

R tools for accessing research literature for text mining

Scott Chamberlain^{*,a}

^a*rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA*

Abstract

Text-mining is a powerful method for answering research questions. However, getting texts to extract information can be a daunting and complicated task. The primary reason for this is publishers. There are thousands of different publishers, each with their own licenses, URL patterns, access options, and more. Layered on top of that is the varied access each user has based on their institutional affiliation. Here, I introduce a suite of software packages in the R programming language for fetching texts. The tapestry of different publishers, access levels, and other factors requires a patchwork of approaches for getting texts to users. The flagship R package called `fulltext` attempts to simplify search and retrieval of texts for text-mining. The `fulltext` package, along with many others, make acquiring texts easier than ever, facilitating answering research questions with text-mining.

*Corresponding author

Email address: `myrmecocystus(at)gmail.com` (Scott Chamberlain)

5 Introduction

6 There's more than 100 million research articles published (Crossref API: [https://github.com/CrossRef/](https://github.com/CrossRef/rest-api-doc)
7 [rest-api-doc](https://github.com/CrossRef/rest-api-doc)), representing an enormous amount of knowledge. In addition to simply reading these
8 articles, the articles contain a vast trove of information of interest to researchers for machine aided
9 questions (Kong & Gerstein, 2018; Usai et al., 2018).

10 For example, many researchers are interested in statistical outcomes of articles: questions about P-values,
11 about effect sizes, and more. With regard to effect sizes, these are of particular interest, as they are
12 often combined in meta-analyses to draw broad conclusions about a particular question.

13 Text-mining is the broad term associated with pulling information out of articles. Given the importance
14 of text-mining, good text-mining tools are needed to make it easier for researchers to do. Graphical user
15 interface (GUI) based text mining tools are available (e.g., Ba & Bossy, 2016) and some research papers
16 have used them (Chaix et al., 2019), but given the urgent recent call to action for more reproducible
17 research (Open Science Collaboration, 2015; Camerer et al., 2016, 2018), we must move away from GUI
18 based tools as fast as possible.

19 In particular, the R programming language is used widely throughout many academic fields and thus
20 tools in R for text mining are of particular importance because they can be adopted by academics
21 rapidly.

22 Here, I present an overview of text-mining tools in the R programming language, not for text-mining
23 analysis, but rather those tools for searching for, acquiring, and extracting particular chunks of texts (e.g.,
24 title, abstract, authors). Most of the packages are part of the rOpenSci suite (<https://ropensci.org/>).

25 Digital articles: technical aspects

26 Those articles that are digital (which in theory includes all articles) can be split into two groups:
27 machine readable and non-machine readable.

28 The machine readable articles are those in XML¹, JSON², or plain text format. The former two, XML
29 and JSON, are the best machine readable types because they are structured data³, whereas plain text
30 has no structure - it's simply a set of characters with line breaks and spaces in between.

¹<https://www.w3.org/TR/xml/>

²<https://tools.ietf.org/html/rfc7159>

³https://en.wikipedia.org/wiki/Data_model

31 Of the non-machine readable types, the most notable is the Portable Document Format (PDF)⁴. These
32 can be broken out into two groups: text based PDFs and scanned PDFs. The former are converted from
33 digital versions of various kinds (MS Word, OpenOffice, LaTeX, markdown, etc.), while the latter are
34 created by scanning print articles to a PDF format. Text-based PDFs are much better for text-mining
35 purposes as plain text can be extracted easily in R with [pdftools](#), a binding to [libpoppler](#). However,
36 with scanned PDFs, text must be extracted using Optimal Character Recognition (OCR; see R package
37 [tesseract](#)), which isn't always a clean solution, especially compared to true text based PDFs.

38 The reality in scholarly publishing is all publishers, if they provide any access to their articles, only
39 provide PDF format. Very few publishers, with some quite large (Elsevier, Pensoft, PLOS), provide
40 XML format. Although most publishers most likely have the XML behind each of their articles, they for
41 some indefensible reason do not share it - making text-mining more difficult. Some provide plain text
42 (Elsevier). I only know of one publisher that provides full text as JSON (PLOS). Thus, text-mining, in
43 most cases, will require extracting text from PDFs.

