

Hands-On Data Science with R

Text Mining

Graham.Williams@togaware.com

5th November 2014

Visit <http://HandsOnDataScience.com/> for more Chapters.

Text Mining or Text Analytics applies analytic tools to learn from collections of text documents like books, newspapers, emails, etc. The goal is similar to humans learning by reading books. Using automated algorithms we can learn from massive amounts of text, much more than a human can. The material could be consist of millions of newspaper articles to perhaps summarise the main themes and to identify those that are of most interest to particular people.

The required packages for this module include:

```
library(tm)           # Framework for text mining.
library(SnowballC)     # Provides wordStem() for stemming.
library(qdap)          # Quantitative discourse analysis of transcripts.
library(qdapDictionaries)
library(dplyr)         # Data preparation and pipes %>%.
library(RColorBrewer)  # Generate palette of colours for plots.
library(ggplot2)       # Plot word frequencies.
library(scales)        # Include commas in numbers.
library(Rgraphviz)     # Correlation plots.
```

As we work through this chapter, new R commands will be introduced. Be sure to review the command's documentation and understand what the command does. You can ask for help using the ? command as in:

```
?read.csv
```

We can obtain documentation on a particular package using the *help=* option of `library()`:

```
library(help=rattle)
```

This chapter is intended to be hands on. To learn effectively, you are encouraged to have R running (e.g., RStudio) and to run all the commands as they appear here. Check that you get the same output, and you understand the output. Try some variations. Explore.

Copyright © 2013-2014 Graham Williams. You can freely copy, distribute, or adapt this material, as long as the attribution is retained and derivative work is provided under the same license.



1 Loading a Corpus

A **corpus** is a collection of texts, usually stored electronically, and from which we perform our analysis. A corpus might be a collection of news articles from Reuters or the published works of Shakespeare. Within each corpus we will have separate articles, stories, volumes, each treated as a separate entity or record.

Documents which we wish to analyse come in many different formats. Quite a few formats are supported by **tm** (Feinerer and Hornik, 2014), the package we will illustrate text mining with in this module. The supported formats include text, PDF, Microsoft Word, and XML.

A number of open source tools are also available to convert most document formats to text files. For our corpus used initially in this module, a collection of PDF documents were converted to text using **pdftotext** from the **xpdf** application which is available for GNU/Linux and MS/Windows and others. On GNU/Linux we can convert a folder of PDF documents to text with:

```
system("for f in *.pdf; do pdftotext -enc ASCII7 -nopgbrk $f; done")
```

The **-enc ASCII7** ensures the text is converted to ASCII since otherwise we may end up with binary characters in our text documents.

We can also convert Word documents to text using **antitword**, which is another application available for GNU/Linux.

```
system("for f in *.doc; do antiword $f; done")
```

1.1 Corpus Sources and Readers

There are a variety of sources supported by `tm`. We can use `getSources()` to list them.

```
getSources()
## [1] "DataframeSource" "DirSource"      "URISource"      "VectorSource"
## [5] "XMLSource"
```

In addition to different kinds of sources of documents, our documents for text analysis will come in many different formats. A variety are supported by `tm`:

```
getReaders()
## [1] "readDOC"          "readPDF"
## [3] "readPlain"        "readRCV1"
## [5] "readRCV1asPlain"  "readReut21578XML"
## [7] "readReut21578XMLasPlain" "readTabular"
## [9] "readXML"
```

1.2 Text Documents

We load a sample corpus of text documents. Our corpus consists of a collection of research papers all stored in the folder we identify below. To work along with us in this module, you can create your own folder called `corpus/txt` and place into that folder a collection of text documents. It does not need to be as many as we use here but a reasonable number makes it more interesting.

```
cname <- file.path(".", "corpus", "txt")
cname
## [1] "./corpus/txt"
```

We can list some of the file names.

```
length(dir(cname))
## [1] 46
dir(cname)
## [1] "acnn96.txt"
## [2] "adm02.txt"
## [3] "ai02.txt"
## [4] "ai03.txt"
....
```

There are 46 documents in this particular corpus.

After loading the `tm` (Feinerer and Hornik, 2014) package into the R library we are ready to load the files from the directory as the source of the files making up the corpus, using `DirSource()`. The source object is passed on to `Corpus()` which loads the documents. We save the resulting collection of documents in memory, stored in a variable called `docs`.

```
library(tm)
docs <- Corpus(DirSource(cname))
docs
## <<VCorpus (documents: 46, metadata (corpus/indexed): 0/0)>>
class(docs)
## [1] "VCorpus" "Corpus"
class(docs[[1]])
## [1] "PlainTextDocument" "TextDocument"
summary(docs)
##
##          Length Class          Mode
## acnn96.txt      2 PlainTextDocument list
## adm02.txt       2 PlainTextDocument list
## ai02.txt        2 PlainTextDocument list
....
```

1.3 PDF Documents

If instead of text documents we have a corpus of PDF documents then we can use the `readPDF()` reader function to convert PDF into text and have that loaded as our Corpus.

```
docs <- Corpus(DirSource(cname), readerControl=list(reader=readPDF))
```

This will use, by default, the `pdftotext` command from `xpdf` to convert the PDF into text format. The `xpdf` application needs to be installed for `readPDF()` to work.

DRAFT

1.4 Word Documents

A simple open source tool to convert Microsoft Word documents into text is **antiword**. The separate **antiword** application needs to be installed, but once it is available it is used by **tm** to convert Word documents into text for loading into R.

To load a corpus of Word documents we use the **readDOC()** reader function:

```
docs <- Corpus(DirSource(cname), readerControl=list(reader=readDOC))
```

Once we have loaded our corpus the remainder of the processing of the corpus within R is then as follows.

The **antiword** program takes some useful command line arguments. We can pass these through to the program from **readDOC()** by specifying them as the character string argument:

```
docs <- Corpus(DirSource(cname), readerControl=list(reader=readDOC("-r -s")))
```

Here, **-r** requests that removed text be included in the output, and **-s** requests that text hidden by Word be included.

2 Exploring the Corpus

We can (and should) inspect the documents using `inspect()`. This will assure us that data has been loaded properly and as we expect.

```
inspect(docs[16])  
## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>  
##  
## [[1]]  
## <<PlainTextDocument (metadata: 7)>>  
## Hybrid weighted random forests for  
## classifying very high-dimensional data  
## Baoxun Xu1 , Joshua Zhexue Huang2 , Graham Williams2 and  
## Yunming Ye1  
## 1  
##  
## Department of Computer Science, Harbin Institute of Technology Shenzhen Gr...  
## School, Shenzhen 518055, China  
## 2  
## Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, S...  
## 518055, China  
## Email: amusing002@gmail.com  
## Random forests are a popular classification method based on an ensemble of a  
## single type of decision trees from subspaces of data. In the literature, t...  
## are many different types of decision tree algorithms, including C4.5, CART...  
## CHAID. Each type of decision tree algorithm may capture different information  
## and structure. This paper proposes a hybrid weighted random forest algorithm,  
## simultaneously using a feature weighting method and a hybrid forest method to  
## classify very high dimensional data. The hybrid weighted random forest alg...  
## can effectively reduce subspace size and improve classification performance  
## without increasing the error bound. We conduct a series of experiments on ...  
## high dimensional datasets to compare our method with traditional random fo...  
....
```

3 Preparing the Corpus

We generally need to perform some pre-processing of the text data to prepare for the text analysis. Example transformations include converting the text to lower case, removing numbers and punctuation, removing stop words, stemming and identifying synonyms. The basic transforms are all available within `tm`.

```
getTransformations()  
## [1] "removeNumbers"      "removePunctuation" "removeWords"  
## [4] "stemDocument"       "stripWhitespace"
```

The function `tm_map()` is used to apply one of these transformations across all documents within a corpus. Other transformations can be implemented using R functions and wrapped within `content_transformer()` to create a function that can be passed through to `tm_map()`. We will see an example of that in the next section.

