

Package ‘octad’

April 10, 2023

Title Open Cancer TherApeutic Discovery (OCTAD)

Version 1.0.0

Description OCTAD provides a platform for virtually screening compounds targeting precise cancer patient groups. The essential idea is to identify drugs that reverse the gene expression signature of disease by tamping down over-expressed genes and stimulating weakly expressed ones. The package offers deep-learning based reference tissue selection, disease gene expression signature creation, pathway enrichment analysis, drug reversal potency scoring, cancer cell line selection, drug enrichment analysis and in silico hit validation. It currently covers ~20,000 patient tissue samples covering 50 cancer types, and expression profiles for ~12,000 distinct compounds.

License Artistic-2.0

Encoding UTF-8

LazyData FALSE

Roxygen list(markdown = TRUE)

RoxygenNote 7.2.1

Depends R (>= 4.2.0), magrittr, dplyr, ggplot2, edgeR, RUVSeq, DESeq2, limma, rhdf5, foreach, Rfast, octad.db, stats, httr, ExperimentHub, AnnotationHub, Biobase, S4Vectors

Imports EDASeq, GSVA, data.table, htmlwidgets, plotly, reshape2, grDevices, utils

Suggests knitr, rmarkdown

VignetteBuilder knitr

biocViews Classification, GeneExpression, Pharmacogenetics, Pharmacogenomics, Software, GeneSetEnrichment

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

git_branch RELEASE_3_16

git_last_commit 9ea3c9a

git_last_commit_date 2022-11-01

Date/Publication 2023-04-10

Author E. Chekalin [aut, cre],
 S. Paithankar [aut],
 B. Zeng [aut],
 B. Glicksberg [ctb],
 P. Newbury [ctb],
 J. Xing [ctb],
 K. Liu [ctb],
 A. Wen [ctb],
 D. Joseph [ctb],
 B. Chen [aut]

Maintainer E. Chekalin <eygen.chekalin@gmail.com>

R topics documented:

computeCellLine	2
computeRefTissue	3
diffExp	5
loadOctadCounts	6
octad	7
octadDrugEnrichment	8
res_example	9
runsRGES	10
sRGES_example	11
topLineEval	12

Index	14
--------------	-----------

computeCellLine	<i>Compute Correlation between cell lines and vector of case ids.</i>
-----------------	-----------------------------------------------------------------------

Description

Select top CCLE cell lines sharing similar expression profiles with input case samples. Input case sample ids and output correlation scores for every cell line and/or output file. The results could be used for in-silico validation of predictions or used to weight cell lines in RGES computation. CellLineCorrelations.csv, correlation between CCLE cell lines and input disease samples.

Usage

```
computeCellLine(case_id = case_id, expSet = NULL, LINCS_overlaps = TRUE,
                source = c("octad.small", "octad.whole", "expSet"),
                file = NULL, output = TRUE,
                outputFolder = NULL)
```

Arguments

case_id	vector of ids from octad database. Ids can be obtained from phenoDF.
output	by default FALSE, if TRUE, file CellLineCorrelations.csv with results are produced in working directory.
outputFolder	Folder to store results.
LINCS_overlaps	vector of cell line ids from octad database. If TRUE, overlap with LINCS cells database will be performed
source	the file for the octad expression matrix. By default, set to octad.small to use only 978 landmark genes profiled in LINCS database. Use octad.whole option to compute DE on the whole transcriptome octad.counts.and.tpm.h5 file. The file should be present in the working directory or the whole path should be included. If source is set to 'side', the expSet matrix is estimated.
expSet	input expression matrix. By default set to NULL since the expSet is created based on cases, controls and source file.
file	if expSet='octad.whole', source path to expSet='octad.counts.and.tpm.h5' file is required if it is not in working directory. By default function seeks for the .h5 file in the working directory.