44 **Digital articles: the access landscape**

45 Access to full-text is the holy grail in text-mining. Some use cases can get by with article metadata
46 (authors, title, etc.), some with abstracts, but many use cases require full-text.

47 The landscape of access to full-text is extremely heterogeneous, with the majority of variation along the
48 publisher axis. The major hurdle is paywalls. The majority of articles are published by the big three
49 publishers - Wiley, Springer, Elsevier - and the majority of their articles are behind paywalls.

50 A promising sign is an increasing number of open access articles, yet open access articles represent a
51 small percent of all articles: an estimate in 2018 said that 28% of the scholarly literature was open
52 access (Piwowar et al., 2018).

53 With respect to paywalled articles, access varies by institution, depending on each institution's publisher
54 contracts. MORE ABOUT THIS ...

55 Some may not realize access to articles varies with IP address so that access from campus vs. from
56 home (if not on a VPN) will drastically differ. Sometimes a VPN is required, and this can provide a
57 significant technical hurdle to users attempting to do text-mining work.

⁴<https://en.wikipedia.org/wiki/PDF>

58 One final hurdle in text-mining comes unsurprisingly from Elsevier. They use so-called “fences” for
59 programmatic access. That is, even if a person trying to get an article programmatically their institution
60 has access to and they have access to, and they are on the correct IP address, they may still not get
61 access to an Elsevier article. Elsevier puts in place these fences and only if you contact their technical
62 team directly can you get these fences removed, and only then on a per institution basis.

63 I can not end this section without mentioning SciHub. This is a last resort option for many probably
64 (or possibly first, depending on your level of access), providing access to full text of articles that are
65 normally paywalled. No tools in this manuscript provide access to SciHub.

66 **The discovery problem**

67 A text-mining project starts with a question. From that question, researchers then attempt to acquire
68 scholarly articles for text-mining. Finding those articles is not altogether straight-forward.

69 There are many places to search for articles; a non-exhaustive list: Google Scholar, Microsoft Academic
70 Research, Scopus, ScienceDirect, Web of Science, Pubmed/Entrez, Europe PMC, Directory of Open
71 Access Journals, Open Knowledge Maps, and more. It’s probably difficult to know where the best place
72 is to search. Some of these are paywalled (e.g., Web of Science), and some are not.

73 The most important aspect about any source for article search with respect to reproducible research is
74 being able to use the data source programmatically. Of those listed above, the following can be used
75 programmatically: Microsoft Academic Research, Scopus, ScienceDirect, Pubmed/Entrez, Europe PMC,
76 and Directory of Open Access Journals. All of these are included in the R package [fulltext](#), discussed
77 further below.

78 **Data sources**

79 There is increasing open access scientific literature content available online. However, only a small
80 proportion of scientific journals provide access to their full content; whereas, most publishers provide
81 open access to their metadata only (most often through Crossref; Table 1). The following is a synopsis
82 of the major data sources and associated R tools.

83 Table 1. Sources of scientific literature, their content type provided via web services, whether rOpenSci
84 has an R packages for the service, and where to find the API documentation.

Data Provider	Content Type	rOpenSci Package	Documentation
Crossref	Metadata only	rcrossref/crminer	5
DataCite	Metadata only	rdatacite	6
Biodiversity Heritage Library	Full content/Metadata	rbhl	7
Public Library of Science (PLOS)	Full text/altmetrics	rplos	8
Scopus (Elsevier)	Full content/Metadata	fulltext	9
arXiv	Full content/Metadata	aRxiv	10
Biomed Central (via Springer)	Full content/Metadata	fulltext	11
bioRxiv	Full content/Metadata	fulltext	12
PMC/Pubmed (via Entrez)	Full content/Metadata	rentrez	13
Europe PMC	Full content/Metadata	europemc	14
Microsoft Academic Search	Metadata	fulltext/microdemic	15
Directory of Open Access Journals	Metadata	jaod	16
JSTOR Data for Research	Full content	jstor	17
ORCID	Metadata	rorcid	18
Wikimedia's Citoid	Citations	rcitoid	19
Open Citation Corpus	Citations	citecorp	20

