

Package ‘crminer’

October 15, 2018

Type Package

Title Fetch 'Scholarly' Full Text from 'Crossref'

Description Text mining client for 'Crossref' (<<https://crossref.org>>). Includes functions for getting getting links to full text of articles, fetching full text articles from those links or Digital Object Identifiers ('DOIs'), and text extraction from 'PDFs'.

Version 0.2.0

License MIT + file LICENSE

URL <https://github.com/ropensci/crminer>

BugReports <https://github.com/ropensci/crminer/issues>

LazyData true

Encoding UTF-8

Imports crul (>= 0.3.4), jsonlite (>= 1.4), xml2 (>= 1.1.0), pdftools (>= 1.2), hoardr (>= 0.5.0)

Suggests roxygen2 (>= 6.0.1), rcrossref, testthat

RoxygenNote 6.1.0

X-schema.org-applicationCategory Literature

X-schema.org-keywords text-mining, literature, publications, crossref, pdf, xml

X-schema.org-isPartOf <https://ropensci.org>

NeedsCompilation no

Author Scott Chamberlain [aut, cre] (<<https://orcid.org/0000-0003-1444-9135>>)

Maintainer Scott Chamberlain <myrmecocystus+r@gmail.com>

Repository CRAN

Date/Publication 2018-10-15 17:20:03 UTC

R topics documented:

crminer-package	2
as_tdmurl	3
auth	4
crm_cache	5
crm_extract	6
crm_html	7
crm_links	8
crm_pdf	9
crm_plain	11
crm_text	12
crm_xml	15
dois_crminer	16
dois_crminer_ccby3	16
dois_elsevier	17
dois_hindawi	17
dois_pensoft	17
dois_wiley	17
Index	18

crminer-package	<i>crminer - Crossref text miner</i>
-----------------	--------------------------------------

Description

crminer - Crossref text miner

Package API (functions)

- `crm_links()` - get full text links from DOIs
- `as_tdmurl()` - coerce a URL to a tdmurl
- `crm_text()` - general purpose interface to request full text
- `crm_pdf()` - request pdf full text
- `crm_plain()` - request plain text full text
- `crm_xml()` - request xml full text
- `crm_extract()` - extract text from a pdf

Authentication

You should first start reading up on authentication (`auth()`) since you are probably here to do text mining, and most of the full text links available via Crossref are behind authentication.

Author(s)

Scott Chamberlain <myrmecocystus+r@gmail.com>

`as_tdmurl`*Coerce a url to a tdmurl with a specific type*

Description

A tmd url is just a URL with some attributes to make it easier to handle within other functions in this package.

Usage

```
as_tdmurl(url, type, doi = NULL, member = NULL,  
          intended_application = NULL)
```

```
## S3 method for class 'tdmurl'  
as_tdmurl(url, type, doi = NULL, member = NULL,  
          intended_application = NULL)
```

```
## S3 method for class 'character'  
as_tdmurl(url, type, doi = NULL, member = NULL,  
          intended_application = NULL)
```

Arguments

<code>url</code>	(character) A URL.
<code>type</code>	(character) One of 'xml' (default), 'html', 'plain', 'pdf', 'unspecified', or 'all'
<code>doi</code>	(character) A DOI, optional, default: NULL
<code>member</code>	(character) Crossref member id. optional
<code>intended_application</code>	(character) intended application string, optional

Examples

```
as_tdmurl("http://downloads.hindawi.com/journals/bmri/2014/201717.xml",  
          "xml")  
as_tdmurl("http://downloads.hindawi.com/journals/bmri/2014/201717.pdf",  
          "pdf")  
out <-  
  as_tdmurl("http://downloads.hindawi.com/journals/bmri/2014/201717.pdf",  
            "pdf", "10.1155/2014/201717")  
attributes(out)  
identical(attr(out, "type"), "pdf")
```

auth

Crossref TDM authentication

Description

Crossref TDM authentication

Authentication

There's a set of publishers that are involved in the Crossref Text and Data Mining (TDM) program (<http://tdmsupport.crossref.org/>), which means essentially the publishers deposit URLs for fulltext in Crossref metadata.

