

Package ‘KEGGprofile’

April 15, 2020

Type Package

Title An annotation and visualization package for multi-types and multi-groups expression data in KEGG pathway

Version 1.28.0

Date 2019-09-10

Author Shilin Zhao, Yan Guo, Yu Shyr

Maintainer Shilin Zhao <shilin.zhao@vanderbilt.edu>

Description KEGGprofile is an annotation and visualization tool which integrated the expression profiles and the function annotation in KEGG pathway maps. The multi-types and multi-groups expression data can be visualized in one pathway map. KEGGprofile facilitated more detailed analysis about the specific function changes inner pathway or temporal correlations in different genes and samples.

License GPL (>= 2)

LazyLoad yes

Imports

AnnotationDbi,png,TeachingDemos,XML,KEGG.db,KEGGREST,biomaRt,RCurl,ggplot2,reshape2

biocViews Pathways, KEGG

RoxygenNote 6.1.1

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

git_branch RELEASE_3_10

git_last_commit 1a1a267

git_last_commit_date 2019-10-29

Date/Publication 2020-04-14

R topics documented:

KEGGprofile-package	2
col_by_value	2
convertId	3
download_KEGGfile	4
download_latest_pathway	4
find_enriched_pathway	5
newIdMatrix	6

parse_XMLfile	7
pho_sites_count	7
plot_pathway	8
plot_pathway_cor	9
plot_pathway_overall	10
plot_profile	11
pro_pho_expr	12

Index 13

KEGGprofile-package *KEGGprofile: An annotation and visualization package for multi-types and multi-groups expression data in KEGG pathway*

Description

KEGGprofile is an annotation and visualization tool which integrated the expression profiles and the function annotation in KEGG pathway maps. The multi-types and multi-groups expression data can be visualized in one pathway map. KEGGprofile facilitated more detailed analysis about the specific function changes inner pathway or temporal correlations in different genes and samples.

col_by_value *col_by_value*

Description

The function will transfer a numeric matrix into a matrix of colors, in which the colors represent the values of numeric matrix

Usage

```
col_by_value(x, col, range = NA, breaks = NA, showColorBar = T)
```

Arguments

x	a numeric matrix
col	colors used to represent the values. (See also 'Details')
range	values out of the range will be modified to in the range.
breaks	a numeric vector of three or more cut points giving the number of intervals into which x is to be cut. See also 'Details'
showColorBar	Logical. Indicates display the colorbar or not. The default value is TRUE.

Details

A colorbar would also be plotted. The returned colors of the function can be used in function plot_profile. if breaks not equal to NA, col must have the same length with breaks-1.

Value

a matrix equal to x, but the values were instead by colors.

Examples

```
data(pho_sites_count)
col<-col_by_value(pho_sites_count,col=colorRampPalette(c('white','khaki2'))(4),
breaks=c(0,1,4,10,Inf))
```

 convertId

convertId

Description

A function to convert ID based on the biomaRt package.

Usage

```
convertId(x, dataset = "hsapiens_gene_ensembl",
  filters = "uniprotswissprot", attributes = c(filters,
  "entrezgene_id"), genesKept = c("foldchange", "first", "random", "var",
  "abs"), keepNoId = T, keepMultipleId = F, verbose = F)
```

Arguments

x	the expression data matrix.
dataset	Dataset you want to use. To see the different datasets available within a biomaRt you can e.g. do: mart = useMart('ensembl'), followed by listDatasets(mart).
filters	Filters (one or more) that should be used in the query. A possible list of filters can be retrieved using the function listFilters.
attributes	Attributes you want to retrieve. A possible list of attributes can be retrieved using the function listAttributes.
genesKept	The method to select target gene in more than one targets. "var"/"foldchange"/"abs" means selecting the gene with largest variation/fold change/absolute value. "first" means selecting the first target and "random" means randomly selection.
keepNoId	Logical. Indicate keep the source IDs without target IDs or not.
keepMultipleId	Logical. Indicate keep the multiple target IDs related to one source ID or not.
verbose	Logical. Indicate report extra information on progress or not.

Details

A function to convert ID based on the biomaRt package..