Value

topline	data.frame with row.names as cell line names and column medcor containing values for correlation between set of samples from case_id and cell lines.
---------	------------------------------------------------------------------------------------------------------------------------------------------------------

See Also

[runsRGES](#)

Examples

```
#load data.frame with samples included in the OCTAD database
phenoDF=get_ExperimentHub_data('EH7274')
HCC_primary=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&
sample.type == 'primary') #select data
case_id=HCC_primary$sample.id #select cases
cell_line_computed=computeCellLine(case_id=case_id,source='octad.small')
```

computeRefTissue	<i>Compute correlating reference control samples.</i>
------------------	-------------------------------------------------------

Description

Compute reference control samples from OCTAD database using precomputed EncoderDF models.

Usage

```
computeRefTissue(case_id = NULL, adjacent = FALSE, source = "octad",
  n_varGenes = 500, method = c("varGenes", 'random'), expSet = NULL,
  control_size = length(case_id),
  outputFolder = NULL, cor_cutoff = "0", output = TRUE)
```

Arguments

case_id	vector of cases used to compute references.
source	by default set octad to use autoencoder results for computation. Any other input like 'side' is expSet defined by users.
adjacent	by default set to FALSE. If TRUE, only tissue with sample.type 'adjacent' from phenoDF would be used instead of 'normal'.
expSet	input for expression matrix. By default NULL, since autoencoder results are used.
n_varGenes	number of genes used to select control cases.
method	one of two options is available. random will take a random number of samples from control subset and varGenes (default) will select control samples based on distance between cases and selected samples.
control_size	number of control samples to be selected.
outputFolder	path to output folder. By default, the function produces result files in working directory.
cor_cutoff	cut-off for correlation values, by default cor_cutoff='0'.
output	if TRUE, two output files are produced.

Value**Return**

control_id a vector of an appropriate set of control samples.

Besides, if output=TRUE, two files are created in the working directory:

case_normal_corMatrix.csv

contains pairwise correlation between case samples vs control samples.

case_normal_median_cor.csv

contains median correlation values with case samples for returned control samples.

See Also

[diffExp](#).

Examples

```
#select data
#load data.frame with samples included in the OCTAD database
phenoDF=get_ExperimentHub_data('EH7274')
HCC_primary=subset(phenoDF,cancer=='Liver Hepatocellular Carcinoma'&
sample.type == 'primary'&data.source == 'TCGA')
#select cases
case_id=HCC_primary$sample.id
#computing reference tissue, by default using small autoEncoder,
#but can use custom expression set,
#by default output=TRUE and outputFolder option is empty,
#which creates control corMatrix.csv to working directory
control_id=computeRefTissue(case_id,outputFolder='',output=TRUE,
expSet = "octad",control_size = 50)
```

diffExp

Compute differential expression

Description

Compute differential expression for case vs control samples. Will produce the file `computedEmpGenes.csv` listing empirically differentially expressed genes used for RNA-Seq normalization.

Usage

```
diffExp(case_id = NULL, control_id = NULL, source = "octad.small",
file = "octad.counts.and.tpm.h5", normalize_samples = TRUE,
k = 1, expSet = NULL, n_topGenes = 500,
DE_method = c("edgeR", 'DESeq2', 'wilcox', 'limma'),
output = FALSE, outputFolder = NULL, annotate = TRUE)
```

Arguments

<code>case_id</code>	vector of cases used for differential expression.
<code>control_id</code>	vector of controls used for differential expression.
<code>source</code>	the file for the octad expression matrix. By default, set to <code>octad.small</code> to use only 978 landmark genes profiled in LINCS database. Use <code>octad.whole</code> option to compute DE on the whole transcriptome <code>octad.counts.and.tpm.h5</code> file. The file should be present in the working directory or the whole path should be included. If <code>source</code> is set to <code>'side'</code> , the <code>expSet</code> matrix is estimated.
<code>expSet</code>	input expression matrix. By default set to <code>NULL</code> since the <code>expSet</code> is created based on cases, controls and source file.
<code>file</code>	if <code>expSet='octad.whole'</code> , source path to <code>expSet='octad.counts.and.tpm.h5'</code> file is required if it is not in working directory. By default function seeks for the <code>.h5</code> file in the working directory.

normalize_samples	if TRUE, RUVSeq normalization is applied to either EdgeR or DESeq. No normalization needed for limma+voom.
k	either k=1 (by default), k=2 or k=3, number of factors used in model matrix construction in RUVSeq normalization if normalize_samples=TRUE.
n_topGenes	number of empirically differentially expressed genes estimated for RUVSeq normalization. Default is 5000.
DE_method	edgeR, DESeq2, limma or wilcox DE analysis.
output	if TRUE, output files is produced.
outputFolder	path to output folder. By default, the function produces result files in working directory.
annotate	if TRUE, annotation by ENSEMBL gene is performed. If TRUE, make sure row.names of the custom input contain ensembl gene ids.