Authentication is applicable only when the publisher you want to get fulltext from requires it. OA publishers shouldn't need it as you'd expect. There's many publishers that don't share links at all, so they are irrelevant here.

For publishers that required authentication, the Crossref TDM program allows for a single token to authenticate across publishers (to make it easier for text miners). The publishers involved with the authentication scheme are really only Elsevier and Wiley.

There's a how to guide for Crossref TDM at <http://tdmsupport.crossref.org/researchers/>. Get your Crossref TDM token by registering at <https://apps.crossref.org/clickthrough/researchers>. Save the token in your `.Renvi ron` file with a new row like `CROSSREF_TDM=your key`. We will read that key in for you - it's best this way rather than passing in a key via a parameter - since you might put that code up on the web somewhere and someone could use your key.

IP addresses

If you don't know what IP addresses are, check out https://en.wikipedia.org/wiki/IP_address. At least Elsevier and I think Wiley also check your IP address in addition to requiring the authentication token. Usually you're good if you're physically at the institution that has access to the publishers content OR on a VPN (i.e., pretending you're there).

If you forget about this, you'll get errors about not being authorized. So check and make sure you're on a VPN if you're not physically located at your institution.

Fences

There's yet another issue to worry about. At least with Elsevier, they have a so-called "fence" - that is, even if an institution has access to Elsevier content, Elsevier doesn't necessarily have the fence turned off - if its not off, you can't get through - if it's off, you can. If you have the right token and you are sure you're on the right IP address, this could be the problem for your lack of access.

HELP!

If you're having trouble with any of this, get in touch with the package maintainer.

Description

Manage cached crminer files with **hoardr**

Details

The default cache directory is `paste0(rappdirs::user_cache_dir(), "/R/crminer")`, but you can set your own path using `cache_path_set()`

`cache_delete` only accepts 1 file name, while `cache_delete_all` doesn't accept any names, but deletes all files. For deleting many specific files, use `cache_delete` in a `lapply()` type call

Useful user functions

- `crm_cache$cache_path_get()` get cache path
- `crm_cache$cache_path_set()` set cache path. You can set the entire path directly via the `full_path` arg like `crm_cache$cache_path_set(full_path = "your/path")`
- `crm_cache$list()` returns a character vector of full path file names
- `crm_cache$files()` returns file objects with metadata
- `crm_cache$details()` returns files with details
- `crm_cache$delete()` delete specific files
- `crm_cache$delete_all()` delete all files, returns nothing

Examples

```
## Not run:
crm_cache

# list files in cache
crm_cache$list()

# delete certain database files
# crm_cache$delete("file path")
# crm_cache$list()

# delete all files in cache
# crm_cache$delete_all()
# crm_cache$list()

## End(Not run)
```

crm_extract	<i>Extract text from a single pdf document</i>
-------------	--

Description

Extract text from a single pdf document

Usage

```
crm_extract(path = NULL, raw = NULL, ...)
```

Arguments

path	(character) path to a file, file must exist
raw	(raw) raw bytes
...	args passed on to <code>pdftools::pdf_info()</code> and <code>pdftools::pdf_text()</code> - any args are passed to both of those function calls, which makes sense

Details

We use **pdftools** under the hood to do pdf text extraction.

You have to supply either path or raw - not both.