Examples

```
temp<-cbind(rnorm(10),rnorm(10))
row.names(temp)<-c("Q04837","P0C0L4","P0C0L5","075379","Q13068","A2MYD1",
"P60709","P30462","P30475","P30479")
colnames(temp)<-c("Exp1","Exp2")
convertId(temp,filters="uniprotswissprot",keepMultipleId=TRUE)
## Not run:
temp<-cbind(rnorm(5000),rnorm(5000),rnorm(5000),rnorm(5000),rnorm(5000),rnorm(5000))
row.names(temp)<-1000:5999
colnames(temp)<-c("Control1","Control2","Control3","Treatment1","Treatment2","Treatment3")
```

```
convertId(temp, filters="entrezgene_id", attributes = c("entrezgene_id", "uniprotswissprot"),
keepNoId=FALSE)
```

```
## End(Not run)
```

```
download_KEGGfile      download_KEGGfile
```

Description

The function download XML files and png files from KEGG website to local disk

Usage

```
download_KEGGfile(pathway_id = "00010", species = "hsa",
target_dir = getwd())
```

Arguments

<code>pathway_id</code>	the KEGG pathway id, such as '00010'
<code>species</code>	the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc
<code>target_dir</code>	the local directory where the downloaded files are saved

Details

If `pathway_id` is set as 'all', all KEGG pathway ids in KEGG.db package will be used and downloaded from KEGG website

Examples

```
download_KEGGfile(pathway_id="00010", species='hsa')
```

```
download_latest_pathway
                          download_latest_pathway
```

Description

The function will download the latest pathway gene link from KEGG website.

Usage

```
download_latest_pathway(species)
```

Arguments

<code>species</code>	the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc
----------------------	---

Details

The function will download the latest pathway gene link from KEGG website.

Value

a list with two parts

name keggpathway2gene

description a list with the genes for each pathway

name pathway2name

description a list with the names for each pathway

Examples

```
## Not run: download_latest_pathway(species="hsa")
```

```
find_enriched_pathway find_enriched_pathway
```

Description

The function will map the genes in KEGG pathway database, and then hypergeometric tests would be used to estimate the significance of enrichment for each pathway

Usage

```
find_enriched_pathway(gene, species = "hsa", returned_pvalue = 0.01,
  returned_adjvalue = 0.05, returned_genenumber = 5,
  download_latest = FALSE, refGene = NULL)
```

Arguments

gene a numeric matrix

species the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc

returned_pvalue

the minimum p value for enriched pathways

returned_adjvalue

the minimum adjusted p value for enriched pathways

returned_genenumber

the minimum number of annotated genes for enriched pathways

download_latest

logical. Indicate if the function will download the latest pathway/gene link from KEGG website. As the KEGG.db package was not updated for a long time due to the KEGG policy change, we provided this parameter so that the users could get the latest KEGG database.

refGene

the names of genes used as reference. If not provided, all genes in KEGG database will be used.

Details

Only the pathways with p value \leq returned_pvalue in hypergeometric tests and number of annotated genes \geq returned_genenumber would be taken as enriched and returned.

Value

a list with two parts

name stastic	description a matrix containing the pathway IDs of enriched pathways, and their names, p values, number of annotated genes
name detail	description a list with the genes annotated for each pathway

Examples

```
data(pho_sites_count)
#the 300 genes with most phosphorylation sites quantified
genes<-names(rev(sort(pho_sites_count[,1]))[1:300])
pho_KEGGresult<-find_enriched_pathway(genes,species='hsa')
```

newIdMatrix

newIdMatrix

Description

A function to convert ID.

Usage

```
newIdMatrix(x, convertIdTable, genesKept = c("var", "foldchange", "abs",
      "first", "random"))
```

Arguments

x	the expression data matrix.
convertIdTable	A vector. The names should be the source IDs, and the values should be the target IDs.
genesKept	The method to select target gene in more than one targets. "var"/"foldchange"/"abs" means selecting the gene with largest variation/fold change/absolute value. "first" means selecting the first target and "random" means randomly selection.

Details

A function to convert ID.