Value

res	data.frame with list of differentially expressed genes.
computedEmpGenes.csv	data.frame listing empirically differentially expressed genes used for RNA-Seq normalization.

See Also

[computeRefTissue](#), [runsRGES](#).

Examples

```
#load data.frame with samples included in the OCTAD database
phenoDF=get_ExperimentHub_data('EH7274')
HCC_primary=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&
sample.type == 'primary') #select data
case_id=HCC_primary$sample.id #select cases
HCC_adjacent=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&
sample.type == 'adjacent'&data.source == 'TCGA') #select data
control_id=HCC_adjacent$sample.id #select cases
res=diffExp(case_id,control_id,source='octad.small',output=FALSE)
```

loadOctadCounts	<i>Load octad expression data</i>
-----------------	-----------------------------------

Description

Create TPM or count expression matrix for the selected samples from OCTAD.

Usage

```
loadOctadCounts(sample_vector='',type='tpm',file='')
```

Arguments

`sample_vector` vector of samples to be selected. Use `phenoDF` object for sample id selection.
`type` either `tpm` (default) or `counts` to be returned.
`file` full path to `octad.counts.and.tpm.h5` file.

Value

`exprData` matrix with either `log2` corrected counts or `tmp` matrix for selected samples.

See Also

[diffExp](#).

Examples

```
#load data.frame with samples included in the OCTAD database
phenoDF=get_ExperimentHub_data('EH7274')
#load expression data for raw counts or tpm values.
HCC_primary=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&
sample.type == 'primary') #select data
#case_id=HCC_primary$sample.id #select cases
#expression_tmp=loadOctadCounts(case_id,type='tpm',
#file='octad.counts.and.tpm.h5')
#expression_log2=loadOctadCounts(case_id,type='counts',
#file='octad.counts.and.tpm.h5')
```

octad

Open Cancer TherApeutic Discovery (OCTAD) database package

Description

Open Cancer TherApeutic Discovery (OCTAD) package implies sRGES approach for the drug discovery. The essential idea is to identify drugs that reverse the gene expression signature of a disease by tamping down over-expressed genes and stimulating weakly expressed ones. The following package contains all required precomputed data for whole OCTAD pipeline computation.

Details

The main functions are:

- `computeRefTissue` - Compute reference control samples from OCTAD database using pre-computed `EncoderDF` models.
- `diffExp` - Compute differential expression for case vs control samples. Will produce the file `computedEmpGenes.csv` listing empirically differentially expressed genes used for RNA-Seq normalization.

- `runsRGES` - Compute sRGES, a score indicating the reversal potency of each drug. It first computes RGES (Reverse Gene Expression Score) for individual instances and then summarizes RGES of individual drugs (one drug may have multiple instances under different treatment conditions).
- `computeCellLine` - Compute Correlation between cell lines and vector of case ids.
- `topLineEval` - Evaluate predictions using pharmacogenomics data. Given a cell line, the function computes the correlation between sRGES and drug sensitivity data taken from CTRP. A higher correlation means a better prediction. The cell line could be computed from `computeCellLine`.
- `octadDrugEnrichment` - Perform enrichment analysis of drug hits based on chemical structures, drug-targets, and pharmacological classifications. An enrichment score calculated using ssGSEA and a p-value computed through a permutation test are provided.

For detailed information on usage, see the package vignette, by typing `vignette('octad')`, or the workflow linked to on the first page of the vignette.

