

Package ‘HPAanalyze’

October 16, 2019

Type Package

Title Retrieve and analyze data from the Human Protein Atlas

Version 1.2.1

Date 2019-06-19

Description Provide functions for retrieving, exploratory analyzing and visualizing the Human Protein Atlas data.

Depends R (>= 3.5.0)

License GPL-3 + file LICENSE

RoxygenNote 6.1.1

Encoding UTF-8

Imports dplyr, openxlsx, ggplot2, readr, tibble, xml2, tidyr, stats, utils, hpar, gridExtra

Suggests knitr, rmarkdown, devtools, BiocStyle

VignetteBuilder knitr

LazyData true

biocViews Proteomics, CellBiology, Visualization, Software

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

git_branch RELEASE_3_9

git_last_commit 75c6aa5

git_last_commit_date 2019-06-19

Date/Publication 2019-10-15

Author Anh Nhat Tran [aut, cre]

Maintainer Anh Nhat Tran <trannhatanh89@gmail.com>

R topics documented:

hpaDownload	2
hpaExport	3
hpaSubset	4
hpaVis	5
hpaVisPatho	6
hpaVisSubcell	7
hpaVisTissue	8

hpaXml	9
hpaXmlAntibody	10
hpaXmlGet	11
hpaXmlProtClass	12
hpaXmlTissueExpr	12
hpaXmlTissueExprSum	13
hpa_downloaded_histology_v18	14

Index	15
--------------	-----------

hpaDownload	<i>Download datasets</i>
-------------	--------------------------

Description

Download the latest version of HPA datasets and import them in R. It is recommended to only download the datasets you need, as some of them may be very big.

Usage

```
hpaDownload(downloadList = "histology", version = "latest")
```

Arguments

`downloadList` A vector or string indicate which datasets to download. Possible value:

- 'Normal tissue'
- 'Pathology'
- 'Subcellular location'
- 'RNA tissue'
- 'RNA cell line'
- 'RNA transcript tissue'
- 'RNA transcript cell line'
- 'all': download everything
- 'histology': same as c('Normal tissue', 'Pathology', 'Subcellular location')
- 'rna': same as c('RNA tissue', 'RNA cell line')
- 'isoform': same as c('RNA transcript tissue', 'RNA transcript cell line')

See <https://www.proteinatlas.org/about/download> for more information.

`version` A string indicate which version to be downloaded. Possible value:

- 'hpar': Load data from the 'hpar' package. Does not contain isoform data.
- 'latest': Download latest version. Require Internet connection.
- 'example': Load the example dataset from 'HPA analyze' ('hpa_downloaded_histology_v18'). Does not contain rna or isoform data.
- 'v?' with '?' is a integer: Download a specific version of the dataset. For example: 'v18' download version 18. Currently support version 17 and above. Require Internet connection.

Value

This function will return a list of tibbles corresponding to requested datasets.

See Also

Other downloadable datasets functions: [hpaExport](#), [hpaSubset](#)

Examples

```
downloadedData <- hpaDownload(downloadList='all', version='example')
summary(downloadData)
```

hpaExport	<i>Export the subset data</i>
-----------	-------------------------------

Description

Export the list object generated by `hpaSubset()` into `xlsx` format. Due to the size of some HPA datasets, as well as the limitation of the output format, exporting the full datasets generated by `hpaDownload()` is not recommended.

Usage

```
hpaExport(data, fileName, fileType = "xlsx")
```

Arguments

<code>data</code>	Input the list object generated by <code>hpaSubset()</code>
<code>fileName</code>	A string indicate the desired output file name. Do not include file extension such as <code>'.xlsx'</code> .
<code>fileType</code>	The format as which the data will be exported. Choose one of these options: <code>'xlsx'</code> , <code>'csv'</code> and <code>'tsv'</code> .