In the following sections we will apply each of the transformations, one-by-one, to remove unwanted characters from the text.

3.1 Simple Transforms

We start with some manual special transforms we may want to do. For example, we might want to replace “/”, used sometimes to separate alternative words, with a space. This will avoid the two words being run into one string of characters through the transformations. We might also replace “@” and “|” with a space, for the same reason.

To create a custom transformation we make use of `content_transformer()` crate a function to achieve the transformation, and then apply it to the corpus using `tm_map()`.

```
toSpace <- content_transformer(function(x, pattern) gsub(pattern, " ", x))
docs <- tm_map(docs, toSpace, "/")
docs <- tm_map(docs, toSpace, "@")
docs <- tm_map(docs, toSpace, "\\|")
```

This can be done with a single call:

```
docs <- tm_map(docs, toSpace, "/|@|\\|")
```

Check the email address in the following.

```
inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## Hybrid weighted random forests for
## classifying very high-dimensional data
## Baoxun Xu1 , Joshua Zhexue Huang2 , Graham Williams2 and
## Yunming Ye1
## 1
##
## Department of Computer Science, Harbin Institute of Technology Shenzhen Gr...
## School, Shenzhen 518055, China
## 2
## Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, S...
## 518055, China
## Email: amusing002 gmail.com
## Random forests are a popular classification method based on an ensemble of a
## single type of decision trees from subspaces of data. In the literature, t...
## are many different types of decision tree algorithms, including C4.5, CART...
## CHAID. Each type of decision tree algorithm may capture different information
## and structure. This paper proposes a hybrid weighted random forest algorithm,
## simultaneously using a feature weighting method and a hybrid forest method to
## classify very high dimensional data. The hybrid weighted random forest alg...
## can effectively reduce subspace size and improve classification performance
## without increasing the error bound. We conduct a series of experiments on ...
## high dimensional datasets to compare our method with traditional random fo...
....
```

3.2 Conversion to Lower Case

```
docs <- tm_map(docs, content_transformer(tolower))

inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests for
## classifying very high-dimensional data
## baoxun xul , joshua zhexue huang2 , graham williams2 and
## yunming ye1
## 1
##
## department of computer science, harbin institute of technology shenzhen gr...
## school, shenzhen 518055, china
## 2
## shenzhen institutes of advanced technology, chinese academy of sciences, s...
## 518055, china
## email: amusing002 gmail.com
## random forests are a popular classification method based on an ensemble of a
## single type of decision trees from subspaces of data. in the literature, t...
## are many different types of decision tree algorithms, including c4.5, cart...
## chaid. each type of decision tree algorithm may capture different information
## and structure. this paper proposes a hybrid weighted random forest algorithm,
## simultaneously using a feature weighting method and a hybrid forest method to
## classify very high dimensional data. the hybrid weighted random forest alg...
## can effectively reduce subspace size and improve classification performance
## without increasing the error bound. we conduct a series of experiments on ...
## high dimensional datasets to compare our method with traditional random fo...
....
```

General character processing functions in R can be used to transform our corpus. A common requirement is to map the documents to lower case, using `tolower()`. As above, we need to wrap such functions with a `content_transformer()`:

3.3 Remove Numbers

```
docs <- tm_map(docs, removeNumbers)

inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests for
## classifying very high-dimensional data
## baoxun xu , joshua zhexue huang , graham williams and
## yunming ye
##
##
## department of computer science, harbin institute of technology shenzhen gr...
## school, shenzhen , china
##
## shenzhen institutes of advanced technology, chinese academy of sciences, s...
## , china
## email: amusing gmail.com
## random forests are a popular classification method based on an ensemble of a
## single type of decision trees from subspaces of data. in the literature, t...
## are many different types of decision tree algorithms, including c., cart, and
## chaid. each type of decision tree algorithm may capture different information
## and structure. this paper proposes a hybrid weighted random forest algorithm,
## simultaneously using a feature weighting method and a hybrid forest method to
## classify very high dimensional data. the hybrid weighted random forest alg...
## can effectively reduce subspace size and improve classification performance
## without increasing the error bound. we conduct a series of experiments on ...
## high dimensional datasets to compare our method with traditional random fo...
....
```

Numbers may or may not be relevant to our analyses. This transform can remove numbers simply.

3.4 Remove Punctuation

```
docs <- tm_map(docs, removePunctuation)

inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests for
## classifying very highdimensional data
## baoxun xu joshua zhexue huang graham williams and
## yunming ye
##
##
## department of computer science harbin institute of technology shenzhen gra...
## school shenzhen china
##
## shenzhen institutes of advanced technology chinese academy of sciences she...
## china
## email amusing gmailcom
## random forests are a popular classification method based on an ensemble of a
## single type of decision trees from subspaces of data in the literature there
## are many different types of decision tree algorithms including c cart and
## chaid each type of decision tree algorithm may capture different information
## and structure this paper proposes a hybrid weighted random forest algorithm
## simultaneously using a feature weighting method and a hybrid forest method to
## classify very high dimensional data the hybrid weighted random forest algo...
## can effectively reduce subspace size and improve classification performance
## without increasing the error bound we conduct a series of experiments on e...
## high dimensional datasets to compare our method with traditional random fo...
....
```

Punctuation can provide gramatical context which supports understanding. Often for initial analyses we ignore the punctuation. Later we will use punctuation to support the extraction of meaning.

3.5 Remove English Stop Words

```
docs <- tm_map(docs, removeWords, stopwords("english"))

inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests
## classifying highdimensional data
## baoxun xu joshua zhexue huang graham williams
## yunming ye
##
##
## department computer science harbin institute technology shenzhen graduate
## school shenzhen china
##
## shenzhen institutes advanced technology chinese academy sciences shenzhen
## china
## email amusing gmailcom
## random forests popular classification method based ensemble
## single type decision trees subspaces data literature
## many different types decision tree algorithms including c cart
## chaid type decision tree algorithm may capture different information
## structure paper proposes hybrid weighted random forest algorithm
## simultaneously using feature weighting method hybrid forest method
## classify high dimensional data hybrid weighted random forest algorithm
## can effectively reduce subspace size improve classification performance
## without increasing error bound conduct series experiments eight
## high dimensional datasets compare method traditional random forest
....
```

Stop words are common words found in a language. Words like *for*, *very*, *and*, *of*, *are*, etc, are common stop words. Notice they have been removed from the above text.

