

Package ‘tm’

September 10, 2009

Title Text Mining Package

Version 0.5

Date 2009-09-08

Author Ingo Feinerer

Maintainer Ingo Feinerer <feinerer@logic.at>

Depends R (>= 2.7.0)

Imports slam (>= 0.1-1)

Suggests filehash, proxy, Rgraphviz, Rmpi, RWeka, snow, Snowball, XML

Description A framework for text mining applications within R.

License GPL (>= 2)

URL <http://tm.r-forge.r-project.org/>

Repository CRAN

Date/Publication 2009-09-10 19:52:53

R topics documented:

acq	3
as.PlainTextDocument	4
convert_UTF_8	4
crude	5
DataframeSource	6
Dictionary	7
DirSource	7
dissimilarity	8
findAssocs	9
findFreqTerms	10
FunctionGenerator	10

getFilters	11
getReaders	12
getSources	12
getTransformations	13
GmaneSource	13
inspect	14
makeChunks	14
materialize	15
meta	16
names	17
number	18
PCorpus	18
PlainTextDocument	19
plot	20
preprocessReut21578XML	21
prescindMeta	22
RCV1Document	23
readDOC	24
readGmane	25
readPDF	26
readPlain	27
readRCV1	28
readReut21578XML	29
readTabular	30
readXML	31
removeNumbers	33
removePunctuation	33
removeSparseTerms	34
removeWords	35
Reuters21578Document	35
ReutersSource	36
searchFullText	37
sFilter	38
Source	39
stemCompletion	40
stemDocument	41
stopwords	41
stripWhitespace	42
TermDocumentMatrix	43
termFreq	44
TextDocument	45
TextRepository	45
tm_cluster	46
tm_combine	47
tm_filter	47
tm_intersect	48
tm_map	49
tm_reduce	50

<i>acq</i>	3
URISource	51
VCorpus	52
VectorSource	53
weightBin	54
WeightFunction	54
weightTf	55
weightTfIdf	56
writeCorpus	56
XMLSource	57
Index	58

<i>acq</i>	<i>50 Exemplary News Articles from the Reuters-21578 XML Data Set of Topic acq</i>
------------	--

Description

This dataset holds 50 news articles with additional meta information from the Reuters-21578 XML data set. All documents belong to the topic *acq* dealing with corporate acquisitions.

Usage

```
data("acq")
```

Format

A corpus of 50 text documents.

Source

Reuters-21578 Text Categorization Collection Distribution 1.0 (XML format).

References

- Lewis, David (1997) *Reuters-21578 Text Categorization Collection Distribution 1.0*. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>
- Luz, Saturnino XML-encoded version of Reuters-21578. <http://modnlp.berlios.de/reuters21578.html>

Examples

```
data("acq")
summary(acq)
```

```
as.PlainTextDocument
```

Create Objects of Class PlainTextDocument

Description

Create objects of class PlainTextDocument.

Usage

```
## S3 method for class 'PlainTextDocument':  
as.PlainTextDocument(x)  
## S3 method for class 'Reuters21578Document':  
as.PlainTextDocument(x)  
## S3 method for class 'RCV1Document':  
as.PlainTextDocument(x)
```

Arguments

x A text document.

See Also

[getTransformations](#) to list available transformation (mapping) functions.

Examples

```
reut21578 <- system.file("texts", "crude", package = "tm")  
r <- Corpus(DirSource(reut21578), readerControl = list(reader = readReut21578XML))  
r[[1]]  
as.PlainTextDocument(r[[1]])
```

```
convert_UTF_8
```

Convert Encoding to UTF-8

Description

Convert the encoding of a character or text document to UTF-8.

Usage

```
convert_UTF_8(x, from = "", sub = NA, ...)
```

Arguments

<code>x</code>	A character or text document.
<code>from</code>	A character string describing the current encoding.
<code>sub</code>	A character string. If not 'NA' it is used to replace any non-convertible bytes in the input.
<code>...</code>	Additional arguments.

Details

This function internally uses [iconv](#) for conversion, and the arguments `from` and `sub` are passed over.

Value

The character or text document `x` converted to UTF-8 encoding.

Author(s)

Ingo Feinerer

See Also

[iconv](#)

Use [getTransformations](#) to list available transformation (mapping) functions.

Examples

```
data("crude")
convert_UTF_8(crude[[1]])
```

<code>crude</code>	<i>20 Exemplary News Articles from the Reuters-21578 XML Data Set of Topic crude</i>
--------------------	--

Description

This data set holds 20 news articles with additional meta information from the Reuters-21578 XML data set. All documents belong to the topic `crude` dealing with crude oil.

Usage

```
data("crude")
```

Format

A corpus of 20 text documents.

Source

Reuters-21578 Text Categorization Collection Distribution 1.0 (XML format).

References

Lewis, David (1997) *Reuters-21578 Text Categorization Collection Distribution 1.0*. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>

Luz, Saturnino XML-encoded version of Reuters-21578. <http://modnlp.berlios.de/reuters21578.html>

Examples

```
data("crude")
summary(crude)
```

DataframeSource	<i>Data Frame Source</i>
-----------------	--------------------------

Description

Constructs a source from a data frame.

Usage

```
DataframeSource(x, encoding = "UTF-8")
```

Arguments

x	A data frame holding the texts.
encoding	A character giving the encoding of x.

Value

An object of class `DataframeSource` which extends the class `Source` representing a data frame interpreting each row as a document.

Author(s)

Ingo Feinerer

Examples

```
docs <- data.frame(docs = c("This is a text.", "This another one."))
(ds <- DataframeSource(docs))
inspect(Corpus(ds))
```

*Dictionary**Dictionary*

Description

Constructs a dictionary from a character vector or a term-document matrix.

Usage

```
## S3 method for class 'character':  
Dictionary(x)  
## S3 method for class 'TermDocumentMatrix':  
Dictionary(x)
```

Arguments

x A character vector or a term-document matrix holding the terms for the dictionary.

Value

An object of class `Dictionary` which extends the class `character` representing a dictionary, i.e., a character vector of terms.

Author(s)

Ingo Feinerer

Examples

```
Dictionary(c("some", "tokens"))  
data(crude)  
Dictionary(TermDocumentMatrix(crude))
```

*DirSource**Directory Source*

Description

Constructs a directory source.

Usage

```
DirSource(directory, encoding = "UTF-8", pattern = NULL, recursive = FALSE, ignore.
```

Arguments

<code>directory</code>	A directory.
<code>encoding</code>	A character giving the encoding of the files in <code>directory</code> .
<code>pattern</code>	An optional regular expression. Only file names which match the regular expression will be returned.
<code>recursive</code>	Logical. Should the listing recurse into directories?
<code>ignore.case</code>	Logical. Should pattern-matching be case-insensitive?

Value

An object of class `DirSource` which extends the class `Source` representing a directory. Each file in this directory is considered to be a document.

Author(s)

Ingo Feinerer

<code>dissimilarity</code>	<i>Dissimilarity</i>
----------------------------	----------------------

Description

Compute the dissimilarity between documents in a term-document matrix or between two explicit documents.

Usage

```
## S3 method for class 'TermDocumentMatrix':
dissimilarity(x, y = NULL, method)
## S3 method for class 'PlainTextDocument':
dissimilarity(x, y = NULL, method)
```

Arguments

<code>x</code>	Either a term-document matrix or a text document.
<code>y</code>	A text document. Only used if <code>x</code> is a text document.
<code>method</code>	Dissimilarity measure. Any method accepted by <code>dist</code> from package proxy can be passed over.

Value

An object of class `dist` representing the dissimilarity between participating documents.

Examples

```
data("crude")
tdm <- TermDocumentMatrix(crude)
dissimilarity(tdm, method = "cosine")
dissimilarity(crude[[1]], crude[[2]], method = "eJaccard")
```

findAssocs

Find Associations in a Term-Document Matrix

Description

Find associations in a term-document matrix.