Value

- `'xlsx'`: return one `.xlsx` file named `'fileName.xlsx'`. One individual sheet for each dataset in the input list object.
- `'csv'`: return `.csv` files, one for each dataset in the input list object, named `'fileName_datasetName.csv'`
- `'tsv'`: return `.tsv` files, one for each dataset in the input list object, named `'fileName_datasetName.tsv'`

See Also

Other downloadable datasets functions: [hpaDownload](#), [hpaSubset](#)

Examples

```
downloadedData <- hpaDownload(downloadList='all', version='example')
geneList <- c('TP53', 'EGFR')
tissueList <- c('breast', 'cerebellum', 'skin 1')
cancerList <- c('breast cancer', 'glioma', 'melanoma')

subsetData <- hpaSubset(data=downloadedData,
                        targetGene=geneList,
                        targetTissue=tissueList,
                        targetCancer=cancerList)

hpaExport(data=subsetData,
           fileName='TP53_EGFR_in_tissue_cancer.xlsx',
           fileType='xlsx')
```

hpaSubset

Subset downloaded data

Description

hpaSubset() subsets data by gene name, tissue, cell type, cancer and/or cell line. The input is the list object generated by hpaDownload() or as the output of another hpaSubset(). Use hpaListParam() to see the list of available parameters for a specific list object. Will not work on isoform data.

hpaListParam() list available variables in downloaded data that can be used as parameters to subset the data via hpaSubset(). This function work with the data object generated by hpaDownload() or a previous call of hpaSubset().

Usage

```
hpaSubset(data = NULL, targetGene = NULL, targetTissue = NULL,
           targetCellType = NULL, targetCancer = NULL, targetCellLine = NULL)
```

```
hpaListParam(data = NULL)
```

Arguments

data	Input the list object generated by hpaDownload() or hpaSubset()
targetGene	Vector of strings of HGNC gene symbols. It will be used to subset every dataset in the list object. You can also mix HGNC gene symbols and ensemnl ids (start with ENSG) and they will be converted to HGNC gene symbols.
targetTissue	Vector of strings of normal tissues. Will be used to subset the normal_tissue and rna_tissue dataset.
targetCellType	Vector of strings of normal cell types. Will be used to subset the normal_tissue dataset.
targetCancer	Vector of strings of cancer types. Will be used to subset the pathology dataset.
targetCellLine	Vector of strings of cell lines. Will be used to subset the rna_cell_line dataset.

Value

hpaSubset will return a list of tibbles as the result of subsetting, depending on the input data.

The output of hpaListParam() is a list of vectors containing all subset parameter for the downloaded data.

See Also

Other downloadable datasets functions: [hpaDownload](#), [hpaExport](#)

Examples

```
downloadedData <- hpaDownload(downloadList='all', version='example')
geneList <- c('TP53', 'EGFR')
tissueList <- c('breast', 'cerebellum', 'skin 1')
cancerList <- c('breast cancer', 'glioma', 'melanoma')

subsetData <- hpaSubset(data=downloadedData,
                        targetGene=geneList,
                        targetTissue=tissueList,
                        targetCancer=cancerList)

downloadedData <- hpaDownload(downloadList='all', version='example')
params <- hpaListParam(data=downloadedData)
params$normal_tissue
```

hpaVis

Visualize data in one function

Description

This function is an universal visualization function that allow calling other hpaVis functions via a single function call. By default, this function will use the dataset bundled with HPAanalyze, and provide a grid of all available plots. The types of plots in the output can be specified via the visType argument. If only one plot type is specified, this function will return the exact same output as the specific hpaVis function used to create the plot.

Usage

```
hpaVis(data = NULL, targetGene = NULL, targetTissue = NULL,
       targetCellType = NULL, targetCancer = NULL, visType = c("Tissue",
       "Patho", "Subcell"), color = c("#ffffb2", "#fecc5c", "#fd8d3c",
       "#e31a1c"), customTheme = FALSE, ...)
```

Arguments

data	Input the list object generated by hpa_download() or hpa_subset(). By default this function use the example dataset bundled with HPAanalyze, which is from version 18 of the HPA program.
targetGene	Vector of strings of HGNC gene symbols. By default it is set to c('TP53', 'EGFR', 'CD44', 'PTEN', '...'). You can also mix HGNC gene symbols and ensembl ids (start with ENSG) and they will be converted to HGNC gene symbols.

targetTissue	Vector of strings of normal tissue names. By default it is set to "breast".
targetCellType	Vector of strings of normal cell types. By default includes all available cell types in the target tissues.
targetCancer	Vector of strings of normal tissues. By default it is set to "breast cancer".
visType	Vector of strings indicating which plots will be generated. Currently available values are "all", "Tissue", "Patho", "Cancer", "Subcell".
color	Vector of 4 colors used to depict different expression levels.
customTheme	Logical argument. If TRUE, the function will return a barebone ggplot2 plot to be customized further.
...	Additional arguments to be passed downstream to other hpaVis functions being called behind the scene. These arguments includes targetTissue, targetCellType, targetCancer. See documentation for individual hpaVis functions for more information.

