

Package ‘ArrayExpress’

March 29, 2021

Title Access the ArrayExpress Microarray Database at EBI and build Bioconductor data structures: ExpressionSet, AffyBatch, NChannelSet

Version 1.50.0

Author Audrey Kauffmann, Ibrahim Emam, Michael Schubert

Maintainer Suhaib Mohammed <suhaib@ebi.ac.uk>

Depends R (>= 2.9.0), Biobase (>= 2.4.0)

Imports XML, oligo, limma

Description Access the ArrayExpress Repository at EBI and build Bioconductor data structures: ExpressionSet, AffyBatch, NChannelSet

License Artistic-2.0

biocViews Microarray, DataImport, OneChannel, TwoChannel

NeedsCompilation no

Suggests affy

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

git_branch RELEASE_3_12

git_last_commit 3a8e996

git_last_commit_date 2020-10-27

Date/Publication 2021-03-29

R topics documented:

ae2bioc	2
ArrayExpress	3
extract.zip	4
getAE	4
getcolproc	5
getcolraw	6
procset	6
queryAE	7
Index	8

`ae2bioc`*Convert MAGE-TAB files from raw data into a Bioconductor object*

Description

`ae2bioc` converts local MAGE-TAB files into a `AffyBatch`, an `ExpressionSet` or a `NChannelSet`.

Usage

```
ae2bioc(mageFiles, dataCols = NULL, drop = TRUE)
```

Arguments

<code>mageFiles</code>	A list as given from getAE function. Containing the following elements: rawFiles all the expression files to use to create the object. The content of the raw.zip MAGE-TAB file. sdrf the name of the sdrf file from MAGE-TAB. idf the name of the idf file from MAGE-TAB. adf the name of the adf file from MAGE-TAB. path is the name of the directory containing these files.
<code>dataCols</code>	by default, the columns are automatically selected according to the scanner type. If the scanner is unknown or if the user wants to use different columns than the default, the argument 'dataCols' can be set. For two colour arrays it must be a list with the fields 'R', 'G', 'Rb' and 'Gb' giving the column names to be used for red and green foreground and background. For one colour arrays, it must be a character string with the column name to be used. These column names must correspond to existing column names of the expression files.
<code>drop</code>	if TRUE and only one platform in series, the platform name will be dropped.

Value

An object of class [AffyBatch](#), [ExpressionSet](#) or [NChannelSet](#) with the raw expression values in the 'assayData' of the object, the information contained in the sdrf file in the 'phenoData', the adf file content in the 'featureData' and the idf file content in the 'experimentData'.

If several array designs are used in the dataset, the output is a list with an object for each array design.

Author(s)

Ibrahim Emam

Maintainer: <iemam@ebi.ac.uk>

See Also

[ArrayExpress](#), [queryAE](#), [getAE](#)

Examples

```
# An example can be found in the help of the getAE function.
```

ArrayExpress *R objects from ArrayExpress database*

Description

ArrayExpress produces an [AffyBatch](#), an [ExpressionSet](#) or a [NChannelSet](#) from a raw dataset from the ArrayExpress database. ArrayExpress needs an Internet connection.

Usage

```
ArrayExpress(accession, path = tempdir(), save = FALSE, dataCols = NULL, drop = TRUE)
```

Arguments

accession	an ArrayExpress experiment identifier.
path	the name of the directory in which the files downloaded on the ArrayExpress repository will be extracted. The default is the current directory.
save	if TRUE, the files downloaded from the database will not be deleted from path after executing the function.
dataCols	by default, for the raw data, the columns are automatically selected according to the scanner type. If the scanner is unknown or if the user wants to use different columns than the default, the argument 'dataCols' can be set. For two colour arrays it must be a list with the fields 'R', 'G', 'Rb' and 'Gb' giving the column names to be used for red and green foreground and background. For one colour arrays, it must be a character string with the column name to be used. These column names must correspond to existing column names of the expression files.
drop	if TRUE and only one platform in series, the platform name will be dropped.

Value

The output is an object of class [AffyBatch](#) or [ExpressionSet](#) or [NChannelSet](#) with the raw expression values in the assayData of the object, the information contained in the .sdrf file in the phenoData, the adf file in the featureData and the idf file content in the experimentData.

If several array designs are used in the data set, the output is a list with an object for each array design.

Author(s)

Audrey Kauffmann, Ibrahim Emam

Maintainer: <iemam@ebi.ac.uk>

See Also

[queryAE](#), [getAE](#), [ae2bioc](#), [getcolproc](#), [procset](#)

Examples

```
ETABM25.affybatch = ArrayExpress("E-TABM-25")
print(ETABM25.affybatch)
sampleNames(ETABM25.affybatch)
colnames(pData(ETABM25.affybatch))
```

extract.zip	<i>Unzip archives in a specified directory</i>
-------------	--

Description

extract.zip extracts the files from a .zip archive in a specific directory.

Usage

```
extract.zip(file, extractpath = dirname(file)[1])
```

Arguments

file	A file name.
extractpath	A path to define where the files are to be extracted.

Value

Success is indicated by returning the directory in which the files have been extracted. If it fails, it returns an empty character string.