The code can be viewed at the GitHub repository, which also lists the contributor code of conduct:

<https://github.com/Bin-Chen-Lab/OCTAD>

References

Zeng, B., Glicksberg, B.S., Newbury, P., Chekalin, E., Xing, J., Liu, K., Wen, A., Chow, C. and Chen, B., 2021. OCTAD: an open workspace for virtually screening therapeutics targeting precise cancer patient groups using gene expression features. *Nature protocols*, 16(2), pp.728-753. <https://www.nature.com/articles/s41596-020-00430-z>

octadDrugEnrichment *Compute Drug enrichment*

Description

Perform enrichment analysis of drug hits based on chemical structures, drug-targets, and pharmacological classifications. An enrichment score calculated using ssGSEA and a p-value computed through a permutation test are provided.

Usage

```
octadDrugEnrichment(sRGES = NULL, target_type = "chembl_targets",
  enrichFolder = "enrichFolder", outputFolder = NULL, outputRank = FALSE)
```

Arguments

sRGES	sRGES data frame produced by runsRGES.
target_type	one or several of 'chembl_targets', 'mesh', 'ChemCluster' databases selected. By default only 'chembl_targets' will be used.
enrichFolder	folder to store output.
outputFolder	path where to store enrichFolder, in case of NULL will be stored in work directory.
outputRank	output detailed rank if TRUE, write sRGES for selected target as vcf.

Value

Following files are created: `enriched*_targets.csv` and `top_enriched*_*_targets.pdf`. In the case of chemical structural analysis, additional files are created: `*drugstructureClusters.csv` and `*misc.csv`. The results provide useful information for following candidate selection and experimental design. For example, if two structurally similar drugs are both predicted as top hits, the chance of each drug as a true positive is high.

`exprData` matrix with either log2 corrected counts or tmp matrix for selected samples.

See Also

[runsRGES](#)

Examples

```
data("sRGES_example", package='octad') #load example sRGES
#run drug enrichment
octadDrugEnrichment(sRGES = sRGES_example, target_type = c('chembl_targets'))
```

res_example	<i>Differential expression example for HCC vs adjacent liver tissue computed in diffExp() function</i>
-------------	--------------------------------------------------------------------------------------------------------

Description

Differential expression example for HCC vs adjacent liver tissue computed in `diffExp()` function

Usage

```
res_example
```

Format

A data.frame with 963 rows and 18 variables:

identifier Ensg ID
log2FoldChange Log2 fold-change
logCPM log CPM value
LR LR value
pvalue p.value
padj FDR
tax_id taxon id
GeneID Gene id
LocusTag Locus tag

chromosome Chromosome
map_location Chromosome location
description Full gene name
type type of gene
Symbol_autho HGNC symbol
other Gene function

Details

To generate this dataset use the following code from the octad package #load data.frame with samples included in the OCTAD database.

```
phenoDF=.eh[['EH7274']]
#select data
HCC_primary=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&sample.type=='primary')
#select cases
case_id=HCC_primary$sample.id
control_id=subset(phenoDF,biopsy.site=='LIVER'&sample.type=='normal')$sample.id[1:50]
res=diffExp(case_id,control_id,source='octad.small',output=FALSE)
```

runsRGES

Compute sRGES

Description

Compute sRGES, a score indicating the reversal potency of each drug. It first computes RGES (Reverse Gene Expression Score) for individual instances and then summarizes RGES of individual drugs (one drug may have multiple instances under different treatment conditions).

Usage

```
runsRGES(dz_signature=NULL,choose_fda_drugs = FALSE,max_gene_size=500,
cells=NULL,output=FALSE,outputFolder='',weight_cell_line=NULL,permutations=10000)
```

Arguments

dz_signature disease signature. Make sure input data frame has a gene Symbol column, otherwise an error is produced. It must be an UPPERCASE gene symbol.

choose_fda_drugs
if TRUE, only FDA approved drugs are used.

max_gene_size maximum number of disease genes used for drug prediction. By default 50 for each side (up/down).

cells cell ids in lincs_sig_info file used for prediction. By default, all cell lines are used.

weight_cell_line	by default NULL, if !NULL, an output object from computeCellLine is estimated (see example).
permutations	number of permutations, by default 10000.
output	if TRUE, output files is produced.
outputFolder	folder path to store drug results, by default write results to working directory.