Value

If multiple visType is chosen, this function will return multiple graphs in one panel. If only one visType is chosen, this function will return a ggplot2 plot object, which can be further modified if desirable. See help file for each of the hpaVis function for more information about individual graphs.

See Also

[hpaDownload](#), [hpaSubset](#)

Other visualization functions: [hpaVisPatho](#), [hpaVisSubcell](#), [hpaVisTissue](#)

Examples

```
hpaVis()
```

hpaVisPatho

Visualize pathology data

Description

Visualize the expression of genes of interest in each cancer.

Usage

```
hpaVisPatho(data = NULL, targetGene = NULL, targetCancer = NULL,
  color = c("#ffffb2", "#fecc5c", "#fd8d3c", "#e31a1c"),
  customTheme = FALSE)
```

Arguments

data	Input the list object generated by <code>hpa_download()</code> or <code>hpa_subset()</code> . Require the pathology dataset. Use HPA version 18 by default.
targetGene	Vector of strings of HGNC gene symbols. By default it is set to <code>c('TP53', 'EGFR', 'CD44', 'PTEN', 'IDH1', 'IDH2', 'CYCS')</code> . You can also mix HGNC gene symbols and ensembl ids (start with ENSG) and they will be converted to HGNC gene symbols.
targetCancer	Vector of strings of normal tissues. The function will plot all available cancer by default.
color	Vector of 4 colors used to depict different expression levels.
customTheme	Logical argument. If TRUE, the function will return a barebone ggplot2 plot to be customized further.

Value

This function will return a ggplot2 plot object, which can be further modified if desirable. The pathology data is visualized as multiple bar graphs, one for each type of cancer. For each bar graph, x axis contains the inquired protein and y axis contains the proportion of patients.

See Also

Other visualization functions: [hpaVisSubcell](#), [hpaVisTissue](#), [hpaVis](#)

Examples

```
data("hpa_downloaded_histology_v18")
geneList <- c('TP53', 'EGFR', 'CD44', 'PTEN', 'IDH1', 'IDH2', 'CYCS')
cancerList <- c('breast cancer', 'glioma', 'melanoma')

## A typical function call
hpaVisPatho(data=hpa_downloaded_histology_v18,
            targetGene=geneList)
```

hpaVisSubcell

Visualize subcellular location data

Description

Visualize the the confirmed subcellular locations of genes of interest.

Usage

```
hpaVisSubcell(data = NULL, targetGene = NULL,
              reliability = c("enhanced", "supported", "approved", "uncertain"),
              color = c("#ffffb2", "#e31a1c"), customTheme = FALSE)
```

Arguments

data	Input the list object generated by <code>hpa_download()</code> or <code>hpa_subset()</code> . Require the <code>subcellular_location</code> dataset. Use HPA version 18 by default.
targetGene	Vector of strings of HGNC gene symbols. By default it is set to <code>c('TP53', 'EGFR', 'CD44', 'PTEN')</code> . You can also mix HGNC gene symbols and ensembl ids (start with <code>ENSG</code>) and they will be converted to HGNC gene symbols.
reliability	Vector of string indicate which reliability scores you want to plot. The default is everything <code>c("enhanced", "supported", "approved", "uncertain")</code> .
color	Vector of 2 colors used to depict if the protein expresses in a location or not.
customTheme	Logical argument. If <code>TRUE</code> , the function will return a barebone <code>ggplot2</code> plot to be customized further.

Value

This function will return a `ggplot2` plot object, which can be further modified if desirable. The subcellular location data is visualized as a tile graph, in which the x axis includes the inquired proteins and the y axis contain the subcellular locations.

