

Package ‘datacollector’

August 29, 2017

Title R package for automatic data collection

Version 0.1

Author Med Dhaker Abdeljawed

Maintainer Med Dhaker Abdeljawed <med.dhaker.abdeljawed@gmail.com>

Description A suite of functions and tools made to automate data collection from the internet, social media networks and webpages.

License Free

Imports rJava, stringr, xlsx

Depends R (>= 3.2.3)

Encoding UTF-8

LazyData true

RoxygenNote 6.0.1

Repository github

Prerequisite

1- Setup java environment including JAVA_HOME and jre/jdk path.

2- Setup Apache Zookeeper Download Zookeeper from
<https://zookeeper.apache.org/releases.html#download>, Start
server instance :
https://zookeeper.apache.org/doc/r3.3.3/zookeeperStarted.html#sc_InstallingSingleMode

3- Setup Apache Drill download drill from
<https://drill.apache.org/download/> and start server instance :
<https://drill.apache.org/docs/starting-drill-in-distributed-mode/>
Note : the package should run in the same node with drillbit
instance

R topics documented:

collectContentByKey	2
collectEmails	2
collectFacebookNodes	3

collectFromYoutube	3
collectMedias	4
collectPhones	4
collectSocialLinks	5
collectWebPageData	5
collectWebSiteData	6
searchFacebookNodes	6
searchFromFlickr	7
searchFromGooglePlus	7
searchFromTumblr	8
searchFromTwitter	8
searchFromYouTube	9
Index	10

collectContentByKey	<i>Collect by keyword links from a web page</i>
---------------------	-------------------------------------------------

Description

This function extract a webpage content by keyword, and return the blocs containing that keyword

Usage

```
collectContentByKey(url, key)
```

Arguments

url	url of the target web page
key	keyword to search for

Value

list of bloc content containing that keyword

collectEmails	<i>Collect emails from a web page</i>
---------------	---------------------------------------

Description

this function extract all emails from a web page and return it as an R list

Usage

```
collectEmails(url)
```

Arguments

url url of the target web page

Value

list of emails

collectFacebookNodes *Collect facebook data*

Description

This function collect data from a facebook node that could be a page, group, event or else based on a given edge which represents the type of data to collect such, posts and feed then parse it and save into an R dataframe and spreadsheet excel file.

Usage

```
collectFacebookNodes(nodeid, edge = "feed", user = FALSE, rootPath)
```

Arguments

nodeid Id of the target facebook node

user FALSE to use an application based call, TRUE to use an Autenticated call (used to fetch private groups and events)

rootPath path to save the generated results

collectFromYoutube *Collect youtube data*

Description

This function collect data about a specific youtube element attached to videos and channels including activities, playlists, subscriptions, comments and video captions tracks then parse it and save it into an R dataframe and an excel spreadsheet file.

Usage

```
collectFromYoutube(type, id, rootPath)
```

Arguments

type element type to fetch ("activity", "playlist", "comment", "subscription", "caption")

id a video id or a channel id, video id for ("caption", "comment") and channel id for the rest

rootPath path to save the generated results

Value

R dataframe representing the result node

collectMedias	<i>Collect media links from a web page</i>
---------------	--------------------------------------------

Description

this function extract all medias such images, videos, audio, documents.. from a web page

Usage

```
collectMedias(url, media)
```

Arguments

url	url of the target web page
media	media type could be : image,video,document,audio

Value

list of medias

collectPhones	<i>Collect phone numbers from a web page</i>
---------------	----------------------------------------------

Description

this function extract all phone numbers from a web page and return it as an R list

Usage

```
collectPhones(url)
```

Arguments

url	url of the target web page
-----	----------------------------

Value

list of phone numbers

collectSocialLinks	<i>Collect SocialLinks</i>
--------------------	----------------------------

Description

this function extract all social medias links from a web page and return it as an R list

Usage

```
collectSocialLinks(url)
```

Arguments

url	url of the target web page
-----	----------------------------

Value

list of social media links

collectWebPageData	<i>Collect data from a web page</i>
--------------------	-------------------------------------

Description

This function scraps data from a web page and save it into an excel sheet.

Usage

```
collectWebPageData(excelpath, url)
```

Arguments

excelpath	Path to the excel file
url	url of the target web page

Details

Data such : Email, phones numbers, medias, links, social media links, and html tables

Value

void

collectWebSiteData	<i>Collect data from a whole website</i>
--------------------	------------------------------------------

Description

collectWebSiteData crawls a whole website for the requested data and save each scraped page into an excel sheet, and download the available media files in a specific directory

Usage

```
collectWebSiteData(url, excelpath, downloadpath)
```

Arguments

url	url of the target web page
excelpath	Path to the excel file
downloadpath	path to the downloaded media files

Value

void

searchFacebookNodes	<i>Search facebook data</i>
---------------------	-----------------------------

Description

This function searches facebook nodes including groups, events and pages; then parse the results and save into an R dataframe and an excel spreadsheet file.

Usage

```
searchFacebookNodes(query, rootPath)
```

Arguments

query	text query to use for the search
rootPath	path to save the generated results

searchFromFlickr	<i>Search Flickr data</i>
------------------	---------------------------

Description

This function search for photos from Flickr based on a text query, then parse the result into an R dataframe and an excel spreadSheet

Usage

```
searchFromFlickr(query, rootPath)
```

Arguments

query	text query to use for the search
rootPath	path to save the generated results

searchFromGooglePlus	<i>Search google plus data</i>
----------------------	--------------------------------

Description

This function searches google plus activities using a text query, then parse the result into an R dataframe and spreadSheet xls

Usage

```
searchFromGooglePlus(query, rootPath)
```

Arguments

query	text query to use for the search
rootPath	path to save the generated results

searchFromTumblr	<i>Search Tumblr data</i>
------------------	---------------------------

Description

This function search for blog posts from Tumblr using a text query, then parse the result into an R dataframe and spreadSheet xls

Usage

```
searchFromTumblr(query, rootPath)
```

Arguments

query	text query to use for the search
rootPath	path to save the generated results

searchFromTwitter	<i>Search Twitter data</i>
-------------------	----------------------------

Description

This function searches Tweets of the past seven days (week) from twitter using a text query, then parse the result into an R dataframe and also creates an excel spreadSheet

Usage

```
searchFromTwitter(query, rootPath)
```

Arguments

query	text query to use for the search
rootPath	path to save the generated results

searchFromYouTube	<i>Search Youtube data</i>
-------------------	----------------------------

Description

This function search for videos and channels from youtube using a text query, then parse the result into an R dataframe and also an excel spreadSheet

Usage

```
searchFromYouTube(query, rootPath)
```

Arguments

query	text query to use for the search
rootPath	path to save the generated results

Index

collectContentByKey, [2](#)
collectEmails, [2](#)
collectFacebookNodes, [3](#)
collectFromYoutube, [3](#)
collectMedias, [4](#)
collectPhones, [4](#)
collectSocialLinks, [5](#)
collectWebPageData, [5](#)
collectWebSiteData, [6](#)

searchFacebookNodes, [6](#)
searchFromFlickr, [7](#)
searchFromGooglePlus, [7](#)
searchFromTumblr, [8](#)
searchFromTwitter, [8](#)
searchFromYouTube, [9](#)