We can list the stop words:

```
length(stopwords("english"))

## [1] 174

stopwords("english")

## [1] "i" "me" "my" "myself" "we"
## [6] "our" "ours" "ourselves" "you" "your"
## [11] "yours" "yourself" "yourselves" "he" "him"
## [16] "his" "himself" "she" "her" "hers"
....
```

3.6 Remove Own Stop Words

```
docs <- tm_map(docs, removeWords, c("department", "email"))

inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests
## classifying highdimensional data
## baoxun xu joshua zhexue huang graham williams
## yunming ye
##
##
## computer science harbin institute technology shenzhen graduate
## school shenzhen china
##
## shenzhen institutes advanced technology chinese academy sciences shenzhen
## china
## amusing gmailcom
## random forests popular classification method based ensemble
## single type decision trees subspaces data literature
## many different types decision tree algorithms including c cart
## chaid type decision tree algorithm may capture different information
## structure paper proposes hybrid weighted random forest algorithm
## simultaneously using feature weighting method hybrid forest method
## classify high dimensional data hybrid weighted random forest algorithm
## can effectively reduce subspace size improve classification performance
## without increasing error bound conduct series experiments eight
## high dimensional datasets compare method traditional random forest
....
```

Previously we used the English stopwords provided by `tm`. We could instead or in addition remove our own stop words as we have done above. We have chosen here two words, simply for illustration. The choice might depend on the domain of discourse, and might not become apparent until we've done some analysis.

3.7 Strip Whitespace

```
docs <- tm_map(docs, stripWhitespace)

inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests
## classifying highdimensional data
## baoxun xu joshua zhexue huang graham williams
## yunming ye
##
##
## computer science harbin institute technology shenzhen graduate
## school shenzhen china
##
## shenzhen institutes advanced technology chinese academy sciences shenzhen
## china
## amusing gmailcom
## random forests popular classification method based ensemble
## single type decision trees subspaces data literature
## many different types decision tree algorithms including c cart
## chaid type decision tree algorithm may capture different information
## structure paper proposes hybrid weighted random forest algorithm
## simultaneously using feature weighting method hybrid forest method
## classify high dimensional data hybrid weighted random forest algorithm
## can effectively reduce subspace size improve classification performance
## without increasing error bound conduct series experiments eight
## high dimensional datasets compare method traditional random forest
....
```

3.8 Specific Transformations

We might also have some specific transformations we would like to perform. The examples here may or may not be useful, depending on how we want to analyse the documents. This is really for illustration using the part of the document we are looking at here, rather than suggesting this specific transform adds value.

```
toString <- content_transformer(function(x, from, to) gsub(from, to, x))
docs <- tm_map(docs, toString, "harbin institute technology", "HIT")
docs <- tm_map(docs, toString, "shenzhen institutes advanced technology", "SIAT")
docs <- tm_map(docs, toString, "chinese academy sciences", "CAS")
```

```
inspect(docs[16])

## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weighted random forests
## classifying highdimensional data
## baoxun xu joshua zhexue huang graham williams
## yunming ye
##
##
## computer science HIT shenzhen graduate
## school shenzhen china
##
## SIAT CAS shenzhen
## china
## amusing gmailcom
## random forests popular classification method based ensemble
## single type decision trees subspaces data literature
## many different types decision tree algorithms including c cart
## chaid type decision tree algorithm may capture different information
## structure paper proposes hybrid weighted random forest algorithm
## simultaneously using feature weighting method hybrid forest method
## classify high dimensional data hybrid weighted random forest algorithm
## can effectively reduce subspace size improve classification performance
## without increasing error bound conduct series experiments eight
## high dimensional datasets compare method traditional random forest
....
```


3.9 Stemming

```
library(SnowballC)
docs <- tm_map(docs, stemDocument)
```

Notice we load the **SnowballC** (Bouchet-Valat, 2014) package which provides stemming.

```
inspect(docs[16])
## <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
##
## [[1]]
## <<PlainTextDocument (metadata: 7)>>
## hybrid weight random forest
## classifi highdimension data
## baoxun xu joshua zhexu huang graham william
## yunm ye
##
##
## comput scienc HIT shenzhen graduat
## school shenzhen china
##
## SIAT CAS shenzhen
## china
## amus gmailcom
## random forest popular classif method base ensembl
## singl type decis tree subspac data literatur
## mani differ type decis tree algorithm includ c cart
## chaid type decis tree algorithm may captur differ inform
## structur paper propos hybrid weight random forest algorithm
## simultan use featur weight method hybrid forest method
## classifi high dimension data hybrid weight random forest algorithm
## can effect reduc subspac size improv classif perform
## without increas error bound conduct seri experi eight
## high dimension dataset compar method tradit random forest
....
```

Stemming uses an algorithm that removes common word endings for English words, such as “es”, “ed” and “s”. The functionality for stemming is provided by `wordStem()` from **SnowballC** (Bouchet-Valat, 2014).

4 Creating a Document Term Matrix

A document term matrix is simply a matrix with documents as the rows and terms as the columns and a count of the frequency of words as the cells of the matrix. We use `DocumentTermMatrix()` to create the matrix:

```
dtm <- DocumentTermMatrix(docs)
dtm

## <<DocumentTermMatrix (documents: 46, terms: 6508)>>
## Non-/sparse entries: 30061/269307
## Sparsity           : 90%
## Maximal term length: 56
## Weighting          : term frequency (tf)
```

We can inspect the document term matrix using `inspect()`. Here, to avoid too much output, we select a subset of inspect.

```
inspect(dtm[1:5, 1000:1005])

## <<DocumentTermMatrix (documents: 5, terms: 6)>>
## Non-/sparse entries: 7/23
## Sparsity           : 77%
## Maximal term length: 9
....
```

The document term matrix is in fact quite sparse (that is, mostly empty) and so it is actually stored in a much more compact representation internally. We can still get the row and column counts.

```
class(dtm)

## [1] "DocumentTermMatrix"      "simple_triplet_matrix"

dim(dtm)

## [1] 46 6508
```

The transpose is created using `TermDocumentMatrix()`:

```
tdm <- TermDocumentMatrix(docs)
tdm

## <<TermDocumentMatrix (terms: 6508, documents: 46)>>
## Non-/sparse entries: 30061/269307
## Sparsity           : 90%
## Maximal term length: 56
## Weighting          : term frequency (tf)
```

We will use the document term matrix for the remainder of the chapter.