See Also

Other visualization functions: [hpaVisPatho](#), [hpaVisTissue](#), [hpaVis](#)

Examples

```
data("hpa_downloaded_histology_v18")
geneList <- c('TP53', 'EGFR', 'CD44', 'PTEN', 'IDH1', 'IDH2', 'CYCS')

## A typical function call
hpaVisSubcell(data=hpa_downloaded_histology_v18,
              targetGene=geneList)
```

hpaVisTissue

Visualize tissue data

Description

Visualize the expression of protein of interest in each target tissue by cell types.

Usage

```
hpaVisTissue(data = NULL, targetGene = NULL, targetTissue = NULL,
             targetCellType = NULL, color = c("#ffffb2", "#fecc5c", "#fd8d3c",
             "#e31a1c"), customTheme = FALSE)
```


Arguments

inputXml	Input can be either one Ensembl gene id (start with ENSG) or a imported XML object resulting from a hpaXmlGet() function call. You can also use HGNC gene symbol and it will be converted to ensembl id.
extractType	A vector of strings indicate which information is desired for extraction. By default this function will call all hpaXml functions available. Other options are 'ProtClass', 'TissueExprSum', 'Antibody', 'TissueExpr'.
...	Additional arguments to be passed downstream to other hpaXml functions being called behind the scene. See help files of other hpaXml functions for more information.

Value

This function returns a list. Each element of the list is information extracted from the XML file specified using other hpaXml functions. See help file for each XML function for more information.

See Also

Other xml functions: [hpaXmlAntibody](#), [hpaXmlGet](#), [hpaXmlProtClass](#), [hpaXmlTissueExprSum](#), [hpaXmlTissueExpr](#)

Examples

```
hpaXml(inputXml='ENSG00000131979', extractType=c('ProtClass', 'TissueExprSum', 'Antibody'))
```

hpaXmlAntibody	<i>Extract antibody information</i>
----------------	-------------------------------------

Description

Extract information about the antibodies used for a specific protein.

Usage

```
hpaXmlAntibody(importedXml)
```

Arguments

importedXml	Input an xml document object resulted from a hpaXmlGet() call.
-------------	--

Value

This function returns a tibble of 4 columns, containing information about the antibodies used in the project for the inquired protein: id, releaseDate, releaseVersion, and RRID.

See Also

Other xml functions: [hpaXmlGet](#), [hpaXmlProtClass](#), [hpaXmlTissueExprSum](#), [hpaXmlTissueExpr](#), [hpaXml](#)

Examples

```
## Not run:
GCH1xml <- hpaXmlGet('ENSG00000131979')
hpaXmlAntibody(GCH1xml)

## End(Not run)
```

hpaXmlGet

Download and import xml file

Description

Download and import individual xml file for a specified protein. This function calls `xml2::read_xml()` under the hood.

Usage

```
hpaXmlGet(targetEnsemblId, version = "latest")
```

Arguments

targetEnsemblId

A string of one ensembl ID, start with ENSG. For example 'ENSG00000131979'. You can also use HGNC gene symbol and it will be converted to ensembl id.

version

A string indicate which version to be downloaded. Possible value:

- 'latest': Download latest version.
- 'v?' with '?' is a integer: Download a specific version of the dataset. For example: 'v18' download version 18. Currently support version 13 and above.

Value

This function return an object of class "xml_document" "xml_node" containing the content of the imported XML file. (See documentations for package xml2 for more information.)

See Also

Other xml functions: [hpaXmlAntibody](#), [hpaXmlProtClass](#), [hpaXmlTissueExprSum](#), [hpaXmlTissueExpr](#), [hpaXml](#)

Examples

```
## Not run:
GCH1xml <- hpaXmlGet('ENSG00000131979')

## End(Not run)
```

hpaXmlProtClass	<i>Extract protein classes</i>
-----------------	--------------------------------

Description

Extract protein class information from imported xml document resulted from `hpaXmlGet()`.

Usage

```
hpaXmlProtClass(importedXml)
```

Arguments

`importedXml` Input an xml document object resulted from a `hpaXmlGet()` call.

Value

This function return a tibble of 4 columns.