Examples

```
convertIdTable<-paste("New",c(1,2,2,2,1,3,4,4,5,5))
names(convertIdTable)<-paste("Old",1:length(convertIdTable))
temp<-matrix(rnorm(20),ncol=2)
row.names(temp)<-names(convertIdTable)
colnames(temp)<-c("Exp1", "Exp2")
newIdMatrix(temp,genesKept="foldchange",convertIdTable)
```

parse_XMLfile	<i>parse_XMLfile</i>
---------------	----------------------

Description

The function parses KEGG XML (KGML) files

Usage

```
parse_XMLfile(pathway_id, species, database_dir = getwd())
```

Arguments

pathway_id	the KEGG pathway id, such as '00010'
species	the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc
database_dir	the directory where the XML files and png files are located

Details

This function will parse the KEGG XML (KGML) file. Then a matrix with genes in this pathway and related informations will be returned. This matrix can be used for plot the expression profiles on the pathway figure.

Value

a matrix containing genes in this pathway, and their names, locations etc, which could be used in the function plot_profile as param KEGG_database

Examples

```
XML2database<-parse_XMLfile(pathway_id="04110",species="hsa",
database_dir=system.file("extdata",package="KEGGprofile"))
```

pho_sites_count	<i>number of phosphorylation sites quantified for each gene</i>
-----------------	---

Description

This data set is a data.frame with number of phosphorylation sites quantified for each gene in the analysis.

Usage

```
pho_sites_count
```

Source

Olsen, J.V., et al. (2010) Quantitative phosphoproteomics reveals widespread full phosphorylation site occupancy during mitosis, Sci Signal, 3, ra3.

plot_pathway	<i>plot_pathway</i>
--------------	---------------------

Description

A wrapper for function `download_KEGGfile`, `parse_XMLfile` and `plot_profile`

Usage

```
plot_pathway(gene_expr, line_col, groups, pathway_id = "00010",
             species = "hsa", pathway_min = 5, database_dir = getwd(),
             speciesRefMap = TRUE, ...)
```

Arguments

<code>gene_expr</code>	the matrix for gene expression, row.names should be NCBI gene ID, such as 67040, 93683
<code>line_col</code>	line color for expression in different samples in the pathway map, valid when <code>type='lines'</code>
<code>groups</code>	a character used to indicate expression values from different types of samples
<code>pathway_id</code>	the KEGG pathway id, such as '00010'
<code>species</code>	the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc
<code>pathway_min</code>	The pathways with number of annotated genes less than <code>pathway_min</code> would be ignored
<code>database_dir</code>	the directory where the XML files and png files are located
<code>speciesRefMap</code>	Logical, use the species specific figure as reference map. if set as FALSE, the reference pathway figure without species information will be used
<code>...</code>	any other Arguments for function <code>plot_profile</code>

Details

This wrapper function is developed to make the visualization process more easier. Firstly the existence of XML file and png file would be checked, if not, the `download_KEGGfile` function would be used to download the files. Then the `parse_XMLfile` function would be used to parse the XML file. At last the `plot_profile` function would be used to generate the pathway map.

See Also

[download_KEGGfile](#), [parse_XMLfile](#), [plot_profile](#)

Examples

```
data(pro_pho_expr)
data(pho_sites_count)
#type='lines'
col<-col_by_value(pho_sites_count,col=colorRampPalette(c('white','khaki2'))(4),
breaks=c(0,1,4,10,Inf))
temp<-plot_pathway(pro_pho_expr,bg_col=col,line_col=c("brown1","seagreen3"),
groups=c(rep("Proteome ",6),rep("Phosphoproteome ",6)),magnify=1.2,species='hsa',
```



```

database_dir=system.file("extdata",package="KEGGprofile"),pathway_id="04110",max_dist=5)
#type='bg'
pho_expr<-pro_pho_expr[,7:12]
temp<-apply(pho_expr,1,function(x) length(which(is.na(x))))
pho_expr<-pho_expr[which(temp==0),]
col<-col_by_value(pho_expr,col=colorRampPalette(c('green','black','red'))(1024),range=c(-6,6))
temp<-plot_pathway(pho_expr,type="bg",bg_col=col,text_col="white",magnify=1.2,species='hsa',
database_dir=system.file("extdata",package="KEGGprofile"),pathway_id="04110")
#Compound and gene data
set.seed(124)
testData1<-rbind(rnorm(6),rnorm(6),rnorm(6),rnorm(6),rnorm(6),rnorm(6),rnorm(6),rnorm(6))
row.names(testData1)<-c("4967","55753","1743","8802","47","50","cpd:C15972","cpd:C16255")
colnames(testData1)<-c("Control0","Control2","Control5","Sample0","Sample2","Sample5")
temp<-plot_pathway(testData1,type="lines",line_col=c("brown1","seagreen3"),
groups=c(rep("Control",3),rep("Sample",3)),magnify=1.2,species='hsa',
database_dir=system.file("extdata",package="KEGGprofile"),pathway_id="00020",max_dist=2)
testData2<-testData1[,4:6]-testData1[,1:3]
col<-col_by_value(testData2,col=colorRampPalette(c('green','black','red'))(1024),range=c(-2,2))
temp<-plot_pathway(testData2,type="bg",bg_col=col,text_col="white",magnify=1.2,species='hsa',
database_dir=system.file("extdata",package="KEGGprofile"),pathway_id="00020")