5 Exploring the Document Term Matrix

We can obtain the term frequencies as a vector by converting the document term matrix into a matrix and summing the column counts:

```
freq <- colSums(as.matrix(dtm))
length(freq)
## [1] 6508
```

By ordering the frequencies we can list the most frequent terms and the least frequent terms:

```
ord <- order(freq)

# Least frequent terms
freq[head(ord)]

## aaaaaaeaceeaeieaeaeiiaiaiciaiicaiaaeaoeneiacaeaaeoouoo
## 1
## aab 1
## 1
## aadrbltn 1
## 1
## aadrhtmliv 1
## 1
## aai 1
## 1
## ...
```

Notice these terms appear just once and are probably not really terms that are of interest to us. Indeed they are likely to be spurious terms introduced through the translation of the original document from PDF to text.

```
# Most frequent terms
freq[tail(ord)]

##      can dataset pattern      use      mine      data
##    709      776      887    1366    1446    3101
```

These terms are much more likely to be of interest to us. Not surprising, given the choice of documents in the corpus, the most frequent terms are: data, mine, use, pattern, dataset, can.

6 Distribution of Term Frequencies

```
# Frequency of frequencies.
head(table(freq), 15)

## freq
##    1    2    3    4    5    6    7    8    9   10   11   12   13   14   15
## 2381 1030  503  311  210  188  134  130   82   83   65   61   54   52   51

tail(table(freq), 15)

## freq
##  483  544  547  555  578  609  611  616  703  709  776  887 1366 1446 3101
##    1    1    1    1    1    1    1    1    1    1    1    1    1    1    1
```

So we can see here that there are 2381 terms that occur just once.

Generate a plot of term frequency.

7 Conversion to Matrix and Save to CSV

We can convert the document term matrix to a simple matrix for writing to a CSV file, for example, for loading the data into other software if we need to do so. To write to CSV we first convert the data structure into a simple matrix:

```
m <- as.matrix(dtm)
dim(m)
## [1] 46 6508
```

For very large corpus the size of the matrix can exceed R's calculation limits. This will manifest itself as an integer overflow error with a message like:

```
## Error in vector(typeof(x$v), nr * nc) : vector size cannot be NA
## In addition: Warning message:
## In nr * nc : NAs produced by integer overflow
```

If this occurs, then consider removing sparse terms from the document term matrix, as we discuss shortly.

Once converted into a standard matrix the usual `write.csv()` can be used to write the data to file.

```
write.csv(m, file="dtm.csv")
```

8 Removing Sparse Terms

We are often not interested in infrequent terms in our documents. Such “sparse” terms can be removed from the document term matrix quite easily using `removeSparseTerms()`:

```
dim(dtm)
## [1] 46 6508
dtms <- removeSparseTerms(dtm, 0.1)
dim(dtms)
## [1] 46 6
```

This has removed most terms!

```
inspect(dtms)
## <<DocumentTermMatrix (documents: 46, terms: 6)>>
## Non-/sparse entries: 257/19
## Sparsity           : 7%
## Maximal term length: 7
....
```

We can see the effect by looking at the terms we have left:

```
freq <- colSums(as.matrix(dtms))
freq
## data graham inform time use william
## 3101 108 467 483 1366 236
table(freq)
## freq
## 108 236 467 483 1366 3101
## 1 1 1 1 1 1
```

9 Identifying Frequent Items and Associations

One thing we often to first do is to get an idea of the most frequent terms in the corpus. We use `findFreqTerms()` to do this. Here we limit the output to those terms that occur at least 1,000 times:

```
findFreqTerms(dtm, lowfreq=1000)
## [1] "data" "mine" "use"
```

So that only lists a few. We can get more of them by reducing the threshold:

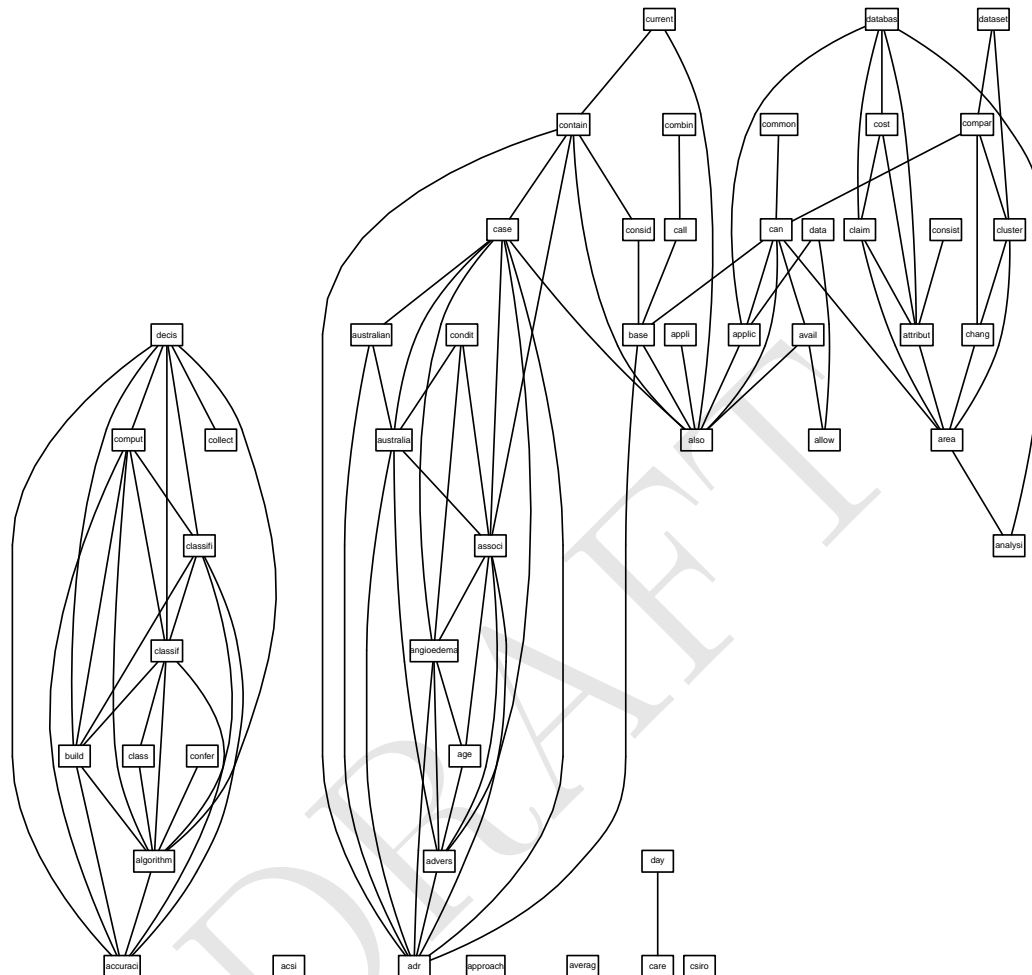
```
findFreqTerms(dtm, lowfreq=100)
## [1] "accuraci" "acsi" "adr" "advers" "age"
## [6] "algorithm" "allow" "also" "analysi" "angioedema"
## [11] "appli" "applic" "approach" "area" "associ"
## [16] "attribut" "australia" "australian" "avail" "averag"
....
```

We can also find associations with a word, specifying a correlation limit.

```
findAssocs(dtm, "data", corlimit=0.6)
##          data
## mine      0.90
## induct    0.72
## challeng 0.70
....
```

If two words always appear together then the correlation would be 1.0 and if they never appear together the correlation would be 0.0. Thus the correlation is a measure of how closely associated the words are in the corpus.