⁵<https://api.crossref.org>

⁶<https://support.datacite.org/docs/api>

⁷<http://bit.ly/KYQ1Rd>

⁸<http://api.plos.org/solr>

⁹<http://bit.ly/J9S616>

¹⁰<https://arxiv.org/help/api/index>

¹¹<https://dev.springer.com/>

¹²<http://www.biorxiv.org/>

¹³<https://www.ncbi.nlm.nih.gov/books/NBK25500>

¹⁴<https://azure.microsoft.com/en-us/services/cognitive-services>

¹⁵<https://dev.labs.cognitive.microsoft.com/docs/services/56332331778daf02acc0a50b/operations/565d9001ca73072048922d97>

¹⁶<https://doaj.org/api/v1/docs>

¹⁷<https://www.jstor.org/dfr/>

¹⁸<https://pub.orcid.org/>

¹⁹https://en.wikipedia.org/api/rest_v1/#/Citation/getCitation

²⁰<http://opencitations.net/>

85 *Crossref/Datacite*

86 Crossref is a non-profit that creates (or “mints”) Digital Object Identifiers (DOIs). In addition, they
87 maintain metadata associated with each DOI. The metadata ranges from simple (including author, title,
88 dates, DOI, type, publisher) to including number of citations to the article, as well as references in the
89 article, and even abstracts. At the time of writing they hold 100 million DOIs.

90 One can search by DOI or search citation data to get citations. In addition, Crossref has a text-mining
91 opt-in program for publishers. The result of this is that some publishers provide URLs for full text
92 content of their articles. The majority of these links are pay-walled, while some are open access. Using
93 any of the various tools for working with Crossref data, you can filter your search to get only articles
94 with full text links, and further to get only articles with full text links that are open access.

95 The main interfaces for Crossref in R are [rcrossref](#) and [crminer](#). Similar interfaces are available in Ruby
96 ([serrano](#)) and Python ([habanero](#)).

97 Datacite is similar to Crossref, but focuses on datasets instead of articles. The main interface for
98 Datacite in R is [rdatacite](#).

99 *Biodiversity Heritage Library*

100 The Biodiversity Heritage Library (BHL) houses scans of biodiversity books, and provides web interfaces
101 and APIs to query and fetch those data. They also provide text of the scanned pages. The main R
102 interace to BHL is through [rbhl](#).

103 *Public Library of Science*

104 The Public Library of Science (PLOS) is one of the largest open access only publishers. They as of this
105 writing have published 2.1 million articles. One of the strong advantages of PLOS is that they provide
106 an API to their Solr instance, which is a very flexible way to search their articles. The main R interace
107 to PLOS is through [rplos](#).

108 *Elsevier/Scopus*

109 Elsevier is one of the largest publishers. Most of their articles are not open access. However, they have a
110 numbrer of advantages if you have access to their articles: they are one of the few publishers to provide

111 machine readable XML (many publishers do have XML versions of articles, but do not provide it); they
112 are one of the few (two) publishers part of Crossref's text and data mining program. The packages
113 [fulltext](#) and [crminer](#) can be used to access Elsevier articles through Crossref's TDM program. There's
114 an interface to Scopus article search within [fulltext](#).

115 *arXiv/bioRxiv*

116 arXiv and bioRxiv are preprint publishers, the former in existence for many years, and the latter new
117 on the scene. You can access articles from these publishers through [fulltext](#). arXiv does provide a web
118 API that we hook into; bioRxiv does not, but we can get you articles nonetheless.

119 *Pubmed/PMC/Europe PMC*

120 Pubmed/PMC is a corpus/website of NIH funded research in the United States; while Europe PMC is
121 an equivalent for the European Union. You can access articles from Pubmed/PMC through [fulltext](#),
122 and for Europe PMC through [europepmc](#).

123 *Microsoft Academic Research*

124 Microsoft Academic Research (MAR) is a search engine for research articles. You can use their GUI
125 web interface to search, and they provide APIs for programmatic access. The R interface for MAR is
126 [microdemic](#); and [fulltext](#) hooks into [microdemic](#) as well for article search and abstract retrieval.