Usage

```
## S3 method for class 'TermDocumentMatrix':
findAssocs(x, term, corlimit)
## S3 method for class 'matrix':
findAssocs(x, term, corlimit)
```

Arguments

x	A term-document matrix.
term	A character holding a term.
corlimit	A numeric for the lower correlation bound limit.

Value

A named numeric with those terms from x which correlate with term more than corlimit.

Examples

```
data("crude")
tdm <- TermDocumentMatrix(crude)
findAssocs(tdm, "oil", 0.7)
```

findFreqTerms	<i>Find Frequent Terms</i>
---------------	----------------------------

Description

Find frequent terms in a term-document matrix.

Usage

```
findFreqTerms(x, lowfreq = 0, highfreq = Inf)
```

Arguments

x	A term-document matrix.
lowfreq	An integer for the lower frequency bound.
highfreq	An integer for the upper frequency bound.

Details

This method works for all numeric weightings but is probably most meaningful for the standard term frequency (tf) weighting of x.

Value

A character vector of terms in x which occur more or equal often than lowfreq times and less or equal often than highfreq times.

Examples

```
data("crude")
tdm <- TermDocumentMatrix(crude)
findFreqTerms(tdm, 2, 3)
```

FunctionGenerator	<i>Function Generator</i>
-------------------	---------------------------

Description

Construct a function generator.

Usage

```
FunctionGenerator(x)
```

Arguments

`x` A generator function which takes some input and constructs and returns a new function based on that input information.

Value

An object of class `FunctionGenerator` which extends the class `function` representing a function generator.

Author(s)

Ingo Feinerer

See Also

Most reader functions (use [getReaders](#) to list available readers) are function generators.

Examples

```
funGen <- FunctionGenerator(function(y, ...) {  
  if (is(y, "integer")) function(x) x+1 else function(x) x-1  
})  
funGen  
funGen(3L)  
funGen("a")
```

`getFilters`*List Available Filters*

Description

List available filters which can be used with [tm_filter](#) and [tm_index](#).

Usage

```
getFilters()
```

Value

A character vector with available filters.

Author(s)

Ingo Feinerer

Examples

```
getFilters()
```

`getReaders`*List Available Readers*

Description

List available readers.

Usage

```
getReaders()
```

Value

A character vector with available readers.

Author(s)

Ingo Feinerer

Examples

```
getReaders()
```

`getSources`*List Available Sources*

Description

Listt available sources.

Usage

```
getSources()
```

Value

A character vector with available sources.

Author(s)

Ingo Feinerer

Examples

```
getSources()
```

getTransformations *List Available Transformations*

Description

List available transformations (mappings) which can be used with [tm_map](#).

Usage

```
getTransformations()
```

Value

A character vector with available transformations.

Author(s)

Ingo Feinerer

Examples

```
getTransformations()
```

GmaneSource *Gmane Source*

Description

Construct a source for a Gmane mailing list RSS feed.

Usage

```
GmaneSource(x, encoding = "UTF-8")
```

Arguments

x	Either a character identifying the file or a connection.
encoding	A character giving the encoding of x.

Value

An object of class `GmaneSource` which extends the class `Source` representing a Gmane mailing list RSS feed.

Author(s)

Ingo Feinerer

inspect

Inspect Objects

Description

Inspect, i.e., display detailed information on a corpus or a term-document matrix.

Usage

```
## S3 method for class 'PCorpus':
inspect(x)
## S3 method for class 'VCorpus':
inspect(x)
## S3 method for class 'TermDocumentMatrix':
inspect(x)
```

Arguments

x Either a corpus or a term-document matrix.

Examples

```
data("crude")
inspect(crude[1:3])
tdm <- TermDocumentMatrix(crude)[1:10, 1:10]
inspect(tdm)
```

makeChunks

Split a Corpus into Chunks

Description

Split a corpus into equally sized chunks conserving document boundaries.

Usage

```
makeChunks(corpus, chunksize)
```

Arguments

corpus The corpus to be split into chunks.
chunksize The chunk size.

Value

A corpus consisting of the chunks. Note that corpus meta data is not passed on to the newly created chunk corpus.

Author(s)

Ingo Feinerer

Examples

```
txt <- system.file("texts", "txt", package = "tm")
ovid <- Corpus(DirSource(txt))
sapply(ovid, length)
ovidChunks <- makeChunks(ovid, 5)
sapply(ovidChunks, length)
```

materialize*Materialize Lazy Mappings*

Description

The function `tm_map` supports so-called lazy mappings, that are mappings which are delayed until the documents' content is accessed. This function triggers the evaluation, i.e., it materializes the documents.

Usage

```
materialize(corpus, range = seq_along(corpus))
```

Arguments

<code>corpus</code>	A document collection with lazy mappings.
<code>range</code>	The indices of documents to be materialized.

Value

A corpus with materialized, i.e., all mappings computed and applied, documents for the requested range.

Author(s)

Ingo Feinerer

See Also[tm_map](#)**Examples**

```
data("crude")
x <- tm_map(crude, stemDocument, lazy = TRUE)
x <- materialize(x)
```

Description

Methods to access and modify meta data of documents, corpora, and repositories. In addition to **tm**'s internal meta data structures, Simple Dublin Core meta data mappings are available.

Usage

```
## S3 method for class 'Corpus':
meta(x, tag, type = c("indexed", "corpus", "local"))
## S3 method for class 'TextDocument':
meta(x, tag, type = NULL)
## S3 method for class 'TextRepository':
meta(x, tag, type = NULL)
DublinCore(x, tag = NULL)
```

Arguments

<code>x</code>	Either a text document, a corpus, or a text repository.
<code>tag</code>	A character identifying the name of the meta datum.
<code>type</code>	A character specifying which meta data of a corpus should be considered.

Details

In general this function can be used to display meta information but also to modify individual meta data:

x = "TextDocument", tag = NULL If no `tag` is given, this method pretty prints all `x`'s meta data. If `tag` is provided its value in the meta data is returned.

x = "Corpus", tag = NULL, type = "indexed" This method investigates the `type` argument. `type` must be either `indexed` (default), `local`, or `corpus`. Former is a shortcut for accessing document level meta data ([DMetaData](#)) stored at the collection level (because it forms an own entity, or for performance reasons, i.e., a form of indexing, hence the name `indexed`), `local` accesses the meta data local to each text document (i.e., meta data in text documents' attributes), and `corpus` is a shortcut for corpus specific meta data ([CMetaData](#)). Depending whether a `tag` is set or not, all or only the meta data identified by the `tag` is displayed or modified.

x = "TextRepository", tag = NULL If no `tag` is given, this method pretty prints all `x`'s meta data. If `tag` is provided its value in the meta data is returned.

Simple Dublin Core meta data is only available locally at each document:

x = "TextDocument", tag = NULL Returns or sets the Simple Dublin Core meta datum named `tag` for `x`. `tag` must be a valid Simple Dublin Core element name (i.e. Title, Creator, Subject, Description, Publisher, Contributor, Date, Type, Format, Identifier, Source, Language, Relation, Coverage, or Rights) or `NULL`. For the latter all Dublin Core meta data are printed.

References

Dublin Core Metadata Initiative. <http://dublincore.org/>

Examples

```
data("crude")
meta(crude[[1]])
DublinCore(crude[[1]])
meta(crude[[1]], tag = "Topics")
meta(crude[[1]], tag = "Comment") <- "A short comment."
meta(crude[[1]], tag = "Topics") <- NULL
DublinCore(crude[[1]], tag = "Creator") <- "Ano Nymous"
DublinCore(crude[[1]], tag = "Format") <- "XML"
DublinCore(crude[[1]])
meta(crude[[1]])
meta(crude)
meta(crude, type = "corpus")
meta(crude, "labels") <- 21:40
meta(crude)
```

names

Row, Column, Dim Names, Document IDs, and Terms

Description

Retrieve the row names, column names, dimnames, document IDs, and terms of a term-document matrix or document-term matrix.

Arguments

x Either a term-document matrix or document-term matrix.