10 Correlations Plots

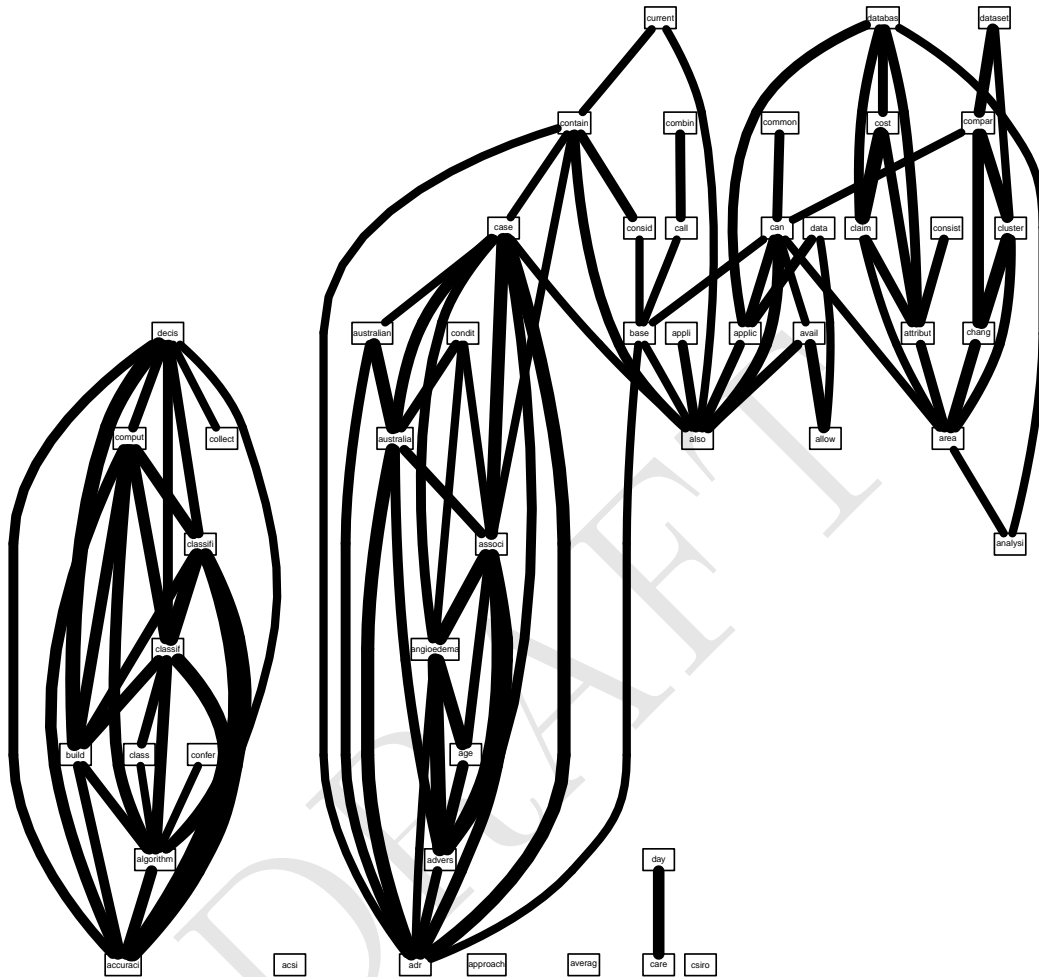


```
plot(dtm,
     terms=findFreqTerms(dtm, lowfreq=100)[1:50],
     corThreshold=0.5)
```

Rgraphviz (Gentry *et al.*, 2014) from the BioConductor repository for R (bioconductor.org) is used to plot the network graph that displays the correlation between chosen words in the corpus. Here we choose 50 of the more frequent words as the nodes and include links between words when they have at least a correlation of 0.5.

By default (without providing terms and a correlation threshold) the plot function chooses a random 20 terms with a threshold of 0.7.

11 Correlations Plot—Options



```
plot(dtm,
     terms=findFreqTerms(dtm, lowfreq=100)[1:50],
     corThreshold=0.5)
```

12 Plotting Word Frequencies

We can generate the frequency count of all words in a corpus:

```
freq <- sort(colSums(as.matrix(dtm)), decreasing=TRUE)
head(freq, 14)

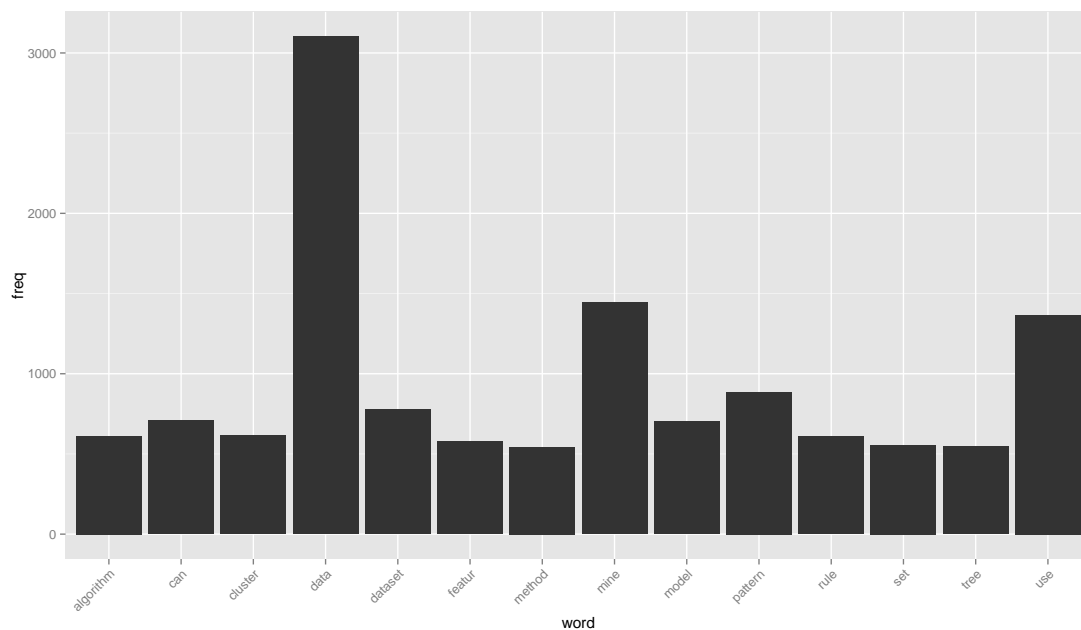
##      data      mine      use  pattern  dataset      can      model
##    3101     1446     1366      887      776      709      703
## cluster algorithm    rule    featur      set      tree    method
##      616      611      609      578      555      547      544
## ....

wf <- data.frame(word=names(freq), freq=freq)
head(wf)

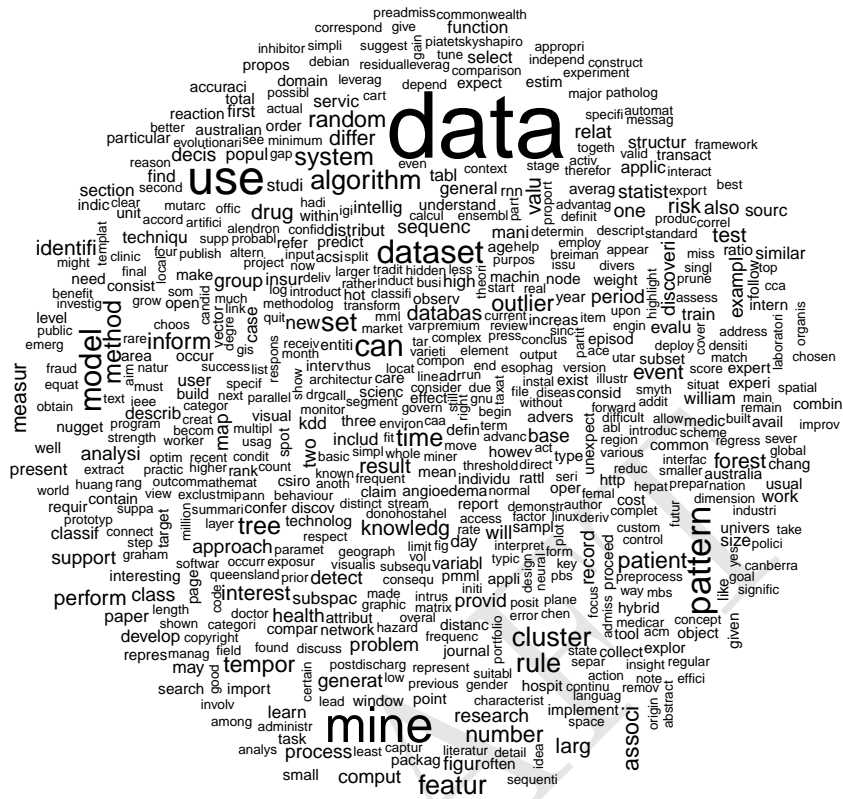
##      word freq
## data    data 3101
## mine    mine 1446
## use     use  1366
## ....
```

