently used and widely accepted, uses a Kolmogorov-Smirnov statistic to test whether the ranks of the p-values of genes in a gene set resemble a uniform distribution (Subramanian et al., 2005).

'padog': pathway analysis with down-weighting of overlapping genes, incorporates gene weights to favor genes appearing in few pathways versus genes that appear in many pathways; implemented in the `PADOG` package.

'roast': rotation gene set test, uses rotation instead of permutation for assessment of gene set significance; implemented in the `limma` and `edgeR` packages for microarray and RNA-seq data, respectively.

'camera': correlation adjusted mean rank gene set test, accounts for inter-gene correlations as implemented in the limma and edgeR packages for microarray and RNA-seq data, respectively.

'gsa': gene set analysis, differs from GSEA by using the maxmean statistic, i.e. the mean of the positive or negative part of gene scores in the gene set; implemented in the GSA package.

'gsva': gene set variation analysis, transforms the data from a gene by sample matrix to a gene set by sample matrix, thereby allowing the evaluation of gene set enrichment for each sample; implemented in the GSVA package.

'globaltest': global testing of groups of genes, general test of groups of genes for association with a response variable; implemented in the globaltest package.

'samgs': significance analysis of microarrays on gene sets, extends the SAM method for single genes to gene set analysis (Dinu et al., 2007).

'ebm': empirical Brown's method, combines p -values of genes in a gene set using Brown's method to combine p -values from dependent tests; implemented in the EmpiricalBrownsMethod package.

'mgsa': model-based gene set analysis, Bayesian modeling approach taking set overlap into account by working on all sets simultaneously, thereby reducing the number of redundant sets; implemented in the mgsa package.

It is also possible to use additional set-based enrichment methods. This requires to implement a function that takes 'eset', 'gs', 'alpha', and 'perm' as arguments and returns a numeric vector 'ps' storing the resulting p -value for each gene set in 'gs'. This vector must be named accordingly (i.e. `names(ps) == names(gs)`). See examples.

Value

`sbea.methods`: a character vector of currently supported methods;

`sbea`: `if(is.null(out.file))`: an enrichment analysis result object that can be detailedly explored by calling `ea.browse` and from which a flat gene set ranking can be extracted by calling `gs.ranking`. If 'out.file' is given, the ranking is written to the specified file.

Author(s)

Ludwig Geistlinger <Ludwig.Geistlinger@bio.ifi.lmu.de>

References

Goeman and Buhlmann (2007) Analyzing gene expression data in terms of gene sets: methodological issues. *Bioinformatics*, 23, 980-7.

Barry et al. (2005) Significance Analysis of Function and Expression. *Bioinformatics*, 21:1943-9.

Subramanian et al. (2005) Gene Set Enrichment Analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci USA*, 102:15545-50.

Dinu et al. (2007) Improving gene set analysis of microarray data by SAM-GS. *BMC Bioinformatics*, 8:242

See Also

Input: `read.eset`, `probe.2.gene.eset` `get.kegg.genesets` to retrieve gene sets from KEGG.

Output: `gs.ranking` to retrieve the ranked list of gene sets. `ea.browse` for exploration of resulting gene sets.

Other: `nbea` to perform network-based enrichment analysis. `comb.ea.results` to combine results from different methods.

Examples

```
# currently supported methods
sbea.methods()

# (1) expression data:
# simulated expression values of 100 genes
# in two sample groups of 6 samples each
eset <- make.example.data(what="eset")
eset <- de.ana(eset)

# (2) gene sets:
# draw 10 gene sets with 15-25 genes
gs <- make.example.data(what="gs", gnames=featureNames(eset))

# (3) make 2 artificially enriched sets:
sig.genes <- featureNames(eset)[fData(eset)$ADJ.PVAL < 0.1]
gs[[1]] <- sample(sig.genes, length(gs[[1]]))
gs[[2]] <- sample(sig.genes, length(gs[[2]]))

# (4) performing the enrichment analysis
ea.res <- sbea(method="ora", eset=eset, gs=gs, perm=0)

# (5) result visualization and exploration
gs.ranking(ea.res)

# using your own tailored function as enrichment method
dummy.sbea <- function(eset, gs, alpha, perm)
{
  sig.ps <- sample(seq(0, 0.05, length=1000), 5)
  nsig.ps <- sample(seq(0.1, 1, length=1000), length(gs)-5)
  ps <- sample(c(sig.ps, nsig.ps), length(gs))
  names(ps) <- names(gs)
  return(ps)
}

ea.res2 <- sbea(method=dummy.sbea, eset=eset, gs=gs)
gs.ranking(ea.res2)
```

Index

annFUN, [13](#)

betweenLaneNormalization, [11](#), [22](#), [23](#)

comb.ea.results, [2](#), [9](#), [12](#), [21](#), [29](#)
compile.grn.from.kegg, [4](#), [12](#), [21](#)
config.ebrowser, [5](#)

de.ana, [6](#), [11](#), [24](#), [25](#)
DESeq, [7](#)
download.kegg.pathways, [4](#), [8](#)

ea.browse, [3](#), [8](#), [12](#), [16](#), [21](#), [29](#)
eBayes, [7](#)
ebrowser, [10](#)
ExpressionSet, [7](#), [15](#), [18](#), [19](#), [25](#), [27](#), [28](#)
exprs.heatmap (plots), [24](#)

get.go.genesets, [12](#)
get.kegg.genesets, [12](#), [13](#), [14](#), [21](#), [29](#)
ggea (nbea), [19](#)
ggea.graph, [15](#)
glmFit, [7](#)
gs.ranking, [2](#), [3](#), [21](#), [29](#)
gs.ranking (ea.browse), [8](#)
gsea (sbea), [27](#)

heatmap, [24](#)

id.types (map.ids), [18](#)

kegg.species.code, [11](#), [12](#)
keggGet, [8](#)
keggLink, [14](#)
keggList, [8](#), [14](#)
KEGGPathway, [4](#), [14](#)
keytypes, [18](#)

lmFit, [11](#)

make.example.data, [16](#)
map.ids, [18](#)
mapIds, [18](#)

nbea, [2](#), [3](#), [9](#), [10](#), [12](#), [16](#), [17](#), [19](#), [29](#)
normalize, [7](#), [11](#), [22](#)
normalizeBetweenArrays, [11](#), [22](#), [23](#)
normalizeWithinArrays, [23](#)

ora (sbea), [27](#)

p.adjust, [7](#), [19](#), [28](#)
parse.genesets.from.GMT, [13](#)
parse.genesets.from.GMT
(get.kegg.genesets), [14](#)
parseKGML, [4](#), [8](#), [14](#)
pdistr (plots), [24](#)
plots, [24](#)
probe.2.gene.eset, [12](#), [15](#), [19](#), [21](#), [25](#), [28](#), [29](#)

read.eset, [7](#), [12](#), [15](#), [19](#), [21–23](#), [25](#), [26](#), [28](#), [29](#)

sbea, [2](#), [3](#), [9](#), [10](#), [12](#), [17](#), [21](#), [27](#)
SeqExpressionSet, [7](#)

truehist, [24](#)

volcano (plots), [24](#)
voom, [7](#)

withinLaneNormalization, [23](#)