```

plot_pathway_cor

plot_pathway_cor

Description

The function will plot the correlation distributions for each enriched pathway (result from `find_enriched_pathway` function), and then Wilcoxon tests would be used to estimate the significance of correlations distribution between genes in each pathway and all genes.

Usage

```
plot_pathway_cor(gene_expr, kegg_enriched_pathway, groups = NULL,
side = c("both", "pos", "neg"), alternative = NULL)
```

Arguments

<code>gene_expr</code>	the matrix for gene expression, row.names should be NCBI gene ID, such as 67040, 93683
<code>kegg_enriched_pathway</code>	The returned value from <code>find_enriched_pathway</code> function, the enriched pathways.
<code>groups</code>	a character used to indicate expression values from different types of samples
<code>side</code>	a character string specifying the correlation directions interested, must be one of "both" (default), "pos" or "neg".
<code>alternative</code>	a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.

Value

p values for Wilcoxon tests in each pathway

Examples

```

data(pro_pho_expr)
data(pho_sites_count)
genes<-row.names(pho_sites_count)[which(pho_sites_count>=10)]
pho_KEGGresult<-find_enriched_pathway(genes,species='hsa')
result<-plot_pathway_cor(gene_expr=pro_pho_expr,kegg_enriched_pathway=pho_KEGGresult)

```

`plot_pathway_overall` *plot_pathway_overall*

Description

The function will plot the correlation distributions for each enriched pathway (result from `find_enriched_pathway` function), and then Wilcoxon tests would be used to estimate the significance of values distribution between genes in each pathway and all other genes.

Usage

```

plot_pathway_overall(gene_values, species = "hsa",
  pathwayNumInFigure = 5, rankByVar = colnames(gene_values)[1])

```

Arguments

<code>gene_values</code>	A data.frame or matrix with gene ID as rownames. Each column represent gene value in one condition.
<code>species</code>	the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc
<code>pathwayNumInFigure</code>	a integer specifying max number of (top) pathways in each direction in the box-plot.
<code>rankByVar</code>	a character string specifying variable (sample) name by which to rank the box-plot.

Value

p values for Wilcoxon tests in each pathway

Examples

```

data(pro_pho_expr)
data(pho_sites_count)
gene_values<-pro_pho_expr[row.names(pho_sites_count)[which(pho_sites_count>=10)],]
plot_pathway_overall(gene_values=gene_values[,1:3])

```

plot_profile	<i>plot_profile</i>
--------------	---------------------

Description

The function plot gene expression profiles on KEGG pathway maps

Usage

```
plot_profile(gene_expr, pathway_name, result_name = paste(pathway_name,
  "_profile_", type, ".png", sep = ""), KEGG_database, groups,
  bg_col = "white", text_col = "black", line_col,
  border_col = "grey", text_cex = 0.25, magnify = 1,
  type = c("lines", "bg"), pathway_min = 5,
  genes_kept = c("foldchange", "first", "random", "var", "abs"),
  species = "hsa", database_dir = getwd(), max_dist, lwd = 1.2,
  speciesRefMap = TRUE)
```