See Also

Other xml functions: [hpaXmlAntibody](#), [hpaXmlGet](#), [hpaXmlTissueExprSum](#), [hpaXmlTissueExpr](#), [hpaXml](#)

Examples

```
## Not run:  
GCH1xml <- hpaXmlGet('ENSG00000131979')  
hpaXmlProtClass(GCH1xml)  
  
## End(Not run)
```

hpaXmlTissueExpr	<i>Extract tissue expression details</i>
------------------	--

Description

Extract tissue expression information for each sample and url to download images from imported xml document resulted from `hpaXmlGet()`.

Usage

```
hpaXmlTissueExpr(importedXml)
```

Arguments

`importedXml` Input an xml document object resulted from a `hpaXmlGet()` call.

Value

This function returns a list of tibbles, each for an antibody. Each tibble contains information about all individual samples and their staining. Due to the variation in amount of information available for these samples, the number of columns differs, but the tibble essentially includes: patientId, age, sex, staining, intensity, quantity, location, imageUrl, snomedCode, and tissueDescription. The last two items may have more than one column each.

See Also

Other xml functions: [hpaXmlAntibody](#), [hpaXmlGet](#), [hpaXmlProtClass](#), [hpaXmlTissueExprSum](#), [hpaXml](#)

Examples

```
## Not run:
GCH1xml <- hpaXmlGet('ENSG00000131979')
hpaXmlTissueExpr(GCH1xml)

## End(Not run)
```

hpaXmlTissueExprSum *Extract tissue expression and download images*

Description

Extract tissue expression information and url to download images from imported xml document resulted from `hpaXmlGet()`.

Usage

```
hpaXmlTissueExprSum(importedXml, downloadImg = FALSE)
```

Arguments

importedXml	Input an xml document object resulted from a <code>hpaXmlGet()</code> call.
downloadImg	Logical argument. The function will download all image from the extracted urls into the working folder.

Value

This function return a list consists of a summary string, which is a very brief description of the protein, and a tibble of 2 columns: tissue (name of tissue available) and imageUrl (link to download the perspective image)

See Also

Other xml functions: [hpaXmlAntibody](#), [hpaXmlGet](#), [hpaXmlProtClass](#), [hpaXmlTissueExpr](#), [hpaXml](#)

Examples

```
## Not run:
GCH1xml <- hpaXmlGet('ENSG00000131979')
hpaXmlTissueExprSum(GCH1xml)

## End(Not run)
```

```
hpa_downloaded_histology_v18
  HPA histology dataset
```

Description

Dataset downloaded with `hpaDownload('histology', version = 'v18')`.

Usage

```
hpa_downloaded_histology_v18
```

Format

A list of 3 tibbles

normal_tissue Normal tissue IHC data

pathology Cancer IHC data

subcellular_location Subcellular location IF data

Details

Links to original data:

- https://v18.proteinatlas.org/download/normal_tissue.tsv.zip
- <https://v18.proteinatlas.org/download/pathology.tsv.zip>
- https://v18.proteinatlas.org/download/subcellular_location.tsv.zip

See Also

[hpaDownload](#)

Examples

```
# load data
data("hpa_downloaded_histology_v18")

# access tibbles
normal_tissue_data <- hpa_downloaded_histology_v18$normal_tissue
cancer_data <- hpa_downloaded_histology_v18$pathology
subcell_location_data <- hpa_downloaded_histology_v18$subcellular_location
```

Index

*Topic **datasets**

hpa_downloaded_histology_v18, [14](#)

hpa_downloaded_histology_v18, [14](#)

hpaDownload, [2](#), [3](#), [5](#), [6](#), [14](#)

hpaExport, [3](#), [3](#), [5](#)

hpaListParam (hpaSubset), [4](#)

hpaSubset, [3](#), [4](#), [6](#)

hpaVis, [5](#), [7–9](#)

hpaVisPatho, [6](#), [6](#), [8](#), [9](#)

hpaVisSubcell, [6](#), [7](#), [7](#), [9](#)

hpaVisTissue, [6–8](#), [8](#)

hpaXml, [9](#), [10–13](#)

hpaXmlAntibody, [10](#), [10](#), [11–13](#)

hpaXmlGet, [10](#), [11](#), [12](#), [13](#)

hpaXmlProtClass, [10](#), [11](#), [12](#), [13](#)

hpaXmlTissueExpr, [10–12](#), [12](#), [13](#)

hpaXmlTissueExprSum, [10–13](#), [13](#)