

Package ‘petro.One’

February 18, 2018

Type Package

Title Statistics and Text Mining for Oil and Gas Papers from OnePetro
Metadata

Version 0.1.3

Description Application that retrieves papers metadata from the OnePetro website. Thousands of papers on oil and gas live in OnePetro. By retrieving metadata from the search queries, a summary of papers that match the query words, can be retrieved for further analysis and text mining. There are some statistics and data mining provided such as word cloud plots, keywords frequency, conversion to corpus document, and removal of common usage words. OnePetro link: <<https://www.onepetro.org/>>.

Depends R (>= 3.2)

License GPL-2

Encoding UTF-8

LazyData true

RoxygenNote 6.0.1

Suggests testthat, knitr, rmarkdown, covr

Imports rvest, xml2, magrittr, tibble, tm, wordcloud, urltools, dplyr,
graph, Rgraphviz, ggplot2, RWeka, cluster, RColorBrewer,
data.table

VignetteBuilder knitr

URL <https://github.com/f0nzie/petro.One>

NeedsCompilation no

Author Alfonso R. Reyes [aut, cre]

Maintainer Alfonso R. Reyes <alfonso.reyes@oilgainsanalytics.com>

Repository CRAN

Date/Publication 2018-02-18 04:24:47 UTC

R topics documented:

| | |
|--------------------------------------|-----------|
| petro.One-package | 2 |
| custom_stopwords | 3 |
| custom_synonyms | 3 |
| get_papers_count | 4 |
| get_term_document_matrix | 4 |
| join_keywords | 5 |
| labels | 5 |
| make_search_url | 6 |
| onepetro_page_to_dataframe | 7 |
| papers_by_publication | 7 |
| papers_by_publisher | 8 |
| papers_by_type | 8 |
| papers_by_year | 9 |
| plot_bars | 9 |
| plot_cluster_dendrogram | 10 |
| plot_relationships | 10 |
| plot_wordcloud | 11 |
| read_multidoc | 11 |
| read_multipage | 12 |
| read_onepetro | 12 |
| summary_by_dates | 13 |
| summary_by_doctype | 13 |
| summary_by_publications | 14 |
| summary_by_publisher | 14 |
| term_frequency | 15 |
| term_frequency_n_grams | 15 |
| use_example | 15 |
| Index | 16 |

petro.One-package *Text mining and statistics for OnePetro papers petro.One*

Description

Text mining and statistics for OnePetro papers petro.One

| | |
|------------------|----------------------------------|
| custom_stopwords | <i>Default custom stop words</i> |
|------------------|----------------------------------|

Description

This is a minimal dataset of custom stopwords. You can supply your own stopwords by editing the file stopwords.txt under 'extdata' and then importing it. The provided dataset is a basic way to start and eliminate common words from the paper titles during classification.

Dataset: stopwords.rda

Source: stopwords.txt

Usage

```
custom_stopwords
```

Format

An object of class NULL of length 0.

| | |
|-----------------|-------------------------|
| custom_synonyms | <i>Default synonyms</i> |
|-----------------|-------------------------|

Description

Synonyms dataset to prevent repetition of equivalent words during classification. Example: 2D/2-D, cased hole/cased-hole, deep water/deepwater, etc. The first column is for the original word; the second column is for the replacement or standard word.

Dataset: synonyms.rda

Source: synonyms.txt

Usage

```
custom_synonyms
```

Format

An object of class data.frame with 33 rows and 2 columns.

get_papers_count *Number of paper for a given query*

Description

Obtains the number of papers being queried by the URL

Usage

```
get_papers_count(url)
```

Arguments

url char a query URL for OnePetro

Examples

```
# Example 1
url_1 <- make_search_url(query = "static gradient survey", how = "all")
get_papers_count(url_1)
#
# Example 2
url_2 <- make_search_url(query = "vertical lift performance", how = "all")
get_papers_count(url_2)
#
# Example 3
url_3 <- make_search_url(query = "inflow performance relationship", how = "all")
get_papers_count(url_3)
```

get_term_document_matrix
 A TermDocumentMatrix corpus objects

Description

Transforms a document into a VCorpus TermDocumentMatrix object plus additional calculated matrix, row sums and words frequency objects

