

# Package ‘taxonomizr’

August 31, 2018

**Maintainer** Scott Sherrill-Mix <shescott@upenn.edu>

**License** GPL (>= 2) | file LICENSE

**Title** Functions to Work with NCBI Accessions and Taxonomy

**Type** Package

**LazyLoad** yes

**Author** Scott Sherrill-Mix [aut, cre]

**BugReports** <https://github.com/sherrillmix/taxonomizr/issues>

## Description

Functions for assigning taxonomy to NCBI accession numbers and taxon IDs based on NCBI's accession2taxid and taxdump files. This package allows the user to download NCBI data dumps and create a local database for fast and local taxonomic assignment.

**Version** 0.5.1

**Date** 2018-08-31

**Suggests** testthat, knitr, rmarkdown

**Depends** R (>= 3.0.0)

**Imports** parallel, RSQLite, data.table, R.utils

**RoxygenNote** 6.0.1

**VignetteBuilder** knitr

**NeedsCompilation** yes

**Repository** CRAN

**Date/Publication** 2018-08-31 20:50:03 UTC

## R topics documented:

accessionToTaxa . . . . .	2
condenseTaxa . . . . .	3
getAccession2taxid . . . . .	4
getAccessions . . . . .	5
getId . . . . .	6
getId2 . . . . .	7

getNamesAndNodes . . . . .	8
getTaxonomy . . . . .	9
getTaxonomy2 . . . . .	10
lastNotNa . . . . .	12
prepareDatabase . . . . .	13
read.accession2taxid . . . . .	14
read.names . . . . .	15
read.names.sql . . . . .	16
read.nodes . . . . .	17
read.nodes.sql . . . . .	18
streamingRead . . . . .	19
taxonomizrSwitch . . . . .	20
trimTaxa . . . . .	20

<b>Index</b>	<b>22</b>
--------------	-----------

---

accessionToTaxa	<i>Convert accessions to taxa</i>
-----------------	-----------------------------------

---

## Description

Convert a vector of NCBI accession numbers to their assigned taxonomy

## Usage

```
accessionToTaxa(accessions, sqlFile, version = c("version", "base"))
```

## Arguments

accessions	a vector of NCBI accession strings to convert to taxa
sqlFile	a string giving the path to a SQLite file screated by <a href="#">read.accession2taxid</a>
version	either 'version' indicating that taxaids are versioned e.g. Z17427.1 or 'base' indicating that taxaids do not have version numbers e.g. Z17427

## Value

a vector of NCBI taxa ids

## References

<ftp://ftp.ncbi.nih.gov/pub/taxonomy/accession2taxid>

## See Also

[getTaxonomy](#), [read.accession2taxid](#)

**Examples**

```

taxa<-c(
  "accession\taccession.version\ttaxid\tgi",
  "Z17427\tZ17427.1\t3702\t16569",
  "Z17428\tZ17428.1\t3702\t16570",
  "Z17429\tZ17429.1\t3702\t16571",
  "Z17430\tZ17430.1\t3702\t16572",
  "X62402\tX62402.1\t9606\t30394"
)
inFile<-tempfile()
sqlFile<-tempfile()
writeLines(taxa,inFile)
read.accession2taxid(inFile,sqlFile,vocal=FALSE)
accessionToTaxa(c("Z17430.1","Z17429.1","X62402.1",'NOTREAL'),sqlFile)

```

---

condenseTaxa	<i>Condense multiple taxonomic assignments to their most recent common branch</i>
--------------	---

---

**Description**

Take a table of taxonomic assignments, e.g. assignments from hits to a read, and condense it to a single vector with NAs where there are disagreements between the hits.