Examples

```
data("crude")
tdm <- TermDocumentMatrix(crude)[1:10, 1:20]
rownames(tdm)
colnames(tdm)
dimnames(tdm)
Docs(tdm)
Terms(tdm)
```

number	<i>The Number of Rows/Columns/Dimensions/Documents/Terms of a Term-Document Matrix</i>
--------	--

Description

Return the number of rows, columns, dimensions, documents, and terms of a term-document matrix or a document-term matrix.

Arguments

`x` Either a term-document matrix or a document-term matrix.

Examples

```
data("crude")
tdm <- TermDocumentMatrix(crude)[1:10,1:20]
ncol(tdm)
nrow(tdm)
dim(tdm)
nDocs(tdm)
nTerms(tdm)
```

PCorpus	<i>Permanent Corpus Constructor</i>
---------	-------------------------------------

Description

Construct a permanent corpus.

Usage

```
PCorpus(x, readerControl = list(reader = x$DefaultReader, language = "eng"), dbControl(x)
DBControl(x)
## S3 method for class 'PCorpus':
DMetaData(x)
```

Arguments

`x` A [Source](#) object for PCorpus, and a corpus for the other functions.

`readerControl` A list with the named components `reader` representing a reading function capable of handling the file format found in `x`, and `language` giving the text's language (preferably in ISO 639-2 format).

dbControl	A list with the named components <code>dbName</code> giving the filename holding the sourced out documents (i.e., the database), and <code>dbType</code> holding a valid database type as supported by package filehash . Under activated database support the <code>tm</code> package tries to keep as few as possible resources in memory under usage of the database.
...	Optional arguments for the reader.

Details

Permanent means that documents are physically stored outside of R (e.g., in a database) and R objects are only pointers to external structures. I.e., changes in the underlying external representation can affect multiple R objects simultaneously.

The constructed corpus object inherits from a `list` and has three attributes containing meta and database management information:

CMetaData Corpus Meta Data contains corpus specific meta data in form of tag-value pairs and information about children in form of a binary tree. This information is useful for reconstructing meta data after e.g. merging corpora.

DMetaData Document Meta Data of class `data.frame` contains document specific meta data for the corpus. This data frame typically encompasses clustering or classification results which basically are metadata for documents but form an own entity (e.g., with its name, the value range, etc.).

DBControl Database control field is a `list` with two named components: `dbName` holds the path to the permanent database storage, and `dbType` stores the database type.

Value

An object of class `PCorpus` which extends the classes `Corpus` and `list` containing a permanent corpus.

Author(s)

Ingo Feinerer

Examples

```
txt <- system.file("texts", "txt", package = "tm")
## Not run: PCorpus(DirSource(txt), dbControl = list(dbName = "myDB.db", dbType = "DB1"))
```

PlainTextDocument *Plain Text Document*

Description

Construct an object representing a plain text document with additional meta information.

Arguments

<code>x</code>	A term-document matrix.
<code>terms</code>	Terms to be plotted. Defaults to 20 randomly chosen terms of the term-document matrix.
<code>corThreshold</code>	Do not plot correlations below this threshold. Defaults to 0.7.
<code>weighting</code>	Define whether the line width corresponds to the correlation.
<code>attrs</code>	Argument passed to the plot method for class <code>graphNEL</code> .
<code>...</code>	Other arguments passed to the <code>graphNEL</code> plot method.

Details

Visualization requires that package **Rgraphviz** is available.

Examples

```
data(crude)
tdm <- TermDocumentMatrix(crude)
set.seed(1234)
plot(tdm, corThreshold = 0.2, weighting = TRUE)
```

```
preprocessReut21578XML
```

Preprocess the Reuters-21578 XML archive.

Description

Preprocess the Reuters-21578 XML archive by correcting invalid UTF-8 encodings and copying each text document into a separate file.

Usage

```
preprocessReut21578XML(input, output, fixEnc = TRUE)
```

Arguments

<code>input</code>	A character describing the input directory.
<code>output</code>	A character describing the output directory.
<code>fixEnc</code>	A logical value indicating whether an invalid UTF-8 encoding in the Reuters-21578 XML dataset should be corrected.

Value

No explicit return value. As a side product the directory `output` contains the corrected dataset.

Author(s)

Ingo Feinerer

References

Lewis, David (1997) *Reuters-21578 Text Categorization Collection Distribution 1.0*. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>

Luz, Saturnino XML-encoded version of Reuters-21578. <http://modnlp.berlios.de/reuters21578.html>

prescindMeta

Prescind Document Meta Data

Description

Extracts meta data from each individual document (either stored in its attributes or in additional user-defined local meta data pairs) of a corpus and creates a data frame which contains both the global meta data information of the corpus plus the extracted (i.e., shifted up) local meta data of the individual text documents.

Usage

```
prescindMeta(x, meta)
```

Arguments

<code>x</code>	A corpus.
<code>meta</code>	A character vector of meta data names to be shifted up.

Value

A data frame constructed from `x` with shifted up `meta` data.

See Also

[DMetaData](#), and [meta](#)

Examples

```
data("crude")
DMetaData(crude)
meta(crude, tag = "ID", type = "local")
prescindMeta(crude, c("ID", "Heading"))
```

RCV1Document	<i>RCV1 Text Document</i>
--------------	---------------------------

Description

Construct an object representing a RCV1 XML text document with meta information.

Usage

```
RCV1Document(x, author = character(), datetimestamp = as.POSIXlt(Sys.time()), tz =
```

Arguments

<code>x</code>	Object of class <code>list</code> containing the content.
<code>author</code>	Object of class <code>character</code> containing the author names.
<code>datetimestamp</code>	Object of class <code>POSIXlt</code> containing the date and time when the document was written.
<code>description</code>	Object of class <code>character</code> containing additional text information.
<code>heading</code>	Object of class <code>character</code> containing the title or a short heading.
<code>id</code>	Object of class <code>character</code> containing an identifier.
<code>origin</code>	Object of class <code>character</code> containing information on the source and origin of the text.
<code>language</code>	Object of class <code>character</code> containing the language of the text (preferably in ISO 639-2 format).
<code>localmetadata</code>	Object of class <code>list</code> containing local meta data in form of tag-value pairs.

Author(s)

Ingo Feinerer

References

Lewis, D. D.; Yang, Y.; Rose, T.; and Li, F (2004). RCV1: A New Benchmark Collection for Text Categorization Research. *Journal of Machine Learning Research*, **5**, 361–397. <http://www.jmlr.org/papers/volume5/lewis04a/lewis04a.pdf>

See Also

[PlainTextDocument](#) and [Reuters21578Document](#)

readDOC

*Read In a MS Word Document***Description**

Return a function which reads in a Microsoft Word document extracting its text.

Usage

```
readDOC(AntiwordOptions = "", ...)
```

Arguments

AntiwordOptions	Options passed over to antiword.
...	Arguments for the generator function.

Details

Formally this function is a function generator, i.e., it returns a function (which reads in a text document) with a well-defined signature, but can access passed over arguments (e.g., options to antiword) via lexical scoping.

Note that this MS Word reader needs the tool `antiword` installed and accessible on your system.

Value

A function with the signature `elem, language, id`:

<code>elem</code>	A list with the two named elements <code>content</code> and <code>uri</code> . The first element must hold the document to be read in, the second element must hold a call to extract this document. The call is evaluated upon a request for load on demand.
<code>language</code>	A character vector giving the text's language.
<code>id</code>	A character vector representing a unique identification string for the returned text document.

The function returns a `PlainTextDocument` representing the text in `content`.

Author(s)

Ingo Feinerer

See Also

[getReaders](#) to list available reader functions.

readGmane	<i>Read In a Gmane RSS Feed</i>
-----------	---------------------------------

Description

Read in a RSS feed as returned by Gmane for mailing lists.

Usage

```
readGmane(elem, language, id)
```

Arguments

elem	A list with the two named elements <code>content</code> and <code>uri</code> . The first element holds the document to be read in, the second element holds a call to extract this document. The call is evaluated upon a request for load on demand.
language	A character vector giving the text's language.
id	A character vector representing a unique identification string for the returned text document.