Usage

```
get_term_document_matrix(df)
```

Arguments

df a dataframe with paper results

Value

a list

| | |
|---------------|---------------------------------------------------------------------------|
| join_keywords | <i>Get paper count and paper dataframe by joining keywords as vectors</i> |
|---------------|---------------------------------------------------------------------------|

Description

Get paper count and paper dataframe by joining keywords as vectors

Usage

```
join_keywords(..., get_papers = TRUE, bool_op = "AND", sleep = 3,  
             verbose = TRUE)
```

Arguments

| | |
|------------|------------------------------------------------|
| ... | input character vectors |
| get_papers | generate or not a dataframe with papers |
| bool_op | boolean operator. It can be AND or OR |
| sleep | seconds to wait before a new query to OnePetro |
| verbose | show progress if TRUE |

| | |
|--------|----------------------------------------------|
| labels | <i>Discipline and Subject labels dataset</i> |
|--------|----------------------------------------------|

Description

dataset containing disciplines and subjects. The purpose is to categorize papers based on the words in their title since OnePetro does not supply keywords or any sort of categorization. File: disciplines.rda Class: data.frame

Usage

```
labels
```

Format

An object of class `data.frame` with 254 rows and 2 columns.

make_search_url *Make a search URL for OnePetro*

Description

Create a URL that works in OnePetro website

Usage

```
make_search_url(query = NULL, start = NULL, from_year = NULL,
  peer_reviewed = NULL, published_between = NULL, rows = NULL,
  to_year = NULL, dc_type = NULL, how = "any")
```

Arguments

| | | |
|-------------------|---------|---------------------------------------------------|
| query | char | any words that will be searched |
| start | int | optional to set the starting paper |
| from_year | int | optional to indicate starting year |
| peer_reviewed | logical | optional, TRUE or FALSE |
| published_between | logical | automatic if from_year or to_year are on |
| rows | int | optional. number of papers to retrieve. max=1000 |
| to_year | int | optional to indicate end year |
| dc_type | char | optional to indicate if journal, conference paper |
| how | char | default="any". "all" will match exact words |

Examples

```
# Example 1
url_1 <- make_search_url(query = "flowing gradient survey", how = "all")
onepetro_page_to_dataframe(url_1)
# Example 2
url_2 <- make_search_url(query = "static gradient survey", how = "all")
onepetro_page_to_dataframe(url_2)
# Example 3
url_3 <- make_search_url(query = "downhole flowrate measurement",
  how = "all", from_year = 1982, to_year = 2017)
onepetro_page_to_dataframe(url_3)
```

`onepetro_page_to_dataframe`*Reads a OnePetro URL and converts it to a dataframe*

Description

A OnePetro URL with a query is read into a HTML page and converted to a dataframe

Usage

```
onepetro_page_to_dataframe(url)
```

Arguments

`url` char a OnePetro type URL

Examples

```
# Example 1
# Search papers with keyword "smartwell"
url_sw <- "https://www.onepetro.org/search?q=smartwell"
onepetro_page_to_dataframe(url_sw)
# Example 2
# Search for exact words ""vertical lift performance"
url_vlp <- "https://www.onepetro.org/search?q=%22vertical+lift+performance%22"
onepetro_page_to_dataframe(url_vlp)
```

`papers_by_publication` *Papers by publication*

Description

Generate a summary by publications. These publications could be World Petroleum Congress, Annual Technical Meeting, SPE Unconventional Reservoirs Conference, etc.