Value

The function returns RGES data.frame containing scores and p.values for every instance. data.frame contains drug id in pert_iname collumn, n contains the number of instances for this drug, mean, median and sd of sRGES RGES sores.

Besides, a number of additional files in the sourced directory:

dz_sig_used.csv	contains genes in the disease signature used for computing reverse gene expression scores.
sRGES.csv	contains the same data as returned data.frame.
all__lincs_score.csv	includes information of RGES.

See Also

[diffExp](#), [octadDrugEnrichment](#), [computeCellLine](#), [topLineEval](#)

Examples

```
#load differential expression example for HCC
#vs adjacent liver tissue computed in diffExp() function
data("res_example",package='octad')
res_example=subset(res_example,abs(log2FoldChange)>1&padj<0.001)[1:10,]
#run sRGES computation
#sRGES=runsRGES(dz_signature=res_example,max_gene_size=100,permutations=1000,output=FALSE)
```

sRGES_example	<i>Data of computed example sRGES for HCC vs liver adjacent tissues on octad.small dataset</i>
---------------	------------------------------------------------------------------------------------------------

Description

Data of computed example sRGES for HCC vs liver adjacent tissues on octad.small dataset

Usage

```
sRGES_example
```

Format

A tibble with 12,442 rows and 6 variables:

pert_iname dbl Year price was recorded
mean mean sRGES for obtained drug if n>1
n times this drug was obtained
median median sRGES for drug if n>1
sd standart deviation for obtained drug if n>1
sRGES sRGES score of the drug

Details

To generate this dataset use the following code from the octad package load differential expression example for HCC vs adjacent liver tissue computed in diffExp() function from res_example.

```
data('res_example', package='octad.db')
res=subset(res_example, abs(log2FoldChange)>1&padj<0.001) #load example expression dataset
sRGES=runsRGES(res, max_gene_size=100, permutations=1000, output=FALSE)
```

topLineEval

Evaluate cell lines

Description

Evaluate predictions using pharmacogenomics data. Given a cell line, the function computes the correlation between sRGES and drug sensitivity data taken from CTRP. A higher correlation means a better prediction. The cell line could be computed from computeCellLine.

Usage

```
topLineEval(topline=NULL, mysRGES=NULL, outputFolder="")
```

Arguments

topline list of cell lines to be used for prediction.
mysRGES sRGES data.frame produced by runsRGES.
outputFolder Path to store results.

Value

The function produces 3 files in the output directory:

CellLineEval*_drug_sensitivity_insilico_results.txt

with drug sensitivity information.

*_auc_insilico_validation.html

correlation between drug AUC and sRGES in a related cell line.

*_ic50_insilico_validation.html

correlation between drug IC50 and sGRES in a related cell line.

See Also

[runsRGES](#)

Examples

```
#load example sRGES computed by runsRGES() function for HCC
#vs liver adjacent tissues on octad.small dataset
data("sRGES_example",package='octad') #load example sRGES
#Pick up cell lines
topLineEval(topline = 'HEPG2',mysRGES = sRGES_example,outputFolder=tempdir())
```

Index

- * **computeRefTissue**
 - computeRefTissue, 3
- * **datasets**
 - res_example, 9
 - sRGES_example, 11
- * **diffExp**
 - diffExp, 5
 - loadOctadCounts, 6
- * **octadDrugEnrichment**
 - computeCellLine, 2
 - octadDrugEnrichment, 8
 - topLineEval, 12
- * **sRGES**
 - runsRGES, 10

computeCellLine, 2, 8, 11
computeCellLine (computeCellLine), 2
computeRefTissue, 3, 6, 7

diffExp, 4, 5, 7, 11

loadOctadCounts, 6

octad, 7
octadDrugEnrichment, 8, 8, 11

res_example, 9
runsRGES, 3, 6, 8, 9, 10, 13

sRGES_example, 11

topLineEval, 8, 11, 12