127 *Directory of Open Access Journals*

128 Directory of Open Access Journals (DOAJ) maintains data on open access journals, as well as some
129 portion of the articles in those journals. Thus, you can search for journals as well as articles with DOAJ.
130 The R interface for DOAJ is [jaod](#).

131 *JSTOR*

132 JSTOR's Data for Research program gives institutions with access to JSTOR, access to full text of
133 articles within JSTOR. There is no way however to make the interaction with JSTOR completely
134 programmatic, thus making reproducible research very difficult. Nonetheless, there is an R package
135 ([jstor](#)) for using data from JSTOR's Data for Research.

136 *ORCID*

137 ORCID (<https://orcid.org/>) is an organization keeping track of identifiers and metadata for researchers
138 around the world. Individuals can optionally maintain metadata on their scholarly works connected to
139 their account with ORCID. Thus, across all of ORCID, a significant cache of metadata is accruing on
140 scholarly works, their funding amounts, collaborators, etc., useful for bibliometrics research and more.
141 The R interface for ORCID is [rorcid](#).

142 *Citoid/Open Citation Corpus*

143 The Open Citation Corpus (<http://opencitations.net/>) holds records of which articles cite which other
144 articles, allowing for all important research on the scholarly web of citation. Citation data has been
145 very closely guarded until recently, but the largest publishers are still not contributing to the Open
146 Citation Corpus. The R interface to the Open Citation Corpus is [rcitoid](#).

147 **fulltext: a swiss army knife for text mining in R**

148 [fulltext](#) is a general purpose R package for the data part of text-mining: search for articles, get links to
149 articles, get article abstracts, and fetch full text of articles. The **fulltext** package is always adding
150 additional data sources as time allows (See Table 1). Starting from searching for articles, the outputs of
151 search can be fed into a function to get links to those articles, or to get abstracts for those articles, or
152 to fetch their full text.

153 The following is a breakdown of the major distinct parts of **fulltext**.

154 *Search*

155 **ft_search()** provides search access to nine different data sources (PLOS, BMC, Crossref, Entrez, arXiv,
156 bioRxiv, Europe PMC, Scopus, Microsoft Academic), creating a mostly unified interface to all data
157 sources. The parts of each data source that are common are for the most part factored out into the
158 parameters of the **ft_search()** function: query term(s), pagination (number of results, result number
159 to start at). In addition, we allow the user to pass on data source specific options to refine the search
160 per data source.

161 With **ft_search()**, you can query any combination of the nine data sources at once. The returned
162 object is a list, with access to results of each data source by its name (e.g., **\$plos**, or **\$crossref**). For

163 each data source, the returned object does vary because the returned data from each data source widely
164 varies; for the most part data.frame's are returned. For those data sources not queried, their slot is
165 empty.

166 One important aspect of the research result we highlight is the licenses in the returned data for each
167 data source.

```
x <- ft_search(query = 'ecology', from = c("plos", "crossref"))
```

168 The results for this PLOS search have all CC-BY licenses

```
x$plos
#> Query: [ecology]
#> Records found, returned: [47257, 10]
#> License: [CC-BY]
#>
#>          id
#> 1 10.1371/journal.pone.0001248
#> 2 10.1371/journal.pone.0059813
#> 3 10.1371/journal.pone.0080763
#> 4 10.1371/journal.pone.0155019
#> 5 10.1371/journal.pone.0175014
#> 6 10.1371/journal.pone.0150648
#> 7 10.1371/journal.pone.0208370
#> 8 10.1371/journal.pcbi.1003594
#> 9 10.1371/journal.pone.0102437
#> 10 10.1371/journal.pone.0166559
```