Usage

```
papers_by_publication(url)
```

Arguments

`url` a OnePetro query URL

Examples

```
#
# Example
my_url <- make_search_url(query = "industrial drilling", how = "all")
papers_by_publication(my_url)
```

papers_by_publisher *Papers by publisher*

Description

Generate a summary by publisher. Know publishers: OTC, SPE, etc.

Usage

```
papers_by_publisher(url)
```

Arguments

url a OnePetro query URL

Examples

```
#  
# Example  
my_url <- make_search_url(query = "shale gas", how = "all")  
papers_by_publisher(my_url)
```

papers_by_type *Summary by document type*

Description

Generate a summary by document type. Types are: conference-paper, journal-paper, presentation, media, other, etc.

Usage

```
papers_by_type(url)
```

Arguments

url a OnePetro page with results

Examples

```
#  
# Example 1  
my_url <- make_search_url(query = "well test", how = "all")  
papers_by_type(my_url)
```

`papers_by_year`*Papers by Year*

Description

Generate a summary by the year the paper was published

Usage

```
papers_by_year(url)
```

Arguments

`url` a OnePetro query URL

Examples

```
#  
# Example  
my_url <- make_search_url(query = "production automation", how = "all")  
papers_by_year(my_url)
```

`plotBars`*Plot frequency distribution with horizontal bars*

Description

Shows a bar plot with words on the y-axis and frequency on the x-axis

Usage

```
plotBars(df, min.freq = 25)
```

Arguments

`df` a dataframe with paper results
`min.freq` minimum frequency of the words to be plotted

plot_cluster_dendrogram
Plot a dendrogram

Description

Plots a clustering diagram of terms

Usage

```
plot_cluster_dendrogram(df)
```

Arguments

df a dataframe with paper results

plot_relationships *Plot a relationship diagram with weights*

Description

Plots a diagram with relationships between words. The lines that link the terms are weighted according to how often they connect together

Usage

```
plot_relationships(df, ..., min.freq = 25, threshold = 0.1)
```

Arguments

df a dataframe with paper results
... additional parameters
min.freq minimum frequency of the words to be plotted
threshold correlation threshold

| | |
|----------------|--------------------------|
| plot_wordcloud | <i>Plot a word cloud</i> |
|----------------|--------------------------|

Description

Plots a cloud plot of words where the size of the words is determined by their frequency

Usage

```
plot_wordcloud(df, ..., max.words = 200, min.freq = 50)
```

Arguments

| | |
|-----------|----------------------------------------|
| df | A dataframe with paper results |
| ... | other parameters |
| max.words | the maximum words to process |
| min.freq | the minimum frequency of words allowed |

| | |
|---------------|--------------------------------------------------------------|
| read_multidoc | <i>Read all OnePetro papers metadata by type of document</i> |
|---------------|--------------------------------------------------------------|

Description

Function iterates through all found document types and extracts papers into a common dataframe

Usage

```
read_multidoc(my_url)
```

Arguments

| | |
|--------|--------------------|
| my_url | OnePetro query URL |
|--------|--------------------|

| | |
|----------------|------------------------------------------------|
| read_multipage | <i>Reads metadata in groups of 1000 papers</i> |
|----------------|------------------------------------------------|

Description

This function will loop over and grab data from the OnePetro results in groups of 1000 papers at a time. OnePetro limits the number of papers to view to 1000 papers and the query in this function automatically sets the start counter to read them in groups.

Usage

```
read_multipage(url, doctype = NULL)
```

Arguments

| | |
|---------|----------------------------------------------------------------------------------------------|
| url | A OnePetro query URL |
| doctype | a OnePetro paper type: conference-paper, journal-paper, general. presentation, chapter, etc. |

| | |
|---------------|-------------------------------------------------|
| read_onepetro | <i>Read OnePetro web page given a query URL</i> |
|---------------|-------------------------------------------------|

Description

Read a OnePetro webpage using a query URL. Uses xml2 function read_html

Usage

```
read_onepetro(url)
```

Arguments

| | |
|-----|-------------------------------|
| url | char a query URL for OnePetro |
|-----|-------------------------------|

| | |
|------------------|------------------------|
| summary_by_dates | <i>Summary by year</i> |
|------------------|------------------------|