Value

An object of class `crm_pdf` with a slot for `info` (pdf metadata essentially), and `text` (the extracted text) - with an attribute (`path`) with the path to the pdf on disk

Examples

```
path <- system.file("examples", "MairChamberlain2014RJournal.pdf",
  package = "crminer")
(res <- crm_extract(path))
res$info
res$text
# with newlines, pretty print
cat(res$text)

# another example
path <- system.file("examples", "ChamberlainEtal2013Ecosphere.pdf",
  package = "crminer")
(res <- crm_extract(path))
res$info
cat(res$text)

# with raw pdf bytes
path <- system.file("examples", "raw-example.rds", package = "crminer")
rds <- readRDS(path)
```

```
class(rds)
crm_extract(raw = rds)
```

crm_html	<i>Get full plain text</i>
----------	----------------------------

Description

Get full plain text

Usage

```
crm_html(url, overwrite_unspecified = FALSE, ...)
```

Arguments

url	A URL (character) or an object of class <code>tdmurl</code> from a call to <code>crm_links()</code> . If you'll be getting text from the publishers are use Crossref TDM (which requires authentication), we strongly recommend using <code>crm_links()</code> first and passing output of that here, as <code>crm_links()</code> grabs the publisher Crossref member ID, which we use to do authentication and other publisher specific fixes to URLs
overwrite_unspecified	(logical) Sometimes the crossref API returns mime type 'unspecified' for the full text links (for some Wiley dois for example). This parameter overrides the mime type to be type.
...	Named curl parameters passed on to <code>curl::HttpClient()</code> , see <code>curl::curl_options()</code> for available curl options

Details

Note that this function is not vectorized. To do many requests use a for/while loop or lapply family calls, or similar.

Note that some links returned will not in fact lead you to full text content as you would understandably think and expect. That is, if you use the `filter` parameter with e.g., `rcrossref::cr_works()` and filter to only full text content, some links may actually give back only metadata for an article. Elsevier is perhaps the worst offender, for one because they have a lot of entries in Crossref TDM, but most of the links that are apparently full text are not in fact full text, but only metadata.

Check out [auth](#) for details on authentication.

Examples

```
## Not run:
link <- crm_links("10.7717/peerj.1545", "html")
crm_html(link)

link <- crm_links("10.7717/peerj.1545")
crm_html(link)
```

```
crm_html("https://peerj.com/articles/1545.html")
## End(Not run)
```

crm_links	<i>Get full text links from a DOI</i>
-----------	---------------------------------------

Description

Get full text links from a DOI

Usage

```
crm_links(doi, type = "all", ...)
```

Arguments

doi	(character) A Digital Object Identifier (DOI). required.
type	(character) One of 'xml', 'html', 'plain', 'pdf', 'unspecified', or 'all' (default). required.
...	Named parameters passed on to <code>curl::HttpClient()</code>

Details

Note that this function is not vectorized.

Some links returned will not in fact lead you to full text content as you would understandably think and expect. That is, if you use the `filter` parameter with e.g., `rcrossref::cr_works()` and filter to only full text content, some links may actually give back only metadata for an article. Elsevier is perhaps the worst offender, for one because they have a lot of entries in Crossref TDM, but most of the links that are apparently full text are not in fact full text, but only metadata. You can get full text if you are part of a subscribing institution to that specific Elsevier content, but otherwise, you're SOL.

Note that there are still some bugs in the data returned from CrossRef. For example, for the publisher eLife, they return a single URL with content-type application/pdf, but the URL is not for a PDF, but for both XML and PDF, and content-type can be set with that URL as either XML or PDF to get that type.

In another example, all Elsevier URLs at time of writing are have http scheme, while those don't actually work, so we have a custom fix in this function for that publisher. Anyway, expect changes...

Value

NULL if no full text links given; a list of `tdmurl` objects if links found. a `tdmurl` object is an S3 class wrapped around a simple list, with attributes for:

- `type`: type, matchin type passed to the function

- doi: DOI
- member: Crossref member ID
- intended_application: intended application, e.g., text-mining