**Usage**

```
condenseTaxa(taxaTable, groupings = rep(1, nrow(taxaTable)))
```

**Arguments**

taxaTable	a matrix or data.frame with hits on the rows and various levels of taxonomy in the columns
groupings	a vector of groups e.g. read queries to condense taxa within

**Value**

a matrix with `ncol(taxaTable)` taxonomy columns with a row for each unique id (labelled on `rownames`) with NAs where there was not complete agreement for an id

**Examples**

```

taxas<-matrix(c(
  'a','b','c','e',
  'a','b','d','e'
),nrow=2,byrow=TRUE)
condenseTaxa(taxas)
condenseTaxa(taxas[c(1,2,2),],c(1,1,2))

```

---

getAccession2taxid      *Download accession2taxid files from NCBI*

---

### Description

Download a nucl\_xxx.accession2taxid.gz from NCBI servers. These can then be used to create a SQLite datanase with [read.accession2taxid](#). Note that if the files already exist in the target directory then this function will not redownload them. Delete the files if a fresh download is desired.

### Usage

```
getAccession2taxid(outDir = ".",
  baseUrl = "ftp://ftp.ncbi.nih.gov/pub/taxonomy/accession2taxid/",
  types = c("nucl_gb", "nucl_est", "nucl_gss", "nucl_wgs"))
```

### Arguments

outDir	the directory to put the accession2taxid.gz files in
baseUrl	the url of the directory where accession2taxid.gz files are located
types	the types if accession2taxid.gz files desired where type is the prefix of xxx.accession2taxid.gz. The default is to download all nucl_ accessions. For protein accessions, try types=c('prot').

### Value

a vector of file path strings of the locations of the output files

### References

<ftp://ftp.ncbi.nih.gov/pub/taxonomy/>, <https://www.ncbi.nlm.nih.gov/Sequin/acc.html>

### See Also

[read.accession2taxid](#)

### Examples

```
## Not run:
  if(readline(
    "This will download a lot data and take a while to process.
    Make sure you have space and bandwidth. Type y to continue: "
  )!='y')
  stop('This is a stop to make sure no one downloads a bunch of data unintentionally')

  getAccession2taxid()

## End(Not run)
```

---

getAccessions	<i>Find all accessions for a taxa</i>
---------------	---------------------------------------

---

### Description

Find accessions numbers for a given taxa ID the NCBI taxonomy. This will be pretty slow unless the database was built with `indexTaxa=TRUE` since the database would not have an index for `taxaId`.

### Usage

```
getAccessions(taxaId, sqlFile, version = c("version", "base"), limit = NULL)
```

### Arguments

<code>taxaId</code>	a vector of taxonomic IDs
<code>sqlFile</code>	a string giving the path to a SQLite file created by <a href="#">read.accession2taxid</a>
<code>version</code>	either 'version' indicating that taxaids are versioned e.g. Z17427.1 or 'base' indicating that taxaids do not have version numbers e.g. Z17427
<code>limit</code>	return only this number of accessions or NULL for no limits

### Value

a vector of character strings giving taxa IDs (potentially comma concatenated for any taxa with ambiguous names)

### See Also

[read.accession2taxid](#)

### Examples

```
taxa<-c(
  "accession\taccession.version\ttaxid\tgi",
  "Z17427\tZ17427.1\t3702\t16569",
  "Z17428\tZ17428.1\t3702\t16570",
  "Z17429\tZ17429.1\t3702\t16571",
  "Z17430\tZ17430.1\t3702\t16572"
)
inFile<-tempfile()
outFile<-tempfile()
writeLines(taxa,inFile)
read.accession2taxid(inFile,outFile)
getAccessions(3702,outFile)
```





---

getNamesAndNodes	<i>Download names and nodes files from NCBI</i>
------------------	---

---

### Description

Download a taxdump.tar.gz file from NCBI servers and extract the names.dmp and nodes.dmp files from it. These can then be used to create a SQLite database with [read.names.sql](#) and [read.nodes.sql](#). Note that if the files already exist in the target directory then this function will not redownload them. Delete the files if a fresh download is desired.