Value

A `PlainTextDocument`.

Author(s)

Ingo Feinerer

See Also

[getReaders](#) to list available reader functions.

Examples

```
gs <- GmaneSource(url("http://rss.gmane.org/gmane.comp.lang.r.general"))
elem <- getElem(stepNext(gs))
(gmane <- readGmane(elem, language = "eng", id = "id1"))
meta(gmane)
```

readPDF

*Read In a PDF Document***Description**

Return a function which reads in a portable document format (PDF) document extracting both its text and its meta data.

Usage

```
readPDF(PdftotextOptions = "", PdftotextOptions = "", ...)
```

Arguments

PdftotextOptions	Options passed over to pdftotext.
PdftotextOptions	Options passed over to pdftotext.
...	Arguments for the generator function.

Details

Formally this function is a function generator, i.e., it returns a function (which reads in a text document) with a well-defined signature, but can access passed over arguments (e.g., options to `pdftotext` or `pdftotext`) via lexical scoping.

Note that this PDF reader needs both the tools `pdftotext` and `pdftotext` installed and accessible on your system.

Value

A function with the signature `elem, language, id`:

<code>elem</code>	A list with the two named elements <code>content</code> and <code>uri</code> . The first element must hold the document to be read in, the second element must hold a call to extract this document. The call is evaluated upon a request for load on demand.
<code>language</code>	A character vector giving the text's language.
<code>id</code>	A character vector representing a unique identification string for the returned text document.

The function returns a `PlainTextDocument` representing the text and meta data in `content`.

Author(s)

Ingo Feinerer

See Also

[getReaders](#) to list available reader functions.

Examples

```
f <- system.file("doc", "tm.pdf", package = "tm")
pdf <- readPDF(PdftotextOptions = "-layout") (elem = list(uri = f), language = "eng", id = "i
pdf[1:13]
```

readPlain

*Read In a Text Document***Description**

Return a function which reads in a text document without knowledge about its internal structure and possible available metadata.

Usage

```
readPlain(...)
```

Arguments

... Arguments for the generator function.

Details

Formally this function is a function generator, i.e., it returns a function (which reads in a text document) with a well-defined signature, but can access passed over arguments via lexical scoping. This is especially useful for reader functions for complex data structures which need a lot of configuration options.

Value

A function with the signature `elem, language, id`:

<code>elem</code>	A list with the two named elements <code>content</code> and <code>uri</code> . The first element must hold the document to be read in, the second element must hold a call to extract this document. The call is evaluated upon a request for load on demand.
<code>language</code>	A character vector giving the text's language.
<code>id</code>	A character vector representing a unique identification string for the returned text document.

The function returns a `PlainTextDocument` representing content.

Author(s)

Ingo Feinerer

See Also

[getReaders](#) to list available reader functions.

Examples

```
docs <- c("This is a text.", "This another one.")
vs <- VectorSource(docs)
elem <- getElem(stepNext(vs))
(result <- readPlain()(elem, "en", "id1"))
meta(result)
```

readRCV1

Read In a Reuters Corpus Volume 1 Document

Description

Read in a Reuters Corpus Volume 1 XML document.

Usage

```
readRCV1(elem, language, id)
```

Arguments

<code>elem</code>	A list with the two named elements <code>content</code> and <code>uri</code> . The first element holds the document to be read in, the second element holds a call to extract this document. The call is evaluated upon a request for load on demand.
<code>language</code>	A character vector giving the text's language.
<code>id</code>	A character vector representing a unique identification string for the returned text document.

Value

An RCV1Document.

Author(s)

Ingo Feinerer

References

Lewis, D. D.; Yang, Y.; Rose, T.; and Li, F (2004). RCV1: A New Benchmark Collection for Text Categorization Research. *Journal of Machine Learning Research*, **5**, 361–397. <http://www.jmlr.org/papers/volume5/lewis04a/lewis04a.pdf>

See Also

Use [getReaders](#) to list available reader functions.

Examples

```
f <- system.file("texts", "rcv1_2330.xml", package = "tm")
rcv1 <- readRCV1(elem = list(content = readLines(f)), language = "eng", id = "id1")
meta(rcv1)
```

readReut21578XML	<i>Read In a Reuters-21578 XML Document</i>
------------------	---

Description

Read in a Reuters-21578 XML document.

Usage

```
readReut21578XML(elem, language, id)
readReut21578XMLasPlain(elem, language, id)
```

Arguments

<code>elem</code>	A list with the two named elements <code>content</code> and <code>uri</code> . The first element holds the document to be read in, the second element holds a call to extract this document. The call is evaluated upon a request for load on demand.
<code>language</code>	A character vector giving the text's language.
<code>id</code>	A character vector representing a unique identification string for the returned text document.

Value

A `Reuters21578Document` for `readReut21578XML`, or a `PlainTextDocument` for `readReut21578XMLasPlain`.

Author(s)

Ingo Feinerer

References

Lewis, David (1997) *Reuters-21578 Text Categorization Collection Distribution 1.0*. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>

Luz, Saturnino XML-encoded version of Reuters-21578. <http://modnlp.berlios.de/reuters21578.html>

See Also

[getReaders](#) to list available reader functions.

readTabular

Read In a Text Document

Description

Return a function which reads in a text document from a tabular data structure (like a data frame or a list matrix) with knowledge about its internal structure and possible available metadata as specified by a so-called mapping.

Usage

```
readTabular(mapping, ...)
```

Arguments

mapping	A named list of characters. The constructed reader will map each character entry to the content or meta datum of the text document as specified by the named list entry. Valid names include <code>Content</code> to access the document's content, any valid attribute name, and characters which are mapped to <code>LocalMetaData</code> entries.
...	Arguments for the generator function.

Details

Formally this function is a function generator, i.e., it returns a function (which reads in a text document) with a well-defined signature, but can access passed over arguments (e.g., the mapping) via lexical scoping.

Value

A function with the signature `elem, language, id`:

elem	A list with the two named elements <code>content</code> and <code>uri</code> . The first element must hold the document to be read in, the second element must hold a call to extract this document. The call is evaluated upon a request for load on demand.
language	A character vector giving the text's language.
id	A character vector representing a unique identification string for the returned text document.

The function returns a `PlainTextDocument` representing content.

Author(s)

Ingo Feinerer

See Also

Vignette ‘Extensions: How to Handle Custom File Formats’.

[getReaders](#) to list available reader functions.

Examples

```
df <- data.frame(contents = c("content 1", "content 2", "content 3"),
                 title    = c("title 1" , "title 2" , "title 3" ),
                 authors  = c("author 1" , "author 2" , "author 3" ),
                 topics   = c("topic 1" , "topic 2" , "topic 3" ),
                 stringsAsFactors = FALSE)
m <- list(Content = "contents", Heading = "title",
          Author = "authors", Topic = "topics")
myReader <- readTabular(mapping = m)
ds <- DataframeSource(df)
elem <- getElem(stepNext(ds))
(result <- myReader(elem, language = "en", id = "idl"))
meta(result)
```

readXML

Read In an XML Document

Description

Return a function which reads in an XML document. The structure of the XML document can be described with a so-called specification.

Usage

```
readXML(spec, doc, ...)
```

Arguments

spec A named list of lists each containing two character vectors. The constructed reader will map each list entry to a attribute or meta datum corresponding to the named list entry. Valid names include `Content` to access the document's content, any valid attribute name, and characters which are mapped to [LocalMetaData](#) entries.

Each list entry must consist of two character vectors: the first describes the type of the second argument, and the second is the specification entry. Valid combinations are:

type = "node", spec = "XPathExpression" The XPath expression `spec` extracts information from an XML node.

type = "attribute", spec = "XPathExpression" The XPath expression `spec` extracts information from an attribute of an XML node.

type = "function", spec = function(tree) ... The function `spec` is called, passing over a tree representation (as delivered by `xmlInternalTreeParse` from package **XML**) of the read in XML document as first argument.

type = "unevaluated", spec = "String" The character vector `spec` is returned without modification.