169 Whereas the results for this Crossref search have mixed licenses

```
x$crossref
#> Query: [ecology]
#> Records found, returned: [164657, 10]
#> License: [variable, see individual records]
#>   archive          container.title    created deposited
```

```

#> 1 Portico Ecology 2006-05-03 2018-08-04
#> 2 Portico Ecology 2006-05-03 2018-08-04
#> 3 NA Ecology 2006-05-03 2018-08-04
#> 4 NA Ecology 2006-05-03 2018-08-04
#> 5 NA Ecology 2006-05-03 2018-08-04
#> 6 NA Ecology 2006-05-03 2018-08-04
#> 7 NA Ecology 2006-05-09 2018-08-01
#> 8 Portico Ecology 2017-04-26 2019-03-08
#> 9 NA Trends in Ecology & Evolution 2002-07-25 2017-06-14
#> 10 NA Journal of Industrial Ecology 2014-11-21 2017-06-23
#> Variables not shown: published.print (chr), published.online (chr), doi
#> (chr), indexed (chr), issn (chr), issue (chr), issued (chr), member
#> (chr), page (chr), prefix (chr), publisher (chr), reference.count
#> (chr), score (chr), source (chr), title (chr), type (chr), url (chr),
#> volume (chr), author (list), link (list), license (list), subject
#> (chr), alternative.id (chr), subtitle (chr), reference (list)

```

170 You can dig into the license field for each article, with URLs holding information on each license

```

vapply(x$crossref$data$license, function(w) w$URL[1], "")
#> [1] "http://doi.wiley.com/10.1002/tdm_license_1.1"
#> [2] "http://doi.wiley.com/10.1002/tdm_license_1.1"
#> [3] "http://doi.wiley.com/10.1002/tdm_license_1"
#> [4] "http://doi.wiley.com/10.1002/tdm_license_1.1"
#> [5] "http://doi.wiley.com/10.1002/tdm_license_1"
#> [6] "http://doi.wiley.com/10.1002/tdm_license_1"
#> [7] "http://doi.wiley.com/10.1002/tdm_license_1"
#> [8] "http://doi.wiley.com/10.1002/tdm_license_1.1"
#> [9] "http://www.elsevier.com/tdm/userlicense/1.0/"
#> [10] "http://doi.wiley.com/10.1002/tdm_license_1.1"

```

171 *Links*

172 `ft_links()` provides two pathways to get links (URLs) for articles, with a choice of four different data
173 sources (PLOS, BMC, Crossref, Entrez). First, you can use `ft_search()`, then pass the output of that
174 function to `ft_links()`.

```
out <- ft_search(query = "ecology", from = "entrez")
ft_links(out)
#> <fulltext links>
#> [Found] 6
#> [IDs] ID_30964001 ID_30962485 ID_30962432 ID_30952928 ID_30674747
#>      ID_30674743 ...
```

175 Second, you can pass DOIs directly to `ft_links()`. Both end up at the same point, links for each
176 article, if they could be found for the user selected data source.

```
# FIXME
ft_links(out$entrez$data$doi)
```

177 The biggest caveat with `ft_links()` is that we can't guarantee that the links will work. Link rot is one
178 way in which the links may not work: link rot is when the URL does not point to the original content
179 anymore, or fails altogether. Additionally, with Crossref, publishers can deposit URLs for articles, but
180 they make change the URLs at some later date but not update the URLs with Crossref.

181 *Abstracts*

182 `ft_abstract()` provides access to article abstracts from four different data sources (PLOS, Scopus,
183 Microsoft Academic Research, Crossref). The only way to use the function is to pass article identifiers,
184 which are for the most DOIs.

185 The advantage of abstracts over full text is that abstracts can often be retrieved even for paywalled
186 articles. That is, you can have much broader coverage of the articles you're targeting relative to full
187 text.

188 If you are after abstracts, and you are already getting or already have full text, and if the articles are in
189 XML format, then you can use [pubchunks](#) to extract out the abstracts.

190 *Fetch full text*

191 `ft_get()` fetches full text of articles from many different data sources. From the DOIs that are passed
192 in to the function, we detect the publisher, and there are specific plugins for certain publishers: AAAS,
193 American Institute of Physics, American Society of Clinical Oncology, American Society for Microbiology,
194 arXiv, bioRxiv, BiomedCentral, Copernicus, Crossref, Elife, Elsevier, Pubmed/PMC via NCBI's Entrez,
195 Frontiers, IEEE, Informa, Instituto de Investigaciones Filologicas, American Medical Association,
196 Microbiology Society, PeerJ, Pensoft, PLOS, PNAS, Royal Society of Chemistry, ScienceDirect, Scientific
197 Societies, and Wiley.