### Usage

```
getNamesAndNodes(outDir = ".",  
  url = "ftp://ftp.ncbi.nih.gov/pub/taxonomy/taxdump.tar.gz",  
  fileNames = c("names.dmp", "nodes.dmp"))
```

### Arguments

outDir	the directory to put names.dmp and nodes.dmp in
url	the url where taxdump.tar.gz is located
fileNames	the filenames desired from the tar.gz file

### Value

a vector of file path strings of the locations of the output files

### References

<ftp://ftp.ncbi.nih.gov/pub/taxonomy/>, <https://www.ncbi.nlm.nih.gov/Taxonomy/taxonomyhome.html/>

### See Also

[read.nodes.sql](#), [read.names.sql](#)

### Examples

```
## Not run:  
  getNamesAndNodes()  
  
## End(Not run)
```











**Description**

A convenience function to return the last value which is not NA in a vector

**Usage**

```
lastNotNa(x, default = "Unknown")
```

**Arguments**

x                    a vector to look for the last value in  
 default            a default value to use when all values are NA in a vector

**Value**

a single element from the last non NA value in x (or the default)

**Examples**

```
lastNotNa(c(1:4,NA,NA))
lastNotNa(c(letters[1:4],NA,'z',NA))
lastNotNa(c(NA,NA))
```

---

prepareDatabase	<i>Download data from NCBI and set up SQLite database</i>
-----------------	---

---

**Description**

Convenience function to do all necessary preparations downloading names, nodes and accession2taxid data from NCBI and preprocessing into a SQLite database for downstream use.

**Usage**

```
prepareDatabase(sqlFile = "nameNode.sqlite", tmpDir = ".", vocal = TRUE,
  ...)
```

**Arguments**

sqlFile            character string giving the file location to store the SQLite database  
 tmpDir            location for storing the downloaded files from NCBI. (Note that it may be useful to store these somewhere convenient to avoid redownloading)  
 vocal            if TRUE output messages describing progress  
 ...            additional arguments to getNamesAndNodes, getAccession2taxid or read.accession2taxid

**Value**

a vector of character string giving the path to the SQLite file

**See Also**

[getNamesAndNodes](#), [getAccession2taxid](#), [read.accession2taxid](#), [read.nodes.sql](#), [read.names.sql](#)

**Examples**

```
## Not run:
  if(readline(
    "This will download a lot data and take a while to process.
    Make sure you have space and bandwidth. Type y to continue: "
  )!='y')
    stop('This is a stop to make sure no one downloads a bunch of data unintentionally')

  prepareDatabase()

## End(Not run)
```

---

read.accession2taxid *Read NCBI accession2taxid files*

---

**Description**

Take NCBI accession2taxid files, keep only accession and taxa and save it as a SQLite database

**Usage**

```
read.accession2taxid(taxaFiles, sqlFile, vocal = TRUE, extraSqlCommand = "",
  indexTaxa = FALSE, overwrite = FALSE)
```

**Arguments**

taxaFiles	a string or vector of strings giving the path(s) to files to be read in
sqlFile	a string giving the path where the output SQLite file should be saved
vocal	if TRUE output status messages
extraSqlCommand	for advanced use. A string giving a command to be called on the SQLite database before loading data e.g. "pragma temp_store = 2;" to keep all temp files in memory (don't do this unless you have a lot (>100 Gb) of RAM)
indexTaxa	if TRUE add an index for taxa ID. This would only be necessary if you want to look up accessions by taxa ID e.g. <a href="#">getAccessions</a>
overwrite	If TRUE, delete accessionTaxa table in database if present and regenerate

**Value**

TRUE if successful

**References**

<ftp://ftp.ncbi.nih.gov/pub/taxonomy/accession2taxid>

**See Also**

[read.nodes.sql](#), [read.names.sql](#)

**Examples**

```
taxa<-c(
  "accession\taccession.version\ttaxid\tgi",
  "Z17427\tZ17427.1\t3702\t16569",
  "Z17428\tZ17428.1\t3702\t16570",
  "Z17429\tZ17429.1\t3702\t16571",
  "Z17430\tZ17430.1\t3702\t16572"
)
inFile<-tempfile()
outFile<-tempfile()
writeLines(taxa,inFile)
read.accession2taxid(inFile,outFile,vocal=FALSE)
db<-RSQLite::dbConnect(RSQLite::SQLite(),dbname=outFile)
RSQLite::dbGetQuery(db,'SELECT * FROM accessionTaxa')
RSQLite::dbDisconnect(db)
```

---

read.names

*Read NCBI names file*

---

**Description**

Take an NCBI names file, keep only scientific names and convert it to a data.table. NOTE: This function is now deprecated for [read.names.sql](#) (using SQLite rather than data.table).