Examples

```
## Not run:
data(dois_crminer)

# pdf link
crm_links(doi = "10.5555/515151", "pdf")

# xml and plain text links
crm_links(dois_crminer[1], "pdf")
crm_links(dois_crminer[6], "xml")
crm_links(dois_crminer[7], "plain")
crm_links(dois_crminer[1]) # all is the default

# pdf link
crm_links(doi = "10.5555/515151", "pdf")
crm_links(doi = "10.3897/phytokeys.52.5250", "pdf")

# many calls, use e.g., lapply
lapply(dois_crminer[1:3], crm_links)

## End(Not run)
```

crm_pdf

Get full text PDFs

Description

Get full text PDFs

Usage

```
crm_pdf(url, overwrite = TRUE, read = TRUE, cache = FALSE,
        overwrite_unspecified = FALSE, ...)
```

Arguments

url	A URL (character) or an object of class <code>tdmurl</code> from a call to <code>crm_links()</code> . If you'll be getting text from the publishers are use Crossref TDM (which requires authentication), we strongly recommend using <code>crm_links()</code> first and passing output of that here, as <code>crm_links()</code> grabs the publisher Crossref member ID, which we use to do authentication and other publisher specific fixes to URLs
overwrite	(logical) Overwrite file if it exists already? Default: TRUE

read	(logical) If reading a pdf, this toggles whether we extract text from the pdf or simply download. If TRUE, you get the text from the pdf back. If FALSE, you only get back the metadata. Default: TRUE
cache	(logical) Use cached files or not. All files are written to your machine locally, so this doesn't affect that. This only states whether you want to use cached version so that you don't have to download the file again. The steps of extracting and reading into R still have to be performed when cache=TRUE. Default: TRUE
overwrite_unspecified	(logical) Sometimes the crossref API returns mime type 'unspecified' for the full text links (for some Wiley dois for example). This parameter overrides the mime type to be type.
...	Named curl parameters passed on to <code>curl::HttpClient()</code> , see <code>curl::curl_options()</code> for available curl options

Details

Note that this function is not vectorized. To do many requests use a for/while loop or lapply family calls, or similar.

Note that some links returned will not in fact lead you to full text content as you would understandably think and expect. That is, if you use the `filter` parameter with e.g., `rcrossref::cr_works()` and filter to only full text content, some links may actually give back only metadata for an article. Elsevier is perhaps the worst offender, for one because they have a lot of entries in Crossref TDM, but most of the links that are apparently full text are not in fact full text, but only metadata.

Check out [auth](#) for details on authentication.

Caching

By default we use `paste0(rappdirs::user_cache_dir(), "/crminer")`, but you can set this directory to something different. Ignored unless getting pdf. See [crm_cache](#) for caching details.

Examples

```
## Not run:
# set a temp dir. cache path
crm_cache$cache_path_set(path = "crminer", type = "tempdir")
## you can set the entire path directly via the `full_path` arg
## like crm_cache$cache_path_set(full_path = "your/path")

## peerj
x <- crm_pdf("https://peerj.com/articles/2356.pdf")

## pensoft
data(dois_pensoft)
(links <- crm_links(dois_pensoft[1], "all"))
### pdf
crm_text(url=links, type="pdf", read = FALSE)
crm_text(links, "pdf")

## hindawi
```

```

data(dois_pensoft)
(links <- crm_links(dois_pensoft[1], "all"))
### pdf
crm_text(links, "pdf", read=FALSE)
crm_text(links, "pdf")

### Caching, for PDFs
# out <- cr_members(2258, filter=c(has_full_text = TRUE), works = TRUE)
# (links <- crm_links(out$data$DOI[10], "all"))
# crm_text(links, type = "pdf", cache=FALSE)
# system.time( cacheyes <- crm_text(links, type = "pdf", cache=TRUE) )
### second time should be faster
# system.time( cacheyes <- crm_text(links, type = "pdf", cache=TRUE) )
# system.time( cacheno <- crm_text(links, type = "pdf", cache=FALSE) )
# identical(cacheyes, cacheno)

## End(Not run)

