

quanteda: An R package for the quantitative analysis of textual data

Kenneth Benoit¹, Kohei Watanabe¹, Haiyan Wang², Paul Nulty³, Adam Obeng⁴, Stefan Müller⁵, and Akitaka Matsuo¹

 Department of Methodology, London School of Economics and Political Science 2 De Beers Inc.
 Centre for Research in Arts, Social Science and Humanities, University of Cambridge 4 Facebook, Inc. (work conducted at the Department of Methodology, London School of Economics and Political Science) 5 Department of Political Science, Trinity College Dublin

DOI: 10.21105/joss.00774

Software

- Review I
- Repository ¹
- Archive C²

Submitted: 30 May 2018 Published: 06 October 2018

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License (CC-BY).

Summary

quanteda is an R package providing a comprehensive workflow and toolkit for natural language processing tasks such as corpus management, tokenization, analysis, and visualization. It has extensive functions for applying dictionary analysis, exploring texts using keywords-in-context, computing document and feature similarities, and discovering multi-word expressions through collocation scoring. Based entirely on sparse operations, it provides highly efficient methods for compiling document-feature matrices and for manipulating these or using them in further quantitative analysis. Using C++ and multi-threading extensively, **quanteda** is also considerably faster and more efficient than other R and Python packages in processing large textual data.

The package is designed for R users needing to apply natural language processing to texts, from documents to final analysis. Its capabilities match or exceed those provided in many end-user software applications, many of which are expensive and not open source. The package is therefore of great benefit to researchers, students, and other analysts with fewer financial resources. While using **quanteda** requires R programming knowledge, its API is designed to enable powerful, efficient analysis with a minimum of steps. By emphasizing consistent design, furthermore, **quanteda** lowers the barriers to learning and using NLP and quantitative text analysis even for proficient R programmers.

Corpus management

quanteda makes it easy to manage texts in the form of a "corpus", which is defined as a collection of texts that includes document-level variables specific to each text, as well as meta-data for documents and for the collection as a whole. With the package, users can easily segment texts by words, paragraphs, sentences, or even user-supplied delimiters and tags, group them into larger documents by document-level variables, or subset them based on logical conditions or combinations of document-level variables.

Natural language processing

quanteda is principally designed to allow users a fast and convenient method to construct a document-feature matrix from a corpus with an ability to perform the most common natural language processing tasks such as tokenizing, stemming, forming n-grams, and



selecting and weighting features. With these functions, users can easily remove stop words and stem words in numerous languages, select words in a dictionary, and convert frequency counts into weights, for instance using tf-idf scoring.

Using the ICU library in the **stringi** package (Gagolewski, 2018) for text processing, **quanteda** correctly handles Unicode character sets for regular expression matching and detecting word boundaries for tokenization. Once texts are tokenized, **quanteda** maps tokens to a hash table of integers to increase processing speed while reducing memory usage. Many of the text processing functions are parallelized using the Intel TBB library via the **RcppParallel** package (Allaire et al., 2018).

Models and textual statistics

quanteda is especially suited to research because it was designed from the outset for the social scientific analysis of textual data. Its textmodel_* functions provide native, highly efficient implementations of several text analytic scaling methods, such as Wordscores (Laver, Benoit, & Garry, 2003), Wordfish (Slapin & Proksch, 2008), class affinity scaling (Perry & Benoit, 2017), and correspondence analysis (Greenacre, 1984). More general textmodel functions include efficient implementations of a multinomial Naive Bayes classifier designed specifically for textual data (Manning, Raghavan, & Schütze, 2008) and latent semantic analysis (Deerwester, Dumais, Furnas, Landauer, & Harshman, 1990). quanteda also works flexibly and efficiently with dictionaries, and is distributed with the 2015 version of the Lexicoder Sentiment Dictionary (Young & Soroka, 2012).

In addition to models, the package provides a variety of text statistics, such as frequency analysis, "keyness", lexical diversity, readability, and similarity and distance of documents or features. These make use of the sparsity of document-feature matrices – often over 90% sparse – and parallelism for efficient, fast computation. **quanteda** also provides methods for statistically scoring collocations, useful in identifying multi-word expressions.

Text visualization

The package provides extensive methods for visualizing textual analyses, via its family of textplot_* functions. These are typically designed to take another package object as an input, to produce a specific form of plot. For instance, from a feature co-occurrence matrix, or fcm, we can directly plot a network using textplot_network():

```
library("quanteda")
```

```
# construct the feature co-occurrence matrix
examplefcm <-
    tokens(data_corpus_irishbudget2010, remove_punct = TRUE) %>%
    tokens_tolower() %>%
    tokens_remove(stopwords("english"), padding = FALSE) %>%
    fcm(context = "window", window = 5, tri = FALSE)
# choose 30 most frequency features
topfeats <- names(topfeatures(examplefcm, 30))
# select the top 30 features only, plot the network
set.seed(100)
textplot_network(fcm_select(examplefcm, topfeats), min_freq = 0.8)</pre>
```



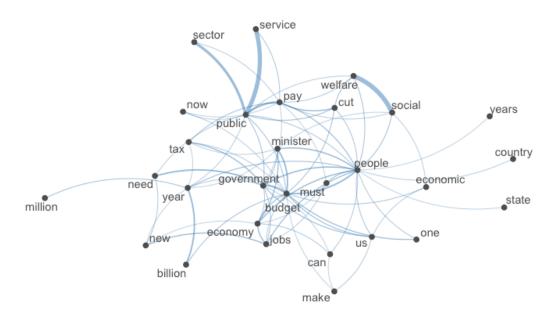


Figure 1: Feature co-occurrence network plot example.