198 If there's no built-in plugin for the publisher already, we use the FTDOI API (<https://ftdoi.org>) to try
199 to get the link for the full text of the article. If the FTDOI API doesn't bear fruit, we search Crossref
200 for a link to the full text. If Crossref doesn't have any full text links, we give up.

201 Since users can go through a lot of article requests, we cache successfully downloaded articles, and keep
202 that knowledge consistent across R sessions; all subsequent requests for the same article just use the
203 cached version. Additionally, all errors in `ft_get()` are collected in a tidy data.frame in the output of
204 the function to help the user quickly determine what went wrong.

205 **How to text mine from R: Three case studies**

206 *Case study 1: Citation mining*

207 In this example, xxxx

208 *Load libraries*

```
library("rcrossref")
library("rplos")
library("rorcid")
library("rcitoid")
library("citecorp")
```

209 *rcrossref*

210 Using `rcrossref` for Crossref data:

```

x <- cr_works(query="NSF")
head(x$data)
#> # A tibble: 6 x 32
#>   alternative.id container.title created deposited published.print doi
#>   <chr>           <chr>           <chr>  <chr>      <chr>           <chr>
#> 1 S106352031630~ Applied and Co~ 2016-0~ 2019-02-~ 2018-03      10.1~
#> 2 <NA>           Biogeosciences~ 2017-0~ 2017-07-~ <NA>         10.5~
#> 3 <NA>           Global Biogeoc~ 2018-0~ 2019-01-~ 2018-10      10.1~
#> 4 <NA>           IEEE Communica~ 2016-1~ 2017-12-~ 2017         10.1~
#> 5 S002178241400~ Journal de Mat~ 2014-0~ 2018-10-~ 2014-10      10.1~
#> 6 123            Light: Science~ 2019-0~ 2019-01-~ 2019-12      10.1~
#> # ... with 26 more variables: indexed <chr>, issn <chr>, issue <chr>,
#> #   issued <chr>, member <chr>, page <chr>, prefix <chr>, publisher <chr>,
#> #   reference.count <chr>, score <chr>, ...

```

211 Case study 2: Abstract mining

212 Sometimes you just need abstracts for your research question. The benefit of only needing abstracts,
 213 and not need full text, is that there's many more articles that will have abstracts available than have
 214 their full text available.

215 As an example, let's say you xxxx

```
library("fulltext")
```

216 *xxxxx*

217 Using fulltext:

```

res <- ft_search("ecology", from = "crossref",
  crossrefopts = list(filter = c(has_abstract = TRUE)))
ids <- res$crossref$data$doi
out <- ft_abstract(x = ids, from = "crossref")
abstracts <- vapply(out$crossref, "[[", "", "abstract")

```

218 Using `quanteda`, read the abstracts into a corpus

```
library("quanteda")
corp <- corpus(abstracts)
docvars(corp) <- ids
```

219 Get a summary of the abstracts

```
summary(corp)
#> Corpus consisting of 10 documents:
#>
#>   Text Types Tokens Sentences          V1
#> text1    143    262         10 10.2458/v22i1.21112
#> text2    117    244          6 10.2458/v17i1.21696
#> text3     75    118          4 10.2458/v25i1.23119
#> text4      5      8          1 10.2458/v1i1.21154
#> text5    105    171          7 10.1155/2011/868426
#> text6    112    181          6 10.1155/2012/273413
#> text7    117    240          8 10.5194/we-13-91-2013
#> text8    140    245          9 10.5194/we-13-95-2013
#> text9    107    202          7 10.1155/2014/198707
#> text10   118    224          6 10.5402/2011/897578
#>
#> Source: /Users/sckott/github/ropensci/textmine/use-cases/* on x86_64 by sckott
#> Created: Thu Apr 11 11:20:19 2019
#> Notes:
```

220 Use the `kwic()` function to see a word in context across the abstracts

```
kwic(corp, pattern = "ecology")
#>
#> [text1, 33] knowledge production within critical political / ecology /
#> [text1, 50] in scientific articles on dryland / ecology /
```