```

crm_plain

Get full plain text

Description

Get full plain text

Usage

```
crm_plain(url, overwrite_unspecified = FALSE, ...)
```

Arguments

url	A URL (character) or an object of class <code>tdmur1</code> from a call to <code>crm_links()</code> . If you'll be getting text from the publishers are use Crossref TDM (which requires authentication), we strongly recommend using <code>crm_links()</code> first and passing output of that here, as <code>crm_links()</code> grabs the publisher Crossref member ID, which we use to do authentication and other publisher specific fixes to URLs
overwrite_unspecified	(logical) Sometimes the crossref API returns mime type 'unspecified' for the full text links (for some Wiley dois for example). This parameter overrides the mime type to be type.
...	Named curl parameters passed on to <code>curl::HttpClient()</code> , see <code>curl::curl_options()</code> for available curl options

Details

Note that this function is not vectorized. To do many requests use a for/while loop or lapply family calls, or similar.

Note that some links returned will not in fact lead you to full text content as you would understandably think and expect. That is, if you use the `filter` parameter with e.g., `rcrossref::cr_works()` and filter to only full text content, some links may actually give back only metadata for an article. Elsevier is perhaps the worst offender, for one because they have a lot of entries in Crossref TDM, but most of the links that are apparently full text are not in fact full text, but only metadata.

Check out [auth](#) for details on authentication.

Examples

```
## Not run:
link <- crm_links("10.1016/j.physletb.2010.10.049", "plain")
crm_plain(link)

# another eg, which requires Crossref TDM authentication, see ?auth
link <- crm_links("10.1016/j.funeco.2010.11.003", "plain")
# crm_plain(link)

## End(Not run)
```

crm_text

Get full text

Description

Get full text

Usage

```
crm_text(url, type = "xml", overwrite = TRUE, read = TRUE,
  cache = FALSE, overwrite_unspecified = FALSE, ...)
```

Arguments

<code>url</code>	A URL (character) or an object of class <code>tdmurl</code> from a call to <code>crm_links()</code> . If you'll be getting text from the publishers are use Crossref TDM (which requires authentication), we strongly recommend using <code>crm_links()</code> first and passing output of that here, as <code>crm_links()</code> grabs the publisher Crossref member ID, which we use to do authentication and other publisher specific fixes to URLs
<code>type</code>	(character) One of 'xml' (default), 'html', 'plain', 'pdf', 'unspecified'
<code>overwrite</code>	(logical) Overwrite file if it exists already? Default: TRUE
<code>read</code>	(logical) If reading a pdf, this toggles whether we extract text from the pdf or simply download. If TRUE, you get the text from the pdf back. If FALSE, you only get back the metadata. Default: TRUE
<code>cache</code>	(logical) Use cached files or not. All files are written to your machine locally, so this doesn't affect that. This only states whether you want to use cached version so that you don't have to download the file again. The steps of extracting and reading into R still have to be performed when <code>cache=TRUE</code> . Default: TRUE

```

overwrite_unspecified
    (logical) Sometimes the crossref API returns mime type 'unspecified' for the
    full text links (for some Wiley dois for example). This parameter overrides the
    mime type to be type.
...
    Named curl parameters passed on to curl::HttpClient(), see curl::curl_options()
    for available curl options

```

Details

Note that this function is not vectorized. To do many requests use a for/while loop or lapply family calls, or similar.

Note that some links returned will not in fact lead you to full text content as you would understandably think and expect. That is, if you use the `filter` parameter with e.g., `rcrossref::cr_works()` and `filter` to only full text content, some links may actually give back only metadata for an article. Elsevier is perhaps the worst offender, for one because they have a lot of entries in Crossref TDM, but most of the links that are apparently full text are not in fact full text, but only metadata.

Check out [auth](#) for details on authentication.

Caching

By default we use `paste0(rappdirs::user_cache_dir(), "/crminer")`, but you can set this directory to something different. Ignored unless getting pdf. See [crm_cache](#) for caching details.