Author(s)

Audrey Kauffmann
 Maintainer: <kauffmann@bergonie.org>

getAE	<i>Download MAGE-TAB files from ArrayExpress in a specified directory</i>
-------	---

Description

getAE downloads and extracts the MAGE-TAB files from an ArrayExpress dataset.

Usage

```
getAE(accession, path = getwd(), type = "full", extract = TRUE, local = FALSE, sourcedir = path)
```

Arguments

accession	is an ArrayExpress experiment identifier.
path	is the name of the directory in which the files downloaded on the ArrayExpress repository will be extracted.
type	can be 'raw' to download and extract only the raw data, 'processed' to download and extract only the processed data or 'full' to have both raw and processed data.
extract	if FALSE, the files are not extracted from the zip archive.
local	if TRUE, files will be read from a local folder specified by sourcedir.
sourcedir	when local = TRUE, files will be read from this directory.

Value

A list with the names of the files that have been downloaded and extracted.

Author(s)

Ibrahim Emam, Audrey Kauffmann

Maintainer: <iemam@ebi.ac.uk>

See Also

[ArrayExpress](#), [ae2bioc](#), [getcolproc](#), [procset](#)

Examples

```
mexp1422 = getAE("E-MEXP-1422", type = "full")

## Build a an ExpressionSet from the raw data
MEXP1422raw = ae2bioc(mageFiles = mexp1422)

## Build a an ExpressionSet from the processed data
cnames = getcolproc(mexp1422)
MEXP1422proc = procset(mexp1422, cnames[2])
```

getcolproc

Return the possible column names from processed MAGE-TAB files

Description

getcolproc extracts the column names from processed MAGE-TAB and return them. The output is needed to call the function procset.

Usage

```
getcolproc(files)
```

Arguments

files A list as given from [getAE](#) function. Containing the following elements:
profile profile is the name of the processed MAGE-TAB file to be read.
path is the name of the directory where to find this file.

Author(s)

Audrey Kauffmann

Maintainer: <iemam@ebi.ac.uk>

See Also

[ArrayExpress](#), [queryAE](#), [getAE](#), [procset](#)

`getcolraw`*Return the possible column names from raw MAGE-TAB files*

Description

`getcolraw` extracts the column names from raw MAGE-TAB and return them. The output can be use to set the argument 'rawcol' of the function `magetab2bioc`.

Usage

```
getcolraw(path, rawfiles)
```

Arguments

`rawfiles` rawfiles are the name of the raw MAGE-TAB files to be read.
`path` is the name of the directory where to find these files.

Author(s)

Audrey Kauffmann
Maintainer: <iemam@ebi.ac.uk>

See Also

[ArrayExpress](#), [queryAE](#), [getAE](#)

`procset`*Convert processed MAGE-TAB files into a Bioconductor object*

Description

`procset` converts local MAGE-TAB files into an [ExpressionSet](#).

Usage

```
procset(files, procol)
```

Arguments

`files` is the list with the names of the processed, the sdrf, the adf and the idf files and the path of the data as given by [getAE](#).
`procol` the name of the column to be extracted from the file. Obtained using [getcolproc](#).

Author(s)

Ibrahim Emam, Audrey Kauffmann
Maintainer: <iemam@ebi.ac.uk>

See Also

[queryAE](#), [getAE](#), [getcolproc](#)

Examples

```
# An example can be found in the help of the getAE function.
```

queryAE	<i>XML query of the ArrayExpress repository</i>
---------	---

Description

queryAE queries the ArrayExpress database with keywords and give a dataframe with ArrayExpress identifiers and related information, as an output.

Usage

```
queryAE(keywords = NULL, species = NULL)
```

Arguments

keywords	the keyword(s) of interest. To use several words, they must be separated by a "+" as shown in the examples.
species	the specie(s) of interest.

Value

A dataframe with all the ArrayExpress dataset identifiers which correspond to the query in the first column. The following columns contain information about these datasets, such as the number of raw files, the number of data processed, the release date on the database, the pubmed ID, the species, the experiment design and the experimental factors.

Author(s)

Ibrahim Emam, Audrey Kauffmann
Maintainer: <iemam@ebi.ac.uk>

See Also

[ArrayExpress](#), [getAE](#)

Examples

```
## To retrieve all the identifiers of pneumonia data sets  
pneumo = queryAE(keywords = "pneumonia")  
  
## To retrieve all the identifiers of pneumonia data sets studied in human  
pneumoHS = queryAE(keywords = "pneumonia", species = "homo+sapiens")
```

Index

* datasets

- ae2bioc, 2
- ArrayExpress, 3
- extract.zip, 4
- getAE, 4
- getcolproc, 5
- getcolraw, 6
- procset, 6
- queryAE, 7

- ae2bioc, 2, 3, 5
- AffyBatch, 2, 3
- ArrayExpress, 2, 3, 5–7

- ExpressionSet, 2, 3, 6
- extract.zip, 4

- getAE, 2, 3, 4, 5–7
- getcolproc, 3, 5, 5, 6, 7
- getcolraw, 6

- NChannelSet, 2, 3

- procset, 3, 5, 6

- queryAE, 2, 3, 5–7, 7