

# rOpenSci tools for accessing science literature for textmining

Scott Chamberlain<sup>\*,a</sup>

<sup>a</sup>*rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA*

## Abstract

Corresponding Author:

Scott Chamberlain

rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA

Email address: [myrmecocystus@gmail.com](mailto:myrmecocystus@gmail.com)

---

\*Corresponding author

Email address: [myrmecocystus\(at\)gmail.com](mailto:myrmecocystus(at)gmail.com) (Scott Chamberlain)

9 Background. xxxx.

10 Methods. xxxx.

11 Results. xxxx.

Discussion. xxxx.

## 12 Introduction

13 There's more than 100 million articles published (source: Crossref API), representing an enormous  
14 amount of knowledge. In addition to simply reading these articles, they contain a vast trove of  
15 information of interest to researchers for machine aided questions.

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

19 Text-mining is the broad term associated with pulling information out of articles. Given the importance  
20 of text-mining, good text-mining tools are needed to make it easier for researchers. In particular, the R  
21 programming language is used widely throughout many academic fields and thus tools in R for text  
22 mining are of particular importance.

23 Here, we present an overview of text-mining tools in the R programming language. We do not cover  
24 analysis tools per se, but rather those tools for searching for, acquiring, and “mashing up” text.

## 25 Digital articles: technical aspects

26 Those articles that are digital can be split into two groups: easily machine readable and non-machine  
27 readable.

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

31 Of the non-machine readable kind, there's PDFs. These can be broken out into two groups: text based  
32 PDFs and scanned PDFs. The former are converted from digital versions of various kinds (MS Word,  
33 OpenOffice, markdown, etc.), while the latter are PDFs created by scanning in print articles for which  
34 there is no digital version.

## 35 Digital articles: the access landscape

36 Acces to full-text is the holy grail in text-mining. Some use cases can get by with article metadata  
37 (authors, title, etc.), some with abstracts, but many use cases need full-text.

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

41 A promising sign is an increasing number of open access publishers, yet these represent a very small  
42 portion of the total articles (XXXXX) (ref.).

43 With respect to paywalled articles, access varies by institution, depending on what each institution  
44 decided to pay for. In addition, some users may not realize access varies with IP address so that access  
45 from campus vs. from home (if not on a VPN) will drastically differ.

46 We can not end this section without mentioning SciHub. This is a last resort option for many probably,  
47 providing access to full text of articles that are normally paywalled. No tools in this manuscript provide  
48 access to SciHub.

## 49 **The discovery problem**

50 XXX

51

52 XXX

## 53 **Data sources**

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

58 Table 1. Sources of scientific literature, their content type provided via web services, whether rOpenSci  
59 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	<a href="#">1</a>
DataCite	Metadata only	rdatacite	<a href="#">2</a>
Biodiversity Heritage Library	Full content/Metadata	rbhl	<a href="#">3</a>
Public Library of Science (PLOS)	Full text/altmetrics	rplos	<a href="#">4</a>
Scopus (Elsevier)	Full content/Metadata	fulltext	<a href="#">5</a>
arXiv	Full content/Metadata	aRxiv	<a href="#">6</a>
Biomed Central (via Springer)	Full content/Metadata	fulltext	<a href="#">7</a>
bioRxiv	Full content/Metadata	fulltext	<a href="#">8</a>
PMC/Pubmed (via Entrez)	Full content/Metadata	rentrez	<a href="#">9</a>
Europe PMC	Full content/Metadata	europepmc	<a href="#">10</a>
Microsoft Academic Search	Metadata	fulltext/microdemic	<a href="#">11</a>
Directory of Open Access Journals	Metadata	jaod	<a href="#">12</a>
JSTOR Data for Research	Full content	jstor	<a href="#">13</a>
ORCID	Metadata	rorcid	<a href="#">14</a>
Wikimedia's Citoid	Citations	rcitoid	<a href="#">15</a>
Open Citation Corpus	Citations	citecorp	<a href="#">16</a>