Examples

```

## Not run:
# set a temp dir. cache path
crm_cache$cache_path_set(path = "crminer", type = "tempdir")
## you can set the entire path directly via the `full_path` arg
## like crm_cache$cache_path_set(full_path = "your/path")

## pensoft
data(dois_pensoft)
(links <- crm_links(dois_pensoft[1], "all"))
### xml
crm_text(url=links, type='xml')
### pdf
crm_text(url=links, type="pdf", read = FALSE)
crm_text(links, "pdf")

## hindawi
data(dois_pensoft)
(links <- crm_links(dois_pensoft[1], "all"))
### xml
crm_text(links, 'xml')
### pdf
crm_text(links, "pdf", read=FALSE)
crm_text(links, "pdf")

## DOIs w/ full text, and with CC-BY 3.0 license
data(dois_crminer_ccby3)

```

```

(links <- crm_links(dois_crminer_ccby3[40], "all"))
# crm_text(links, 'pdf')

## You can use crm_xml, crm_plain, and crm_pdf to go directly to
## that format
(links <- crm_links(dois_crminer_ccby3[5], "all"))
crm_xml(links)

### Caching, for PDFs
if (requireNamespace("rcrossref")) {
  library("rcrossref")
  out <- cr_members(2258, filter=c(has_full_text = TRUE), works = TRUE)
  # (links <- crm_links(out$data$DOI[10], "all"))
  # crm_text(links, type = "pdf", cache=FALSE)
  # system.time( cacheyes <- crm_text(links, type = "pdf", cache=TRUE) )
  ### second time should be faster
  # system.time( cacheyes <- crm_text(links, type = "pdf", cache=TRUE) )
  # system.time( cacheno <- crm_text(links, type = "pdf", cache=FALSE) )
  # identical(cacheyes, cacheno)
}

## elsevier
## requires authentication
### load some Elsevier DOIs
data(dois_elsevier)

## set key first - OR set globally - See ?auth
# Sys.setenv(CROSSREF_TDM_ELSEVIER = "your-key")
## XML
link <- crm_links("10.1016/j.funeco.2010.11.003", "xml")
# res <- crm_text(url = link, type = "xml")
## plain text
link <- crm_links("10.1016/j.funeco.2010.11.003", "plain")
# res <- crm_text(url = link, "plain")

## Wiley
## requires authentication
### load some Wiley DOIs
data(dois_wiley)

## set key first - OR set globally - See ?auth
# Sys.setenv(CROSSREF_TDM = "your-key")

### all wiley
tmp <- crm_links("10.1111/apt.13556", "all")
# crm_text(url = tmp, type = "pdf", cache=FALSE,
#   overwrite_unspecified = TRUE)

#### older dates for Wiley
# tmp <- crm_links(dois_wiley$set2[1], "all")
# crm_text(tmp, type = "pdf", cache=FALSE,
#   overwrite_unspecified=TRUE)

```

```
### Wiley paper with CC By 4.0 license
# tmp <- crm_links("10.1113/jp272944", "all")
# crm_text(tmp, type = "pdf", cache=FALSE)

## End(Not run)
```

 crm_xml

Get full text XML

Description

Get full text XML

Usage

```
crm_xml(url, overwrite_unspecified = FALSE, ...)
```

Arguments

url	A URL (character) or an object of class <code>tdmurl</code> from a call to <code>crm_links()</code> . If you'll be getting text from the publishers are use Crossref TDM (which requires authentication), we strongly recommend using <code>crm_links()</code> first and passing output of that here, as <code>crm_links()</code> grabs the publisher Crossref member ID, which we use to do authentication and other publisher specific fixes to URLs
overwrite_unspecified	(logical) Sometimes the crossref API returns mime type 'unspecified' for the full text links (for some Wiley dois for example). This parameter overrides the mime type to be type.
...	Named curl parameters passed on to <code>curl::HttpClient()</code> , see <code>curl::curl_options()</code> for available curl options

Details

Note that this function is not vectorized. To do many requests use a for/while loop or lapply family calls, or similar.

Note that some links returned will not in fact lead you to full text content as you would understandably think and expect. That is, if you use the `filter` parameter with e.g., `rcrossref::cr_works()` and filter to only full text content, some links may actually give back only metadata for an article. Elsevier is perhaps the worst offender, for one because they have a lot of entries in Crossref TDM, but most of the links that are apparently full text are not in fact full text, but only metadata.