We can then plot the frequency of those words that occur at least 500 times in the corpus:

```
library(ggplot2)
subset(wf, freq>500) %>%
  ggplot(aes(word, freq)) +
  geom_bar(stat="identity") +
  theme(axis.text.x=element_text(angle=45, hjust=1))
```



Word Clouds



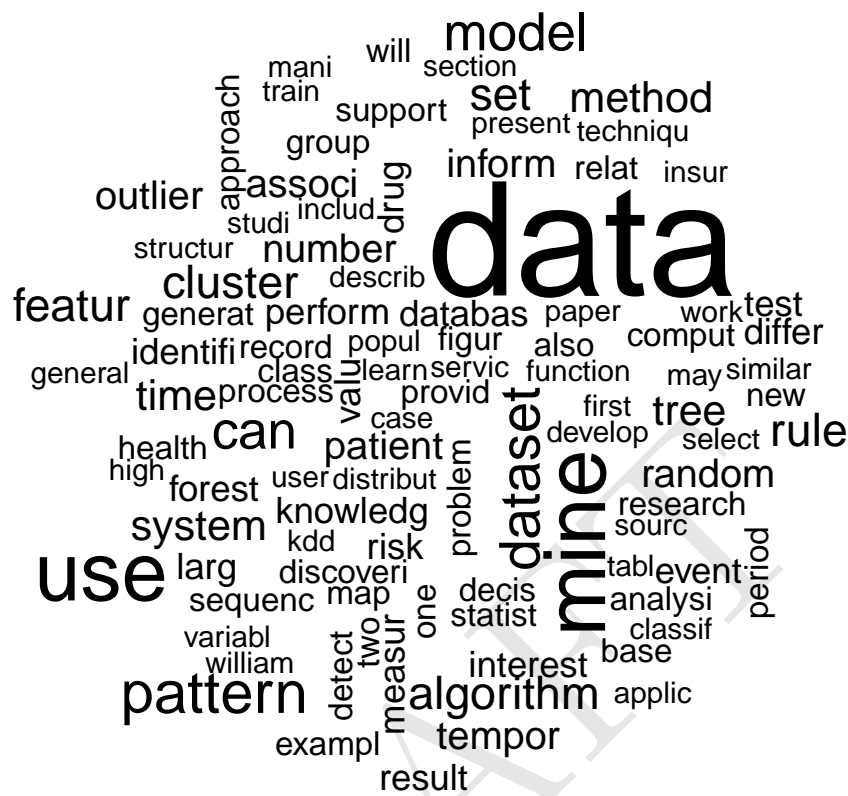
We can generate a word cloud as an effective alternative to providing a quick visual overview of the frequency of words in a corpus.

The wordcloud (?) package provides the required function.

```
library(wordcloud)
set.seed(123)
wordcloud(names(freq), freq, min.freq=40)
```

Notice the use of `set.seed()` only so that we can obtain the same layout each time—otherwise a random layout is chosen, which is not usually an issue.

13.1 Reducing Clutter With Max Words



To increase or reduce the number of words displayed we can tune the value of `max.words=`. Here we have limited the display to the 100 most frequent words.

```
set.seed(142)
wordcloud(names(freq), freq, max.words=100)
```


[illegible]

```
set.seed(142)
wordcloud(names(freq), freq, min.freq=100, colors=brewer.pal(6, "Dark2"))
```

[illegible]

```
set.seed(142)
wordcloud(names(freq), freq, min.freq=100, scale=c(5, .1), colors=brewer.pal(6, "Dark2"))
```

[illegible]

```
set.seed(142)
dark2 <- brewer.pal(6, "Dark2")
wordcloud(names(freq), freq, min.freq=100, rot.per=0.2, colors=dark2)
```


14 Quantitative Analysis of Text

The `qdap` (Rinker, 2014) package provides an extensive suite of functions to support the quantitative analysis of text.