`doc` An (empty) document of some subclass of `TextDocument`

`...` Arguments for the generator function.

Details

Formally this function is a function generator, i.e., it returns a function (which reads in a text document) with a well-defined signature, but can access passed over arguments (e.g., the specification) via lexical scoping.

Value

A function with the signature `elem, language, id`:

`elem` A list with the two named elements `content` and `uri`. The first element must hold the document to be read in, the second element must hold a call to extract this document. The call is evaluated upon a request for load on demand.

`language` A character vector giving the text's language.

`id` A character vector representing a unique identification string for the returned text document.

The function returns `doc` augmented by the parsed information out of the XML file as described by `spec`.

Author(s)

Ingo Feinerer

See Also

Vignette 'Extensions: How to Handle Custom File Formats'.

[getReaders](#) to list available reader functions.

Examples

```
## Not run:
readReut21578XML <- readXML(
  spec = list(Author = list("node", "/REUTERS/TEXT/AUTHOR"),
    DateTimeStamp = list("function", function(node)
      strptime(sapply(XML::getNodeSet(node, "/REUTERS/DATE"), XML::xmlValue),
        format = "
          tz = "GMT")),
    Description = list("unevaluated", ""),
    Heading = list("node", "/REUTERS/TEXT/TITLE"),
    ID = list("attribute", "/REUTERS/@NEWID"),
```



```

Origin = list("unevaluated", "Reuters-21578 XML"),
Topics = list("node", "/REUTERS/TOPICS/D"),
doc = Reuters21578Document()
## End(Not run)

```

removeNumbers	<i>Remove Numbers from a Text Document</i>
---------------	--

Description

Strip any numbers from a text document.

Usage

```

## S3 method for class 'PlainTextDocument':
removeNumbers(x)

```

Arguments

`x` A text document.

Value

The text document with any numbers in it removed.

See Also

[getTransformations](#) to list available transformation (mapping) functions.

Examples

```

data("crude")
crude[[1]]
removeNumbers(crude[[1]])

```

removePunctuation	<i>Remove Punctuation Marks from a Text Document</i>
-------------------	--

Description

Remove all punctuation marks from a text document.

Usage

```

## S3 method for class 'PlainTextDocument':
removePunctuation(x)

```

Arguments

`x` A text document.

Value

The text document with any punctuation marks in it removed.

See Also

[getTransformations](#) to list available transformation (mapping) functions.

Examples

```
data("crude")
crude[[1]]
removePunctuation(crude[[1]])
```

`removeSparseTerms` *Remove Sparse Terms from a Term-Document Matrix*

Description

Remove sparse terms from a term-document matrix.

Usage

```
removeSparseTerms(x, sparse)
```

Arguments

`x` A term-document matrix.
`sparse` A numeric for the maximal allowed sparsity.

Value

A term-document matrix where those terms from `x` are removed which have at least a `sparse` percentage of empty (i.e., terms occurring 0 times in a document) elements. I.e., the resulting matrix contains only terms with a sparse factor of less than `sparse`.

Examples

```
data("crude")
tdm <- TermDocumentMatrix(crude)
removeSparseTerms(tdm, 0.2)
```

`removeWords`*Remove Words from a Text Document*

Description

Remove a set of words from a text document.

Usage

```
## S3 method for class 'PlainTextDocument':  
removeWords(x, words)
```

Arguments

<code>x</code>	A text document.
<code>words</code>	A character vector list the words to be removed.

Value

The text document with the specified words in it removed.

See Also

[getTransformations](#) to list available transformation (mapping) functions.

Examples

```
data("crude")  
crude[[1]]  
removeWords(crude[[1]], stopwords("english"))
```

`Reuters21578Document`*Reuters-21578 Text Document*

Description

Construct an object representing a Reuters-21578 XML text document with meta information.

Usage

```
Reuters21578Document(x, author = character(0), datetimestamp = as.POSIXlt(Sys.time
```

Arguments

<code>x</code>	Object of class <code>list</code> containing the content.
<code>author</code>	Object of class <code>character</code> containing the author names.
<code>datetimestamp</code>	Object of class <code>POSIXlt</code> containing the date and time when the document was written.
<code>description</code>	Object of class <code>character</code> containing additional text information.
<code>heading</code>	Object of class <code>character</code> containing the title or a short heading.
<code>id</code>	Object of class <code>character</code> containing an identifier.
<code>origin</code>	Object of class <code>character</code> containing information on the source and origin of the text.
<code>language</code>	Object of class <code>character</code> containing the language of the text (preferably in ISO 639-2 format).
<code>localmetadata</code>	Object of class <code>list</code> containing local meta data in form of tag-value pairs.

Author(s)

Ingo Feinerer

References

Lewis, David (1997) *Reuters-21578 Text Categorization Collection Distribution 1.0*. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>

Luz, Saturnino XML-encoded version of Reuters-21578. <http://modnlp.berlios.de/reuters21578.html>

See Also

[PlainTextDocument](#) and [RCV1Document](#)

ReutersSource

Reuters-21578 XML Source

Description

Construct a source for an input containing several Reuters-21578 XML documents.

Usage

```
ReutersSource(x, encoding = "UTF-8")
```

Arguments

<code>x</code>	Either a character identifying the file or a connection.
<code>encoding</code>	A character giving the encoding of <code>x</code> .

Value

An object of class `XMLSource` which extends the class `Source` representing a Reuters-21578 XML document.

Author(s)

Ingo Feinerer

References

Lewis, David (1997) *Reuters-21578 Text Categorization Collection Distribution 1.0*. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>

Luz, Saturnino XML-encoded version of Reuters-21578. <http://modnlp.berlios.de/reuters21578.html>

Examples

```
reuters21578 <- system.file("texts", "reuters-21578.xml", package = "tm")
rs <- ReutersSource(reuters21578)
inspect(Corpus(rs)[1:2])
```

searchFullText	<i>Full Text Search</i>
----------------	-------------------------

Description

Perform a full text search in text documents.

Usage

```
## S3 method for class 'PlainTextDocument':
searchFullText(x, pattern)
```

Arguments

<code>x</code>	A text document.
<code>pattern</code>	A regular expression as accepted by gsub .

Value

TRUE if the regular expression `pattern` matches in `x`, otherwise FALSE.

See Also

[getFilters](#) to list available filter functions.

Examples

```
data("crude")
searchFullText(crude[[1]], "co[m]?pany")
```

sFilter

*Statement Filter***Description**

Filter meta data by user-defined statements.

Usage

```
sFilter(x, s)
```

Arguments

x A Corpus.
s A statement of format "tag1 == 'expr1' & tag2 == 'expr2' & ...".

Details

The statement *s* models a simple query language. It consists of an expression as passed over to a data frame for subsetting. Tags in *s* represent meta data variables. Variables only available at document level are shifted up to the data frame if necessary. Note that the meta data tags for the slots Author, DateTimeStamp, Description, ID, Origin and Heading are author, datetimestamp, description, id, origin and heading, respectively, to avoid name conflicts.

Value

A logical vector to represent the subset of the DMetaData (extended for shifted up variables) data frame as specified by the statement.

Author(s)

Ingo Feinerer

See Also

[getFilters](#) to list available filter functions.

Examples

```
data("crude")
sFilter(crude, "id == '127' & heading == 'DIAMOND SHAMROCK (DIA) CUTS CRUDE PRICES'")
```

Source

*Access Sources***Description**

Methods to access sources which abstract input locations, like a directory, a connection, or simply an R vector.