Check out [auth](#) for details on authentication.

Examples

```
## Not run:
## peerj
x <- crm_xml("https://peerj.com/articles/2356.xml")

## pensoft
data(dois_pensoft)
(links <- crm_links(dois_pensoft[1], "all"))
### xml
crm_xml(url=links)

## hindawi
data(dois_pensoft)
(links <- crm_links(dois_pensoft[1], "all"))
### xml
crm_xml(links)

## End(Not run)
```

dois_crminer	<i>A character vector of 500 DOIs from Crossref</i>
--------------	---

Description

Obtained via `rcrossref::cr_works(filter = c(has_full_text = TRUE), limit = 500)`

Format

A character vector of length 500

dois_crminer_ccby3	<i>A character vector of 100 DOIs from Crossref with license CC-BY 3.0</i>
--------------------	--

Description

Obtained via `rcrossref::cr_works(filter = list(has_full_text = TRUE, license_url="http://creativecommons.org/licenses/by/3.0/"))`

Format

A character vector of length 100

dois_elsevier *A character vector of 100 Elsevier DOIs from Crossref*

Description

Obtained via `rcrossref::cr_members(78, filter = c(has_full_text = TRUE), works = TRUE, limit = 100)`

Format

A character vector of length 100

dois_hindawi *A character vector of 50 Hindawi DOIs from Crossref*

Description

Obtained via `rcrossref::cr_members(98, filter = c(has_full_text = TRUE), works = TRUE, limit = 50)`

Format

A character vector of length 50

dois_pensoft *A character vector of 100 Pensoft DOIs from Crossref*

Description

Obtained via `rcrossref::cr_members(2258, filter = c(has_full_text = TRUE), works = TRUE, limit = 100)`

Format

A character vector of length 100

dois_wiley *A list of 3 character vectors totaling 250 Wiley DOIs from Crossref*

Description

- set1: Obtained via `rcrossref::cr_members(311, filter = c(has_full_text = TRUE), works = TRUE, limit = 100)`
- set2 (a set with older dates): Obtained via `rcrossref::cr_members(311, filter=c(has_full_text = TRUE, type = "older"), works = TRUE, limit = 100)`
- set3 (with CC By 4.0 license): Obtained via `rcrossref::cr_members(311, filter=c(has_full_text = TRUE, license = "cc-by"), works = TRUE, limit = 100)`

Format

A list of length 3, set1 with a character vector of length 100, set2 with a character vector of length 100, and set3 with a character vector of length 50.

Index

*Topic **datasets**

- dois_crminer, 16
- dois_crminer_ccby3, 16
- dois_elsevier, 17
- dois_hindawi, 17
- dois_pensoft, 17
- dois_wiley, 17

*Topic **package**

- crminer-package, 2

as_tdmurl, 3

as_tdmurl(), 2

auth, 4, 7, 10, 12, 13, 15

auth(), 2

crm_cache, 5, 10, 13

crm_extract, 6

crm_extract(), 2

crm_html, 7

crm_links, 8

crm_links(), 2, 7, 9, 11, 12, 15

crm_pdf, 9

crm_pdf(), 2

crm_plain, 11

crm_plain(), 2

crm_text, 12

crm_text(), 2

crm_xml, 15

crm_xml(), 2

crminer (crminer-package), 2

crminer-package, 2

curl::HttpClient(), 7, 8, 10, 11, 13, 15

curl::curl_options(), 7, 10, 11, 13, 15

dois_crminer, 16

dois_crminer_ccby3, 16

dois_elsevier, 17

dois_hindawi, 17

dois_pensoft, 17

dois_wiley, 17

lapply(), 5

pdftools::pdf_info(), 6

pdftools::pdf_text(), 6

rcrossref::cr_works(), 7, 8, 10, 12, 13, 15