We can obtain simple summaries of a list of words, and to do so we will illustrate with the terms from our Term Document Matrix *tdm*. We first extract the shorter terms from each of our documents into one long word list. To do so we convert *tdm* into a matrix, extract the column names (the terms) and retain those shorter than 20 characters.

```
words <- dtm                                     %>%
  as.matrix                                       %>%
  colnames                                         %>%
  (function(x) x[nchar(x) < 20])
```

We can then summarise the word list. Notice, in particular, the use of `dist_tab()` from `qdap` to generate frequencies and percentages.

```
length(words)
## [1] 6456

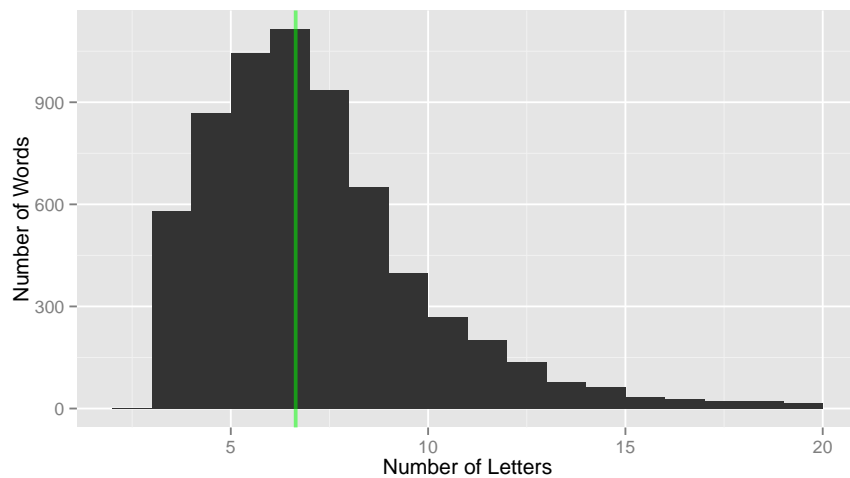
head(words, 15)
## [1] "aaai"      "aab"      "aad"      "aadbhtml" "aadbhtmln"
## [6] "aadbhtmliv" "aai"      "aam"      "aba"      "abbrev"
## [11] "abbrevi"   "abc"      "abcd"     "abdul"    "abel"

summary(nchar(words))
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      3.000  5.000   6.000   6.644  8.000  19.000

table(nchar(words))
##
##      3      4      5      6      7      8      9     10     11     12     13     14     15     16     17
## 579  867 1044 1114  935  651  397  268  200  138   79   63   34   28   22
##    18    19
##    21    16

dist_tab(nchar(words))
##      interval freq cum.freq percent cum.percent
## 1           3  579      579    8.97      8.97
## 2           4  867     1446   13.43     22.40
## 3           5 1044     2490   16.17     38.57
## 4           6 1114     3604   17.26     55.82
## 5           7  935     4539   14.48     70.31
## 6           8  651     5190   10.08     80.39
## 7           9  397     5587    6.15     86.54
## 8          10  268     5855    4.15     90.69
## 9          11  200     6055    3.10     93.79
## 10         12  138     6193    2.14     95.93
## ....
```

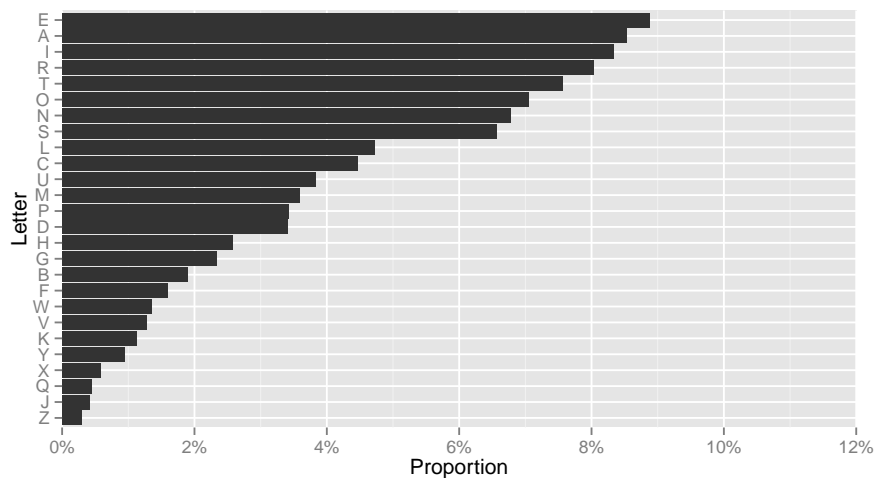
14.1 Word Length Counts



A simple plot is then effective in showing the distribution of the word lengths. Here we create a single column data frame that is passed on to `ggplot()` to generate a histogram, with a vertical line to show the mean length of words.

```
data.frame(nletters=nchar(words)) %>%  
  ggplot(aes(x=nletters)) +  
  geom_histogram(binwidth=1) +  
  geom_vline(xintercept=mean(nchar(words)),  
             colour="green", size=1, alpha=.5) +  
  labs(x="Number of Letters", y="Number of Words")
```

14.2 Letter Frequency



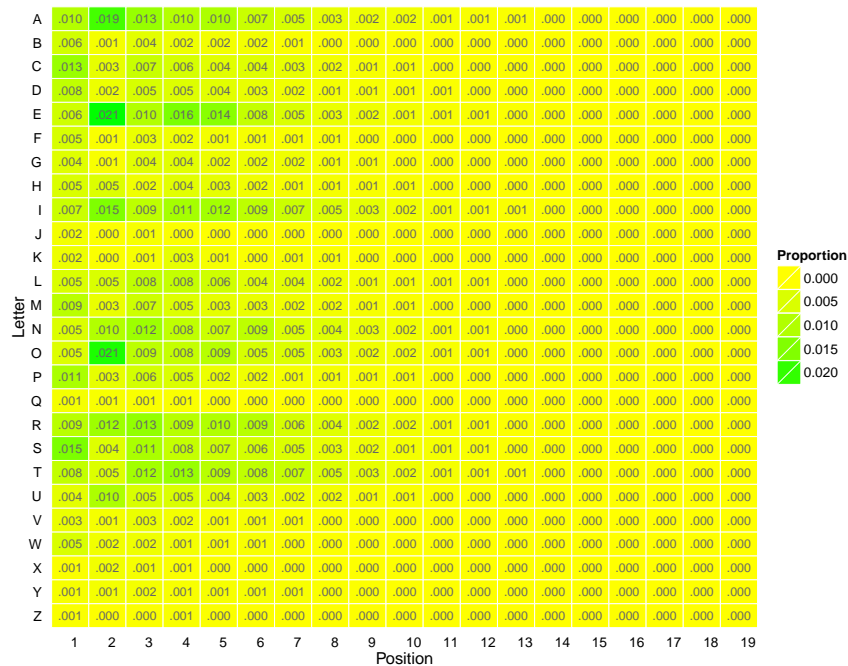
Next we want to review the frequency of letters across all of the words in the discourse. Some data preparation will transform the vector of words into a list of letters, which we then construct a frequency count for, and pass this on to be plotted.