Usage

```
## S3 method for class 'DataframeSource':
eoi(x)
## S3 method for class 'DirSource':
eoi(x)
## S3 method for class 'URISource':
eoi(x)
## S3 method for class 'VectorSource':
eoi(x)
## S3 method for class 'XMLSource':
eoi(x)
## S3 method for class 'DataframeSource':
getElem(x)
## S3 method for class 'DirSource':
getElem(x)
## S3 method for class 'URISource':
getElem(x)
## S3 method for class 'VectorSource':
getElem(x)
## S3 method for class 'XMLSource':
getElem(x)
## S3 method for class 'DataframeSource':
pGetElem(x)
## S3 method for class 'DirSource':
pGetElem(x)
## S3 method for class 'VectorSource':
pGetElem(x)
## S3 method for class 'Source':
stepNext(x)
```

Arguments

x A source.

Details

The class `Source` is implemented as a list with following components:

DefaultReader Object of class `function` holding a default reader.

Encoding Object of class `character` holding the encoding of the texts delivered by the source.

Length Object of class `numeric` denoting the number of the elements delivered by the source.
If the number cannot be determined in advance it is set to zero.

LoDSupport Object of class `logical` indicating whether this source supports load on demand.

Position object of class `numeric` indicating the position in the source.

Vectorized object of class `logical` indicating the ability for parallel element access.

The function `eof` returns `TRUE` if the end of input of the source is reached. `getElement` fetches the element at the current position, whereas `pGetElement` retrieves all elements in parallel at once. `stepNext` increases the position in the source to the next element.

Author(s)

Ingo Feinerer

See Also

[getSource](#) to list available sources.

stemCompletion	<i>Complete Stems</i>
----------------	-----------------------

Description

Heuristically complete stemmed words.

Usage

```
stemCompletion(x, words, type = c("prevalent", "first"))
```

Arguments

<code>x</code>	A Corpus to be searched for possible completions.
<code>words</code>	A character vector of stems to be completed.
<code>type</code>	A character naming the heuristics to be used: <code>prevalent</code> is default and takes the most frequent match as completion, whereas <code>first</code> takes the first found completion.

Value

A character vector with completed words.

Author(s)

Ingo Feinerer

Examples

```
data("crude")
stemCompletion(crude, c("compan", "entit", "suppl"))
```

stemDocument	<i>Stem Words</i>
--------------	-------------------

Description

Stem words in a text document using Porter's stemming algorithm.

Usage

```
## S3 method for class 'PlainTextDocument':
stemDocument(x, language = "english")
```

Arguments

x	A text document.
language	A character setting the language to be used for stemming.

Details

The argument language is passed over to [SnowballStemmer](#) as the name of the snowball stemmer.

Examples

```
data("crude")
crude[[1]]
stemDocument(crude[[1]])
```

stopwords	<i>Multilingual Stopwords</i>
-----------	-------------------------------

Description

Return stopwords in different languages.

Usage

```
stopwords(language = "eng")
```

Arguments

language	A character giving the desired language.
----------	--

Details

At the moment supported languages are danish, dutch, english, finnish, french, german, hungarian, italian, norwegian, portuguese, russian, spanish, and swedish. Alternatively, their ISO 639-2 code may be used.

Value

A character vector containing the requested stopwords.

Examples

```
stopwords("eng")
```

stripWhitespace	<i>Strip Whitespace from a Text Document</i>
-----------------	--

Description

Strip extra whitespace from a text document. Multiple white space characters are collapsed to a single blank.

Usage

```
## S3 method for class 'PlainTextDocument':
stripWhitespace(x)
```

Arguments

x A text document.

Value

The text document with multiple white space characters collapse to a single blank.

See Also

[getTransformations](#) to list available transformation (mapping) functions.

Examples

```
data("crude")
crude[[1]]
stripWhitespace(crude[[1]])
```

TermDocumentMatrix *Term-Document Matrix*

Description

Constructs a term-document matrix or a document-term matrix.

Usage

```
TermDocumentMatrix(x, control = list())  
DocumentTermMatrix(x, control = list())
```

Arguments

<code>x</code>	a corpus
<code>control</code>	a named list of control options. The component <code>weighting</code> must be a weighting function capable of handling a <code>TermDocumentMatrix</code> . It defaults to <code>weightTf</code> for term frequency weighting. All other options are delegated internally to a termFreq call.

Value

An object of class `TermDocumentMatrix` or class `DocumentTermMatrix` containing a sparse term-document matrix or document-term matrix. The attribute `Weighting` contains the weighting applied to the matrix.

Author(s)

Ingo Feinerer

See Also

The documentation of [termFreq](#) gives an extensive list of possible options.

Available weighting functions shipped with the **tm** package are [weightTf](#), [weightTfIdf](#), and [weightBin](#).

Examples

```
data("crude")  
tdm <- TermDocumentMatrix(crude, control = list(weighting = weightTfIdf, stopwords = TRUE))  
dtm <- DocumentTermMatrix(crude, control = list(weighting = weightTfIdf, stopwords = TRUE))  
inspect(tdm[165:170,1:5])  
inspect(dtm[1:5,165:170])
```

termFreq	<i>Term Frequency Vector</i>
----------	------------------------------

Description

Generate a term frequency vector from a text document.

Usage

```
termFreq(doc, control = list())
```

Arguments

<code>doc</code>	An object inheriting from <code>TextDocument</code> .
<code>control</code>	A list of control options. Possible settings are <ul style="list-style-type: none"> tolower A function converting characters to lower case. Defaults to <code>base::tolower</code>. tokenize A function tokenizing documents to single tokens. Defaults to AlphabeticTokenizer (package RWeka), where tokens are to be formed only from contiguous alphabetic sequences. removeNumbers A logical value indicating whether numbers should be removed from <code>doc</code>. Defaults to <code>FALSE</code>. stemming Either a Boolean value indicating whether tokens should be stemmed or a stemming function. Defaults to <code>FALSE</code>. stopwords Either a Boolean value indicating stopword removal using default language specific stopword lists shipped with this package or a character vector holding custom stopwords. Defaults to <code>FALSE</code>. dictionary A character vector to be tabulated against. No other terms will be listed in the result. Terms from the dictionary not occurring in the document at all will be skipped for performance reasons. Defaults to no action (i.e., all terms are considered). minDocFreq An integer value. Words that appear less often in <code>doc</code> than this number are discarded. Defaults to 1 (i.e., every token will be used). minWordLength An integer value. Words smaller than this number are discarded. Defaults to length 3.

Value

A named integer vector with term frequencies as values and tokens as names.

Examples

```
data("crude")
termFreq(crude[[1]])
termFreq(crude[[1]], control = list(stemming = TRUE, minWordLength = 4))
```

TextDocument

*Access and Modify Text Documents***Description**

Access the meta data of text documents, and access and modify the content of text documents.

Details

The class `TextDocument` provides an abstraction over the concept of text documents and attached meta data which is stored in following attributes:

Author Object of class `character` containing the author names.

DateTimeStamp Object of class `POSIXlt` containing the date and time when the document was written.

Description Object of class `character` containing additional text information.

ID Object of class `character` containing an identifier.

Origin Object of class `character` containing information on the source and origin of the text.

Heading Object of class `character` containing the title or a short heading.

Language Object of class `character` containing the language of the text.

LocalMetaData Object of class `list` containing additional meta data in form of tag-value pairs which is local to each individual text document.

Author(s)

Ingo Feinerer

See Also

[meta](#) and [DublinCore](#) which provide a unified interface for meta data management.

TextRepository

*Text Repository***Description**

Construct a text repository for corpora. A repository is designed to keep track of multiple corpora, either different ones, or corpora with the same underlying texts but at different preprocessing stages.

Usage

```
TextRepository(x, meta = list(created = as.POSIXlt(Sys.time(), tz = "GMT")))
```

Arguments

x	A corpus.
meta	An initial list of tag-value pairs for the repository meta data.

Value

An object of class `TextRepository` which extends the class `list` containing corpora. Meta data annotations are stored in the attribute `RepoMetaData` in form of tag-value pairs (i.e., a named list).

Author(s)

Ingo Feinerer

Examples

```
data("crude")
repo <- TextRepository(crude)
summary(repo)
RepoMetaData(repo)
```

tm_cluster

Allow 'tm' to Use a Cluster

Description

tm can use a (MPI) cluster if available on your system.