**Usage**

```
read.names(nameFile, onlyScientific = TRUE)
```

**Arguments**

**nameFile** string giving the path to an NCBI name file to read from (both gzipped or uncompressed files are ok)

**onlyScientific** If TRUE, only store scientific names. If FALSE, synonyms and other types are included (increasing the potential for ambiguous taxonomic assignments).

**Value**

a data.table with columns id and name with a key on id











---

taxonomizrSwitch	<i>Switch from data.table to SQLite</i>
------------------	---

---

### Description

In version 0.5.0, taxonomizr switched from data.table to SQLite name and node lookups. See below for more details.

### Details

Version 0.5.0 marked a change for name and node lookups from using data.table to using SQLite. This was necessary to increase performance (10-100x speedup for [getTaxonomy](#)) and create a simpler interface (a single SQLite database contains all necessary data). Unfortunately, this switch requires a couple breaking changes:

- [getTaxonomy](#) changes from `getTaxonomy(ids, namesDT, nodesDT)` to `getTaxonomy(ids, sqlFile)`
- [getId](#) changes from `getId(taxa, namesDT)` to `getId(taxa, sqlFile)`
- [read.names](#) is deprecated, instead use [read.names.sql](#). For example, instead of calling `names<-read.names('names.dmp')` in every session, simply call `read.names.sql('names.dmp', 'accessionTaxa.')` once (or use the convenient [prepareDatabase](#))).
- [read.nodes](#) is deprecated, instead use [read.names.sql](#). For example, instead of calling `nodes<-read.names('nodes.dmp')` in every session, simply call `read.nodes.sql('nodes.dmp', 'accessionTaxa.')` once (or use the convenient [prepareDatabase](#))).

I've tried to ease any problems with this by overloading [getTaxonomy](#) and [getId](#) to still function (with a warning) if passed a data.table names and nodes argument and providing a simpler [prepareDatabase](#) function for completing all setup steps (hopefully avoiding direct calls to [read.names](#) and [read.nodes](#) for most users).

I plan to eventually remove data.table functionality to avoid a split codebase so please switch to the new SQLite format in all new code.

### See Also

[getTaxonomy](#), [read.names.sql](#), [read.nodes.sql](#), [prepareDatabase](#), [getId](#)

---

trimTaxa	<i>Trim columns from taxa file</i>
----------	------------------------------------

---

### Description

A simple script to delete the first row and then delete the first and fourth column of a four column tab delimited file and write to another file.

**Usage**

```
trimTaxa(inFile, outFile, desiredCols = c(2, 3))
```

**Arguments**

<code>inFile</code>	a single string giving the 4 column tab separated file to read from
<code>outFile</code>	a single string giving the file path to write to
<code>desiredCols</code>	the integer IDs for columns to pull out from file

# Index

## \*Topic **internal**

taxonomizrSwitch, 20

accessionToTaxa, 2

condenseTaxa, 3

getAccession2taxid, 4, 14

getAccessions, 5, 14

getId, 6, 7, 20

getId2, 7

getNamesAndNodes, 8, 14

getTaxonomy, 2, 6, 9, 11, 20

getTaxonomy2, 10

lastNotNa, 12

prepareDatabase, 13, 20

read.accession2taxid, 2, 4, 5, 14, 14

read.names, 7, 11, 15, 17, 20

read.names.sql, 6, 8, 9, 14–16, 16, 18, 20

read.nodes, 11, 16, 17, 17, 20

read.nodes.sql, 8, 9, 14, 15, 17, 18, 20

streamingRead, 19

taxonomizrSwitch, 20

trimTaxa, 20