#> [text1, 204] to equilibrium models in range / ecology /

#> [text1, 246] communal areas.Keywords: Critical political / ecology /

#> [text1, 255] , scientific models, rangeland / ecology /

#> [text2, 5] < jats:p> Political / ecology /

#> [text2, 23] manifestations of political economy and / ecology /

#> [text2, 45] I try to extend political / ecology /

#> [text2, 149] , in dialogue with political / ecology /

#> [text2, 177] people and resources that political / ecology /

#> [text2, 229] indigeneity scholars.Key words: political / ecology /

#> [text3, 71] an analysis from a political / ecology /

#> [text3, 114] system, supermarkets, political / ecology /

#> [text6, 134] was observed when allopatry and / ecology /

#> [text7, 167] ecosystem should be considered for / ecology /

#> [text7, 185] the" four-color issue of / ecology /

#> [text7, 201] step toward advancing knowledge in / ecology /

#> [text9, 195] or for theoretical studies integrating / ecology /

#>

#> . This article is a

#> , and investigates the functions

#> , and the fence-line photographs

#> , fence-line photography, scientific

#> , Southern Africa</

#> has expanded in multiple new

#> in the" problem"

#> to engage with ethnic studies

#> approaches to better understand the

#> focuses on cannot be adequately

#> , coloniality, Maidu,

#> standpoint allows a different interpretation

#> </ jats:p>

#> act together, leading to

```
#> "? Here, I
#> ", and propose that
#> and conservation biology. In
#> and biogeography.</
```

221 *Case study 3: Full text mining*

222 In this example, xxxx

```
library("fulltext")
# library("crminer")
```

223 *Search for articles*

224 Search for the term *ecology* in PLOS journals.

```
(res1 <- ft_search(query = 'ecology', from = 'plos'))
#> Query:
#> [ecology]
#> Found:
#> [PLOS: 47272; BMC: 0; Crossref: 0; Entrez: 0; arxiv: 0; biorxiv: 0; Europe PMC: 0; Scopus:
#> Returned:
#> [PLOS: 10; BMC: 0; Crossref: 0; Entrez: 0; arxiv: 0; biorxiv: 0; Europe PMC: 0; Scopus: 0; .
```

225 Each publisher/search-engine has a slot with metadata and data

```
res1$plos
#> Query: [ecology]
#> Records found, returned: [47272, 10]
#> License: [CC-BY]
#>
#> id
#> 1 10.1371/journal.pone.0001248
#> 2 10.1371/journal.pone.0059813
#> 3 10.1371/journal.pone.0080763
```



```
#> 4 10.1371/journal.pone.0155019
#> 5 10.1371/journal.pone.0175014
#> 6 10.1371/journal.pone.0150648
#> 7 10.1371/journal.pone.0208370
#> 8 10.1371/journal.pcbi.1003594
#> 9 10.1371/journal.pone.0102437
#> 10 10.1371/journal.pone.0166559
```

226 *Get full text*

227 Using the results from `ft_search()` we can grab full text of some articles

```
(out <- ft_get(res1))
#> <fulltext text>
#> [Docs] 10
#> [Source] ext - /Users/sckott/Library/Caches/R/fulltext
#> [IDs] 10.1371/journal.pone.0001248 10.1371/journal.pone.0059813
#>      10.1371/journal.pone.0080763 10.1371/journal.pone.0155019
#>      10.1371/journal.pone.0175014 10.1371/journal.pone.0150648
#>      10.1371/journal.pone.0208370 10.1371/journal.pcbi.1003594
#>      10.1371/journal.pone.0102437 10.1371/journal.pone.0166559 ...
```

228 *Extract text from pdfs*

229 Ideally for text mining you have access to XML or other text based formats. However, sometimes you
 230 only have access to PDFs. In this case you want to extract text from PDFs. `fulltext` can help with
 231 that.