<sup>1</sup><https://api.crossref.org>

<sup>2</sup><https://support.datacite.org/docs/api>

<sup>3</sup><http://bit.ly/KYQ1Rd>

<sup>4</sup><http://api.plos.org/solr>

<sup>5</sup><http://bit.ly/J9S616>

<sup>6</sup><https://arxiv.org/help/api/index>

<sup>7</sup><https://dev.springer.com/>

<sup>8</sup><http://www.biorxiv.org/>

<sup>9</sup><https://www.ncbi.nlm.nih.gov/books/NBK25500>

<sup>10</sup><https://azure.microsoft.com/en-us/services/cognitive-services>

<sup>11</sup><https://dev.labs.cognitive.microsoft.com/docs/services/56332331778daf02acc0a50b/operations/565d9001ca73072048922d97>

<sup>12</sup><https://doaj.org/api/v1/docs>

<sup>13</sup><https://www.jstor.org/dfr/>

<sup>14</sup><https://pub.orcid.org/>

<sup>15</sup>[https://en.wikipedia.org/api/rest\\_v1/#/Citation/getCitation](https://en.wikipedia.org/api/rest_v1/#/Citation/getCitation)

<sup>16</sup><http://opencitations.net/>

## 60 *Crossref/Datacite*

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

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

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

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

## 74 *Biodiversity Heritage Library*

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

## 78 *Public Library of Science*

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

## 83 *Elsevier/Scopus*

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

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

#### 90 *arXiv/bioRxiv*

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

#### 94 *Pubmed/PMC/Europe PMC*

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

#### 98 *Microsoft Academic Research*

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

#### 102 *Directory of Open Access Journals*

103 XXXXX

#### 104 *JSTOR*

105 XXXXX

#### 106 *ORCID*

107 XXXXX

#### 108 *Citoid/Open Citation Corpus*

109 XXX

## 110 How to text mine from R: Three case studies

### 111 *Case study 1: Citation mining*

112 In this example, xxxx

#### 113 *Load libraries*

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

#### 114 *rcrossref*

115 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>, ...
```

### 116 *Case study 2: Abstract mining*

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



119 their full text available.

120 As an example, let's say you xxxx

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

121 *xxxxx*

122 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")
```

123 Using [quanteda](#), read the abstracts into a corpus

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

124 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: Fri Apr 5 11:36:04 2019
#> Notes:

```

125 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.</

```

126 *Case study 3: Full text mining*

127 In this example, xxxx

```

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

```

128 *Search for articles*

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

```

(res1 <- ft_search(query = 'ecology', from = 'plos'))
#> Query:
#> [ecology]
#> Found:
#> [PLoS: 47337; 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; ...]
```

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

```
res1$plos
```

```
#> Query: [ecology]
```

```
#> Records found, returned: [47337, 10]
```

```
#> License: [CC-BY]
```

```
#> id
```

```
#> 1 10.1371/journal.pone.0001248
```

```
#> 2 10.1371/journal.pone.0059813
```

```
#> 3 10.1371/journal.pone.0155019
```

```
#> 4 10.1371/journal.pone.0080763
```

```
#> 5 10.1371/journal.pone.0208370
```

```
#> 6 10.1371/journal.pone.0150648
```

```
#> 7 10.1371/journal.pcbi.1003594
```

```
#> 8 10.1371/journal.pone.0102437
```

```
#> 9 10.1371/journal.pone.0175014
```

```
#> 10 10.1371/journal.pone.0166559
```

131 *Get full text*

132 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.0155019 10.1371/journal.pone.0080763
```

```
#> 10.1371/journal.pone.0208370 10.1371/journal.pone.0150648
```

```
#> 10.1371/journal.pcbi.1003594 10.1371/journal.pone.0102437
```

```
#> 10.1371/journal.pone.0175014 10.1371/journal.pone.0166559 ...
```

133 *Extract text from pdfs*

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

137 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
```

138 *Extract text chunks*

139 Requires the `pubchunks` library. Here, we'll search for some PLOS articles, then get their full text, then  
140 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.0155019`
#>                                doi history.received history.accepted
#> 1 10.1371/journal.pone.0155019      2015-09-22      2016-04-22
#> .publisher
#> 1      plos
```

## 141 Future directions

142 XXXX

## 143 Acknowledgments

144 XXXX

## 145 Data Accessibility

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

## 150 References