Package design

quanteda has been carefully designed with several key aims in mind.

Consistency. quanteda functions and objects are named systematically such that corpus(), tokens() and dfm() construct those object types, and that corpus_*(), tokens_*() and dfm_*() functions return a modified version of these objects. Naming consistency applies also to the extensive built-in data objects in the package, whose names always start with data_* followed by object types. This not only gives the users a clear overview of the package, but also makes the package more reliable for other packages that depend on it.

Accessibility. **quanteda** contains extensive manual pages structured around the naming rules. Furthermore, there are references, package vignettes, examples, and tutorials on the website at https://quanteda.io. These materials help beginner users understand how to use these functions for basic operations and expert users how to combine the functions for advanced text processing and analysis.

Performance. **quanteda**'s performance is enhanced by token hashing and parallel computation implemented in C++, permitting large and fast text analysis even on computers with relatively limited resources (such as laptop computers). Built to use sparse data structures, **quanteda** can efficiently performs complex textual data analyses, such as computing distances, calculating feature discrimination statistics (keyness), or model fitting, even for large document-feature matrices.

Transparency and reproducibility. **quanteda** is designed to facilitate rigorous, transparent, and reproducible scientific analysis of text. Being open-source software, its source code can be scrutinized and corrected by other experts. Its functions are designed to encourage a reproducible workflow by linking successive processing tasks in a clear, readable manner.

Compatibility with other packages. For analysis not provided by built-in functions, users can move **quanteda** objects seamlessly to other packages, such as the **stm** package for structural topic models (Roberts, Stewart, & Tingley, 2018) or word embedding packages like **text2vec** (Selivanov & Wang, 2018). **quanteda** also works well with companion packages such as **spacyr** (Benoit & Matsuo, 2018), an R wrapper to the spaCy Python



library (Honnibal & Montani, 2017), and **readtext** (Benoit & Obeng, 2018), a package for converting and importing text files into R.

Funding and Support

Created at the London School of Economics in 2013 with funding from the European Research Council (ERC-2011-StG 283794-QUANTESS), **quanteda** is now supported by the Quanteda Initiative, a non-profit organization founded in 2018 to provide ongoing support for the "quanteda ecosystem" of open-source text analysis software.

References

Allaire, J., Francois, R., Ushey, K., Vandenbrouck, G., Geelnard, M., & Intel. (2018). *RcppParallel: Parallel programming tools for "rcpp"*. Retrieved from https://CRAN. R-project.org/package=RcppParallel

Benoit, K., & Matsuo, A. (2018). *Spacyr: Wrapper to the spaCy NLP library*. Retrieved from https://CRAN.R-project.org/package=spacyr

Benoit, K., & Obeng, A. (2018). *Readtext: Import and handling for plain and formatted text files.* Retrieved from https://CRAN.R-project.org/package=readtext

Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6), 391–407. doi:10.1002/(SICI)1097-4571(199009)41:6%3C391::AID-ASI1%3E3.0.CO;2-9

Gagolewski, M. (2018). *R package stringi: Character string processing facilities*. doi:10.5281/zenodo.1292492

Greenacre, M. J. (1984). Theory and applications of correspondence analysis. London: Academic Press.

Honnibal, M., & Montani, I. (2017). SpaCy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. doi:10.5281/zenodo.1212304

Laver, M., Benoit, K., & Garry, J. (2003). Extracting policy positions from political texts using words as data. *American Political Science Review*, 97(2), 311–331. doi:10.1017/S0003055403000698

Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to information retrieval*. Cambridge: Cambridge University Press. doi:10.1017/CBO9780511809071

Perry, P. O., & Benoit, K. (2017). Scaling Text with the Class Affinity Model. ArXiv e-prints. Retrieved from https://arxiv.org/abs/1710.08963

Roberts, M. E., Stewart, B. M., & Tingley, D. (2018). Stm: R package for structural topic models. Retrieved from http://www.structuraltopicmodel.com

Selivanov, D., & Wang, Q. (2018). *Text2vec: Modern text mining framework for r.* Re-trieved from https://CRAN.R-project.org/package=text2vec

Slapin, J. B., & Proksch, S.-O. (2008). A scaling model for estimating time-series party positions from texts. *American Journal of Political Science*, 52(3), 705–722. doi:10.1111/j.1540-5907.2008.00338.x

Young, L., & Soroka, S. (2012). Affective news: The automated coding of sentiment in political texts. *Political Communication*, 29(2), 205–231. doi:10.1080/10584609.2012.671234