Usage

```
tm_startCluster()
tm_stopCluster()
```

Details

`tm_startCluster` first investigates the MPI environment and tries to attach to a running MPI instance. If no MPI instance is found the function starts a new one. On success **tm** functions automatically use the cluster.

`tm_stopCluster` shuts down a running MPI instance.

Author(s)

Ingo Feinerer

tm_combine

*Combine Corpora and Documents***Description**

Combine several corpora into a single one or combine multiple documents into a corpus.

Usage

```
## S3 method for class 'Corpus':
c(x, ..., recursive = FALSE)
## S3 method for class 'TextDocument':
c(x, ..., recursive = FALSE)
```

Arguments

x	A corpus or a text document.
...	Corpora or text documents.
recursive	Logical. Provided by generic function definition but not used.

Details

Meta data from input objects (corpora or documents) is preserved during concatenation and intelligently merged into the newly created corpus. Although we use a sophisticated merging strategy (by using a binary tree for corpus specific meta data and by joining document level specific meta data in data frames) you should check the newly created meta data for consistency when merging corpora with (partly) identical meta data. However, in most cases the meta data merging strategy will produce validly combined and arranged meta data structures.

Examples

```
data("acq")
data("crude")
summary(c(acq, crude))
summary(c(acq[[30]], crude[[10]]))
```

tm_filter

*Filter and Index Functions on Corpora***Description**

Interface to apply filter and index functions to corpora.

Usage

```
## S3 method for class 'Corpus':
tm_filter(x, ..., FUN = searchFullText, doclevel = TRUE, useMeta = FALSE)
## S3 method for class 'Corpus':
tm_index(x, ..., FUN = searchFullText, doclevel = TRUE, useMeta = FALSE)
```

Arguments

x	A corpus.
...	Arguments to FUN.
FUN	A filter function returning a logical value.
doclevel	Logical. If the document level flag is set FUN is applied to each element of x, otherwise FUN is applied to x itself. If FUN has an attribute doclevel its value will be automatically used.
useMeta	Logical. Should DMetaData be passed over to FUN as argument?

Value

tm_filter returns a corpus containing documents where FUN matches, whereas tm_index only returns the corresponding indices.

See Also

[sFilter](#) for a filter using a simple statement query language, and [getFilters](#) to list available filter and index functions.

Examples

```
data("crude")
attr(searchFullText, "doclevel")
tm_filter(crude, FUN = searchFullText, "company")
tm_index(crude, FUN = searchFullText, "company")
```

tm_intersect

Intersection between Documents and Words

Description

Perform intersection on text documents and words.

Usage

```
## S3 method for class 'PlainTextDocument':
tm_intersect(x, y)
```


Arguments

x	A text document.
y	A character vector of words to be intersected with.

Value

A logical value indicating whether a word in y appears in x.

See Also

[getFilters](#) to list available filter functions.

Examples

```
data("crude")
crude[[1]]
tm_intersect(crude[[1]], c("crude", "oil"))
tm_intersect(crude[[1]], "acquisition")
```

tm_map

Transformations on Corpora

Description

Interface to apply transformation functions (also denoted as mappings) to corpora.

Usage

```
## S3 method for class 'PCorpus':
tm_map(x, FUN, ..., useMeta = FALSE, lazy = FALSE)
## S3 method for class 'VCorpus':
tm_map(x, FUN, ..., useMeta = FALSE, lazy = FALSE)
```

Arguments

x	A corpus.
FUN	A transformation function returning a text document.
...	Arguments to FUN.
useMeta	Logical. Should DMetaData be passed over to FUN as argument?
lazy	Logical. Lazy mappings are mappings which are delayed until the documents' content is accessed. Lazy mapping is useful when working with large corpora but only few documents will be accessed, as it avoids the computationally expensive application of the mapping to all elements in the corpus.

Value

A corpus with FUN applied to each document in x. In case of lazy mappings only annotations are stored which are evaluated upon access of individual documents which trigger the execution of the corresponding transformation function.

Note

Please be aware that lazy transformations are an experimental feature and change R's standard evaluation semantics.

See Also

[getTransformations](#) for available transformations, and [materialize](#) for manually triggering the materialization of documents with pending lazy transformations.

Examples

```
data("crude")
tm_map(crude, stemDocument)
## Generate a custom transformation function which takes the heading as new content
headings <- function(x) PlainTextDocument(Heading(x), id = ID(x), language = Language(x))
inspect(tm_map(crude, headings))
```

tm_reduce

Combine Transformations

Description

Fold multiple transformations (mappings) into a single one.

Usage

```
tm_reduce(x, tmFuns, ...)
```

Arguments

x	A corpus.
tmFuns	A list of tm transformations.
...	Arguments to the individual transformations.

Value

A single **tm** transformation function.

Author(s)

Ingo Feinerer

See Also

Reduce for R's internal folding/accumulation mechanism, and [getTransformations](#) to list available transformation (mapping) functions.

Examples

```
data(crude)
crude[[1]]
skipWords <- function(x) removeWords(x, c("it", "the"))
funs <- list(tolower, removePunctuation, skipWords, stripWhitespace)
tm_map(crude, FUN = tm_reduce, tmFuns = funs)[[1]]
```

URISource

Uniform Resource Identifier Source

Description

Constructs a source which represents a single document located by a uniform resource identifier.

Usage

```
URISource(x, encoding = "UTF-8")
```

Arguments

x	The Uniform Resource Identifier, i.e., either a character identifying the file or a connection.
encoding	A character giving the encoding of x.

Value

An object of class `URISource` which extends the class `Source` representing a single document located by a URI.

Author(s)

Ingo Feinerer

See Also

[DirSource](#) for accessing multiple files, and [getSources](#) to list available sources.

Examples

```
loreipsum <- system.file("texts", "loremipsum.txt", package = "tm")
us <- URISource(loreipsum)
inspect(Corpus(us))
```

VCorpus

*Volatile Corpus***Description**

Data structures and operators for volatile corpora.

Usage

```
Corpus(x, readerControl = list(reader = x$DefaultReader, language = "eng"), ...)
VCorpus(x, readerControl = list(reader = x$DefaultReader, language = "eng"), ...)
## S3 method for class 'VCorpus':
DMetaData(x)
## S3 method for class 'Corpus':
CMetaData(x)
```

Arguments

<code>x</code>	A Source object for <code>Corpus</code> and <code>VCorpus</code> , and a corpus for the other functions.
<code>readerControl</code>	A list with the named components <code>reader</code> representing a reading function capable of handling the file format found in <code>x</code> , and <code>language</code> giving the text's language (preferably in ISO 639-2 format).
<code>...</code>	Optional arguments for the <code>reader</code> .

Details

Volatile means that the corpus is fully kept in memory and thus all changes only affect the corresponding R object. In contrast there is also a corpus implementation available providing a permanent semantics (see [PCorpus](#)).

The constructed corpus object inherits from a `list` and has two attributes containing meta information:

CMetaData Corpus Meta Data contains corpus specific meta data in form of tag-value pairs and information about children in form of a binary tree. This information is useful for reconstructing meta data after e.g. merging corpora.

DMetaData Document Meta Data of class `data.frame` contains document specific meta data for the corpus. This data frame typically encompasses clustering or classification results which basically are metadata for documents but form an own entity (e.g., with its name, the value range, etc.).

Value

An object of class `VCorpus` which extends the classes `Corpus` and `list` containing a collection of text documents.

Author(s)

Ingo Feinerer

Examples

```
reut21578 <- system.file("texts", "crude", package = "tm")
(r <- Corpus(DirSource(reut21578), readerControl = list(reader = readReut21578XMLasPlain)))
```

VectorSource

Vector Source

Description

Constructs a source for a vector as input.

Usage

```
VectorSource(x, encoding = "UTF-8")
```

Arguments

<code>x</code>	A vector.
<code>encoding</code>	A character giving the encoding of <code>x</code> .