Description

Generate a summary by the year the paper was published

Usage

```
summary_by_dates(result)
```

Arguments

result a OnePetro page with results

Examples

```
#  
# Example  
my_url <- make_search_url(query = "production automation", how = "all")  
result <- read_onepetro(my_url)  
summary_by_dates(result)
```

| | |
|--------------------|---------------------------------|
| summary_by_doctype | <i>Summary by document type</i> |
|--------------------|---------------------------------|

Description

Generate a summary by document type. Types are: conference-paper, journal-paper, presentation, media, other, etc.

Usage

```
summary_by_doctype(result)
```

Arguments

result a OnePetro page with results

Examples

```
#  
# Example 1  
my_url <- make_search_url(query = "well test", how = "all")  
result <- read_onepetro(my_url)  
summary_by_doctype(result)
```

summary_by_publications

Summary by publication

Description

Generate a summary by publications. These publications could be World Petroleum Congress, Annual Technical Meeting, SPE Unconventional Reservoirs Conference, etc.

Usage

```
summary_by_publications(result)
```

Arguments

result a OnePetro page with results

Examples

```
#  
# Example  
my_url <- make_search_url(query = "industrial drilling", how = "all")  
result <- read_onepetro(my_url)  
summary_by_publications(result)
```

summary_by_publisher *Summary by publisher*

Description

Generate a summary by publisher. Know publishers: OTC, SPE, etc.

Usage

```
summary_by_publisher(result)
```

Arguments

result a OnePetro page with results

Examples

```
#  
# Example  
my_url <- make_search_url(query = "shale gas", how = "all")  
page <- read_onepetro(my_url)  
summary_by_publisher(page)
```

| | |
|----------------|---------------------------------|
| term_frequency | <i>Word Frequency Dataframe</i> |
|----------------|---------------------------------|

Description

Returns a dataframe of words vs frequency

Usage

```
term_frequency(df)
```

Arguments

| | |
|----|--------------------------------|
| df | a dataframe with paper results |
|----|--------------------------------|

| | |
|------------------------|----------------------------------------------------------|
| term_frequency_n_grams | <i>Find the frequency for two or more words together</i> |
|------------------------|----------------------------------------------------------|

Description

Use this function when trying to find frequency of two or more words

Usage

```
term_frequency_n_grams(df, gram.min = 2, gram.max = 2)
```

Arguments

| | |
|----------|----------------------------------|
| df | a dataframe with paper results |
| gram.min | minimum amount of words together |
| gram.max | maximum amount of words together |

| | |
|-------------|--------------------------|
| use_example | <i>Unpack an example</i> |
|-------------|--------------------------|

Description

Examples are zipped to save disk space and prevent R complaining while creating the package

Usage

```
use_example(which_one = NULL)
```

Arguments

| | |
|-----------|-----------------------|
| which_one | example number to use |
|-----------|-----------------------|

Index

*Topic **datasets**

- custom_stopwords, 3
- custom_synonyms, 3
- labels, 5

- custom_stopwords, 3
- custom_synonyms, 3

- get_papers_count, 4
- get_term_document_matrix, 4

- join_keywords, 5

- labels, 5

- make_search_url, 6

- onepetro_page_to_dataframe, 7

- papers_by_publication, 7
- papers_by_publisher, 8
- papers_by_type, 8
- papers_by_year, 9
- petro.One-package, 2
- plot_bars, 9
- plot_cluster_dendrogram, 10
- plot_relationships, 10
- plot_wordcloud, 11

- read_multidoc, 11
- read_multipage, 12
- read_onepetro, 12

- summary_by_dates, 13
- summary_by_doctype, 13
- summary_by_publications, 14
- summary_by_publisher, 14

- term_frequency, 15
- term_frequency_n_grams, 15

- use_example, 15