We again use a pipeline to string together the operations on the data. Starting from the vector of words stored in *word* we split the words into characters using `str_split()` from *stringr* (Wickham, 2012), removing the first string (an empty string) from each of the results (using `sapply()`). Reducing the result into a simple vector, using `unlist()`, we then generate a data frame recording the letter frequencies, using `dist_tab()` from *qdap*. We can then plot the letter proportions.

```
library(dplyr)
library(stringr)

words %>%
  str_split("") %>%
  sapply(function(x) x[-1]) %>%
  unlist %>%
  dist_tab %>%
  mutate(Letter=factor(toupper(interval),
                       levels=toupper(interval[order(freq)]))) %>%
  ggplot(aes(Letter, weight=percent)) +
  geom_bar() +
  coord_flip() +
  ylab("Proportion") +
  scale_y_continuous(breaks=seq(0, 12, 2),
                     label=function(x) paste0(x, "%"),
                     expand=c(0,0), limits=c(0,12))
```

14.3 Letter and Position Heatmap



The `qheat()` function from `qdap` provides an effective visualisation of tabular data. Here we transform the list of words into a position count of each letter, and constructing a table of the proportions that is passed on to `qheat()` to do the plotting.

```
words %>%
  lapply(function(x) sapply(letters, gregexpr, x, fixed=TRUE)) %>%
  unlist %>%
  (function(x) x[x!=-1]) %>%
  (function(x) setNames(x, gsub("\\d", "", names(x)))) %>%
  (function(x) apply(table(data.frame(letter=toupper(names(x)),
                                     position=unname(x))),
                     1, function(y) y/length(x))) %>%
  qheat(high="green", low="yellow", by.column=NULL,
        values=TRUE, digits=3, plot=FALSE) +
  ylab("Letter") +
  xlab("Position") +
  theme(axis.text.x=element_text(angle=0)) +
  guides(fill=guide_legend(title="Proportion"))
```

14.4 Miscellaneous Functions

We can generate gender from a name list, using the `genderdata` (?) package

```
devtools::install_github("lmullen/gender-data-pkg")
```

```
name2sex(qcv(graham, frank, leslie, james, jacqui, jack, kerry, kerrie))  
## [1] M M F M F M F F  
## Levels: F M
```

DRAFT

15 Review—Preparing the Corpus

Here in one sequence is collected the code to perform a text mining project. Notice that we would not necessarily do all of these steps so pick and choose as is appropriate to your situation.

```
# Required packages

library(tm)
library(wordcloud)

# Locate and load the Corpus.

cname <- file.path(".", "corpus", "txt")
docs <- Corpus(DirSource(cname))

docs
summary(docs)
inspect(docs[1])

# Transforms

toSpace <- content_transformer(function(x, pattern) gsub(pattern, " ", x))
docs <- tm_map(docs, toSpace, "/|@|\\|")

docs <- tm_map(docs, content_transformer(tolower))
docs <- tm_map(docs, removeNumbers)
docs <- tm_map(docs, removePunctuation)
docs <- tm_map(docs, removeWords, stopwords("english"))
docs <- tm_map(docs, removeWords, c("own", "stop", "words"))
docs <- tm_map(docs, stripWhitespace)

toString <- content_transformer(function(x, from, to) gsub(from, to, x))
docs <- tm_map(docs, toString, "specific transform", "ST")
docs <- tm_map(docs, toString, "other specific transform", "OST")

docs <- tm_map(docs, stemDocument)
```

16 Review—Analysing the Corpus

```
# Document term matrix.

dtm <- DocumentTermMatrix(docs)
inspect(dtm[1:5, 1000:1005])

# Explore the corpus.

findFreqTerms(dtm, lowfreq=100)
findAssocs(dtm, "data", corlimit=0.6)

freq <- sort(colSums(as.matrix(dtm)), decreasing=TRUE)
wf <- data.frame(word=names(freq), freq=freq)

library(ggplot2)

p <- ggplot(subset(wf, freq>500), aes(word, freq))
p <- p + geom_bar(stat="identity")
p <- p + theme(axis.text.x=element_text(angle=45, hjust=1))

# Generate a word cloud

library(wordcloud)
wordcloud(names(freq), freq, min.freq=100, colors=brewer.pal(6, "Dark2"))
```

17 Further Reading and Acknowledgements

The [Rattle Book](#), published by Springer, provides a comprehensive introduction to data mining and analytics using Rattle and R. It is available from [Amazon](#). Other documentation on a broader selection of R topics of relevance to the data scientist is freely available from <http://datamining.togaware.com>, including the [Datamining Desktop Survival Guide](#).

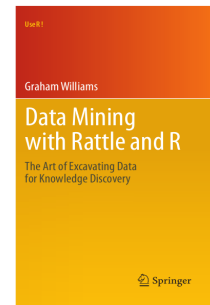
This chapter is one of many chapters available from <http://HandsOnDataScience.com>. In particular follow the links on the website with a * which indicates the generally more developed chapters.

Other resources include:

- The Journal of Statistical Software article, *Text Mining Infrastructure in R* is a good start <http://www.jstatsoft.org/v25/i05/paper>
- [Bilisoly \(2008\)](#) presents methods and algorithms for text mining using Perl.

Thanks also to Tony Nolan for suggestions of some of the examples used in this chapter.

Some of the qdap examples were motivated by <http://trinkerrstuff.wordpress.com/2014/10/31/exploration-of-letter-make-up-of-english-words/>.



18 References

- Bilisoly R (2008). *Practical Text Mining with Perl*. Wiley Series on Methods and Applications in Data Mining. Wiley. ISBN 9780470382851. URL <http://books.google.com.au/books?id=YkMFVbsrdzkC>.
- Bouchet-Valat M (2014). *SnowballC: Snowball stemmers based on the C libstemmer UTF-8 library*. R package version 0.5.1, URL <http://CRAN.R-project.org/package=SnowballC>.
- Feinerer I, Hornik K (2014). *tm: Text Mining Package*. R package version 0.6, URL <http://CRAN.R-project.org/package=tm>.
- Gentry J, Long L, Gentleman R, Falcon S, Hahne F, Sarkar D, Hansen KD (2014). *Rgraphviz: Provides plotting capabilities for R graph objects*. R package version 2.6.0.
- Neuwirth E (2011). *RColorBrewer: ColorBrewer palettes*. R package version 1.0-5, URL <http://CRAN.R-project.org/package=RColorBrewer>.
- R Core Team (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- Rinker T (2014). *qdap: Bridging the Gap Between Qualitative Data and Quantitative Analysis*. R package version 2.2.0, URL <http://CRAN.R-project.org/package=qdap>.
- Wickham H (2012). *stringr: Make it easier to work with strings*. R package version 0.6.2, URL <http://CRAN.R-project.org/package=stringr>.
- Williams GJ (2009). "Rattle: A Data Mining GUI for R." *The R Journal*, **1**(2), 45–55. URL http://journal.r-project.org/archive/2009-2/RJournal_2009-2_Williams.pdf.
- Williams GJ (2011). *Data Mining with Rattle and R: The art of excavating data for knowledge discovery*. Use R! Springer, New York. URL http://www.amazon.com/gp/product/1441998896/ref=as_li_qf_sp_asin_tl?ie=UTF8&tag=togaware-20&linkCode=as2&camp=217145&creative=399373&creativeASIN=1441998896.

This document, sourced from TextMiningO.Rnw revision 531, was processed by KnitR version 1.7 of 2014-10-13 and took 36.8 seconds to process. It was generated by gjw on nix running Ubuntu 14.04.1 LTS with Intel(R) Xeon(R) CPU W3520 @ 2.67GHz having 8 cores and 12.3GB of RAM. It completed the processing 2014-11-05 19:22:29.