Value

An object of class `VectorSource` which extends the class `Source` representing a vector where each entry is interpreted as a document.

Author(s)

Ingo Feinerer

See Also

[getSources](#) to list available sources.

Examples

```
docs <- c("This is a text.", "This another one.")
(vs <- VectorSource(docs))
inspect(Corpus(vs))
```

weightBin

Weight Binary

Description

Binary weight a term-document matrix.

Usage

```
weightBin(m)
```

Arguments

`m` A [TermDocumentMatrix](#) in term frequency format.

Details

Formally this function is of class `WeightingFunction` with the additional attributes `Name` and `Acronym`.

Value

The weighted matrix.

Author(s)

Ingo Feinerer

WeightFunction

Weighting Function

Description

Construct a weighting function for term-document matrices.

Usage

```
WeightFunction(x, name, acronym)
```

Arguments

`x` A function which takes a [TermDocumentMatrix](#) with term frequencies as input, weights the elements, and returns the weighted matrix.

`name` A character naming the weighting function.

`acronym` A character giving an acronym for the name of the weighting function.

Value

An object of class `WeightFunction` which extends the class `function` representing a weighting function.

Author(s)

Ingo Feinerer

Examples

```
weightCutBin <- WeightFunction(function(m, cutoff) m > cutoff, "binary with cutoff", "bincut")
```

`weightTf`*Weight by Term Frequency*

Description

Weight a term-document matrix by term frequency.

Usage

```
weightTf(x)
```

Arguments

`x` A `TermDocumentMatrix` in term frequency format.

Details

Formally this function is of class `WeightingFunction` with the additional attributes `Name` and `Acronym`.

This function acts as the identity function since the input matrix is already in term frequency format.

Value

The weighted matrix.

Author(s)

Ingo Feinerer

weightTfIdf	<i>Weight by Term Frequency - Inverse Document Frequency</i>
-------------	--

Description

Weight a term-document matrix by term frequency - inverse document frequency.

Usage

```
weightTfIdf(m)
```

Arguments

m	A <code>TermDocumentMatrix</code> in term frequency format.
---	---

Details

Formally this function is of class `WeightingFunction` with the additional attributes `Name` and `Acronym`.

Value

The weighted matrix.

Author(s)

Ingo Feinerer

writeCorpus	<i>Write a Corpus to Disk</i>
-------------	-------------------------------

Description

Write a plain text representation of a corpus to multiple files on disk corresponding to the individual documents in the corpus.

Usage

```
writeCorpus(x, path = ".", filenames = NULL)
```

Arguments

x	A corpus.
path	A character listing the directory to be written into.
filenames	Either <code>NULL</code> or a character vector. In case no filenames are provided, filenames are automatically generated by using the documents' ID strings in x.

Examples

```
data("crude")
## Not run: writeCorpus(crude, path = ".", filenames = paste(seq_along(crude), ".txt", sep =
```

XMLSource

XML Source

Description

Constructs a source for an XML file.

Usage

```
XMLSource(x, parser, reader, encoding = "UTF-8")
```

Arguments

x	Either a character identifying a file or a connection.
parser	A function accepting an XML tree (as delivered by <code>xmlTreeParse</code> in package XML) as input and returning a list of XML elements.
reader	A function capable of turning XML elements as returned by <code>parser</code> into a subclass of TextDocument .
encoding	A character giving the encoding of x.

Value

An object of class `XMLSource` which extends the class [Source](#) representing an XML file.

Author(s)

Ingo Feinerer

See Also

Vignette 'Extensions: How to Handle Custom File Formats'.

Index

*Topic **datasets**

acq, 2
crude, 4

*Topic **file**

stopwords, 40

*Topic **math**

termFreq, 43

*Topic **methods**

removePunctuation, 32

acq, 2

AlphabeticTokenizer, 43

as.PlainTextDocument, 2

Author (TextDocument), 44

c.Corpus (tm_combine), 46

c.TextDocument (tm_combine), 46

CMetaData, 15

CMetaData (VCorpus), 51

colnames.DocumentTermMatrix
(names), 16

colnames.TermDocumentMatrix
(names), 16

Content (TextDocument), 44

Content<- (TextDocument), 44

convert_UTF_8, 3

Corpus (VCorpus), 51

crude, 4

CSVSource (DataframeSource), 5

DataframeSource, 5

DateTimeStamp (TextDocument), 44

DBControl (PCorpus), 17

Description (TextDocument), 44

Dictionary, 5

dim.DocumentTermMatrix (number),
17

dim.TermDocumentMatrix (number),
17

dimnames.DocumentTermMatrix
(names), 16

dimnames.TermDocumentMatrix
(names), 16

DirSource, 6, 50

dissimilarity, 7

DMetaData, 15, 21, 47, 48

DMetaData (VCorpus), 51

DMetaData.PCorpus (PCorpus), 17

DMetaData<- (VCorpus), 51

DMetaData<- .PCorpus (PCorpus), 17

Docs (names), 16

DocumentTermMatrix
(TermDocumentMatrix), 42

DublinCore, 44

DublinCore (meta), 15

DublinCore<- (meta), 15

eoi (Source), 38

findAssocs, 7

findFreqTerms, 8

FunctionGenerator, 9

getElem (Source), 38

getFilters, 10, 36, 37, 47, 48

getReaders, 9, 10, 23–28, 30, 31

getSources, 11, 39, 50, 52

getTransformations, 3, 4, 11, 32–34,
41, 49, 50

GmaneSource, 12

graphNEL, 20

gsub, 36

Heading (TextDocument), 44

iconv, 3, 4

ID (TextDocument), 44

inspect, 12

Language (TextDocument), 44

- LocalMetaData, 29, 30
- LocalMetaData (TextDocument), 44
- makeChunks, 13
- materialize, 14, 49
- meta, 15, 21, 44
- meta<- (meta), 15
- names, 16
- ncol.DocumentTermMatrix (number), 17
- ncol.TermDocumentMatrix (number), 17
- nDocs (number), 17
- nrow.DocumentTermMatrix (number), 17
- nrow.TermDocumentMatrix (number), 17
- nTerms (number), 17
- number, 17
- Origin (TextDocument), 44
- PCorpus, 17, 51
- pGetElem (Source), 38
- PlainTextDocument, 18, 22, 24, 35
- plot, 19
- preprocessReut21578XML, 20
- prescindMeta, 21
- RCV1Document, 22, 35
- readDOC, 23
- readGmane, 24
- readPDF, 25
- readPlain, 26
- readRCV1, 27
- readReut21578XML, 28
- readReut21578XMLasPlain (readReut21578XML), 28
- readTabular, 29
- readXML, 30
- removeNumbers, 32
- removePunctuation, 32
- removeSparseTerms, 33
- removeWords, 34
- RepoMetaData (TextRepository), 44
- Reuters21578Document, 22, 34
- ReutersSource, 35
- rownames.DocumentTermMatrix (names), 16
- rownames.TermDocumentMatrix (names), 16
- searchFullText, 36
- sFilter, 37, 47
- SnowballStemmer, 40
- Source, 17, 38, 51, 56
- stemCompletion, 39
- stemDocument, 40
- stepNext (Source), 38
- stopwords, 40
- stripWhitespace, 41
- TermDocMatrix (TermDocumentMatrix), 42
- TermDocumentMatrix, 42, 53–55
- termFreq, 42, 43
- Terms (names), 16
- TextDocument, 44, 56
- TextRepository, 44
- tm_cluster, 45
- tm_combine, 46
- tm_filter, 10, 46
- tm_index, 10
- tm_index (tm_filter), 46
- tm_intersect, 47
- tm_map, 11, 14, 48
- tm_reduce, 49
- tm_startCluster (tm_cluster), 45
- tm_stopCluster (tm_cluster), 45
- URISource, 50
- VCorpus, 51
- VectorSource, 52
- weightBin, 42, 53
- WeightFunction, 53
- weightTf, 42, 54
- weightTfIdf, 42, 55
- writeCorpus, 55
- XMLSource, 56