Arguments

gene_expr	the matrix for gene expression, row.names should be NCBI gene ID, such as 67040, 93683
pathway_name	the species id and KEGG pathway id, such as 'hsa00010'
result_name	the name of figure file generated by KEGGprofile. The default name is pathway_name+'_profile_'+type+'.png', such as 'hsa04110_profile_lines.png'
KEGG_database	the matrix returned by function parse_XMLfile, which contains genes in this pathway, and their names, locations etc
groups	a character used to indicate expression values from different types of samples
bg_col	background color for gene rectangles in the pathway map
text_col	the colors for text in the pathway map. A color matrix generated by function col_by_value can be used here
line_col	line color for expression in different samples in the pathway map, valid when type='lines'
border_col	border color for gene rectangles in the pathway map. A color matrix generated by function col_by_value can be used here
text_cex	cex for text in the pathway map. A color matrix generated by function col_by_value can be used here
magnify	the coefficient used to magnify the gene rectangles
type	the type of pathway map visulization, could be 'bg' or 'lines'. Default is 'bg'. See also 'Details'
pathway_min	The pathways with number of annotated genes less than pathway_min would be ignored
genes_kept	methods used for choosing genes when several genes corresponding to one location in pathway map. Default is 'foldchange', which kept the gene with largest fold changes. 'first' kept the first gene. 'random' chosed gene random. 'var' kept the gene with largest variation. 'abs' kept the gene with largest absolute value

species	the species id in KEGG database, 'hsa' means human, 'mmu' means mouse, 'rno' means rat, etc
database_dir	the directory where the XML files and png files are located
max_dist	The expression changes that represented by the distance from the bottom to the top of gene rectangle, valid when type='lines'. This param is used to ensure the dynamic changes of lines in different gene polygon represent equal variation. It would be calculated from the maximum changes of genes in this pathway by default. If max_dist=NA, then the lines would be plotted from top to bottom in each gene rectangle
lwd	The line width when type='lines'
speciesRefMap	Logical, use the species specific figure as reference map. if set as FALSE, the reference pathway figure without species information will be used

Details

There are two visualization methods to represent gene expression profiles: 'background' and 'lines'. The first one is applicable for analysis with only one sample or one type of data, which divides the gene polygon into several sub-polygons to represent different time points. And each sub-polygon has a specific background color to represent expression changes in that time point. The second method plots lines with different colors in the gene polygon to represent different samples or different types of data. The dynamic changes of lines mean the profiles of genes in different time points.

Value

a matrix containing genes mapped in this pathway, and their names, expressions

Examples

```
XML2database<-parse_XMLfile(pathway_id="04110",species="hsa",
database_dir=system.file("extdata",package="KEGGprofile"))
data(pro_pho_expr)
temp<-plot_profile(pro_pho_expr,pathway_name="hsa04110",KEGG_database=XML2database,
line_col=c("brown1","seagreen3"),groups=c(rep("Proteome ",6),rep("Phosphoproteome ",6)),
magnify=1.2,database_dir=system.file("extdata",package="KEGGprofile"),max_dist=5)
```

pro_pho_expr	<i>expression profiles in proteome and phosphoproteome</i>
--------------	--

Description

This data set is from a previously published data of proteome and phosphoproteome analysis in different cell phase. The column 1-6 are proteome data and column 7-12 are phosphoproteome data in this data.frame. The 6 time points are G1, G1/S, Early S, Late S, G2, Mitosis.

Usage

```
pro_pho_expr
```

Source

Olsen, J.V., et al. (2010) Quantitative phosphoproteomics reveals widespread full phosphorylation site occupancy during mitosis, *Sci Signal*, 3, ra3.

Index

`col_by_value`, [2](#), [11](#)
`convertId`, [3](#)

`download_KEGGfile`, [4](#), [8](#)
`download_latest_pathway`, [4](#)

`find_enriched_pathway`, [5](#)

`KEGGprofile` (`KEGGprofile-package`), [2](#)
`KEGGprofile-package`, [2](#)

`newIdMatrix`, [6](#)

`parse_XMLfile`, [7](#), [8](#)
`pho_sites_count`, [7](#)
`plot_pathway`, [8](#)
`plot_pathway_cor`, [9](#)
`plot_pathway_overall`, [10](#)
`plot_profile`, [8](#), [11](#)
`pro_pho_expr`, [12](#)