232 You can extract from any pdf from a file path, like:

```
path <- system.file("examples", "example1.pdf", package = "fulltext")
ft_extract(path)
#> <document>/Library/Frameworks/R.framework/Versions/3.5/Resources/library/fulltext/examples/ex
#> Title: Suffering and mental health among older people living in nursing homes---a mixed-met
```

```
#> Producer: pdfTeX-1.40.10
#> Creation date: 2015-07-17
```

233 *Extract text chunks*

234 Requires the [pubchunks](#) library. Here, we'll search for some PLOS articles, then get their full text, then
235 extract various parts of each article with `pub_chunks()`.

```
library("pubchunks")
res <- ft_search(query = "ecology", from = "plos", limit = 3)
x <- ft_get(res)
x %>% ft_collect() %>% pub_chunks(c("doi", "history")) %>% pub_tabularize()

#> $plos
#> $plos$`10.1371/journal.pone.0001248`
#>
#> doi history.received history.accepted
#> 1 10.1371/journal.pone.0001248 2007-07-02 2007-11-06
#> .publisher
#> 1 plos
#>
#> $plos$`10.1371/journal.pone.0059813`
#>
#> doi history.received history.accepted
#> 1 10.1371/journal.pone.0059813 2012-09-16 2013-02-19
#> .publisher
#> 1 plos
#>
#> $plos$`10.1371/journal.pone.0080763`
#>
#> doi history.received history.accepted
#> 1 10.1371/journal.pone.0080763 2013-08-15 2013-10-16
#> .publisher
#> 1 plos
```

236 **Future directions**

237 To make text mining easier:

1. publishers should provide XML if they have it
2. publishers should not change URL patterns so often, or at all
3. publishers should keep their Crossref metadata up to date
4. remove all paywalls (Easy, yes?)
5. ...

Acknowledgments

XXXX

Data Accessibility

All scripts and data used in this paper can be found in the permanent data archive Zenodo under the digital object identifier (DOI). This DOI corresponds to a snapshot of the GitHub repository at <https://github.com/ropensci/textmine>. Software can be found at <https://github.com/ropensci/xxx>, xxxx, all under MIT licenses.

References

- Ba M., Bossy R. 2016. Interoperability of corpus processing work-flow engines: The case of alvisnlp/ml in openminted. In: *Proceedings of the workshop on cross-platform text mining and natural language processing interoperability (interop 2016) at IREC*. 15–18.
- Camerer CF., Dreber A., Forsell E., Ho T-H., Huber J., Johannesson M., Kirchler M., Almenberg J., Altmejd A., Chan T., Heikensten E., Holzmeister F., Imai T., Isaksson S., Nave G., Pfeiffer T., Razen M., Wu H. 2016. Evaluating replicability of laboratory experiments in economics. *Science* 351:1433–1436.
- Camerer CF., Dreber A., Holzmeister F., Ho T-H., Huber J., Johannesson M., Kirchler M., Nave G., Nosek BA., Pfeiffer T., Altmejd A., Buttrick N., Chan T., Chen Y., Forsell E., Gampa A., Heikensten E., Hummer L., Imai T., Isaksson S., Manfredi D., Rose J., Wagenmakers E-J., Wu H. 2018. Evaluating the replicability of social science experiments in nature and science between 2010 and 2015. *Nature Human Behaviour* 2:637–644.

263 Chaix E., Deléger L., Bossy R., Nédellec C. 2019. Text mining tools for extracting information about
264 microbial biodiversity in food. *Food Microbiology* 81:63–75.

265 Kong X., Gerstein MB. 2018. Text mining systems biology: Turning the microscope back on the
266 observer. *Current Opinion in Systems Biology* 11:117–122.

267 Open Science Collaboration. 2015. Estimating the reproducibility of psychological science. *Science*
268 349:aac4716–aac4716.

269 Piwowar H., Priem J., Larivière V., Alperin JP., Matthias L., Norlander B., Farley A., West J., Haustein
270 S. 2018. The state of OA: A large-scale analysis of the prevalence and impact of open access articles.
271 *PeerJ* 6:e4375.

272 Usai A., Pironti M., Mital M., Mejri CA. 2018. Knowledge discovery out of text data: A systematic
273 review via text mining. *Journal of Knowledge Management* 22:1471–1488.