

# Introduction to Narrative

*Christopher Lovell*

*July 6, 2015*

## Introduction

The *Narrative* package is designed to extend on the ubiquitous *tm* package, making time series analysis of corpora easier, as well additional functionality.

If you have not done any text mining in R before I recommend reading the [tm introduction](#) vignette first. The following document provides only a quick introduction to the *tm* package before moving on to demonstrate the additional features provided by *Narrative*.

## tm

For this demo we'll be using a corpus of Reuters documents provided with the *tm* package. To load your own data, *Narrative* provides the `readSeparateText()` function for reading a directory full of .txt files. This is discussed in detail later.

```
library(tm)
```

```
## Loading required package: NLP
```

```
## Warning: package 'NLP' was built under R version 3.2.1
```

```
reut21578 <- system.file("texts", "crude", package = "tm")
reuters <- VCorpus(DirSource(reut21578), readerControl = list(reader = readReut21578XMLasPlain))
```

```
reuters
```

```
## <<VCorpus>>
## Metadata:  corpus specific: 0, document level (indexed): 0
## Content:   documents: 20
```

The corpus contains 20 documents. To look at a document in detail, use `as.character()`.

```
as.character(reuters[[1]])
```

```
## [1] "Diamond Shamrock Corp said that\neffective today it had cut its contract prices for crude oil by\n1.5
```

The meta data contains the title, language, some high level topics, place of origin, and a time stamp; this will be important later when we begin to use the *Narrative* package.

```
meta(reuters[[1]])
```

```
## author      : character(0)
## timestamp   : 1987-02-26 17:00:56
## description :
## heading     : DIAMOND SHAMROCK (DIA) CUTS CRUDE PRICES
## id          : 127
## language    : en
## origin      : Reuters-21578 XML
## topics      : YES
## lewissplit  : TRAIN
## cgisplit    : TRAINING-SET
## oldid       : 5670
## topics_cat  : crude
## places      : usa
## people      : character(0)
## orgs        : character(0)
## exchanges   : character(0)
```

After reading in your data the next stage is to cleanse it and stem it. *tm* provides a number of functions for doing so.

```
reuters <- tm_map(reuters, content_transformer(tolower))
reuters <- tm_map(reuters, content_transformer(removeNumbers))
reuters <- tm_map(reuters, content_transformer(removePunctuation))
reuters <- tm_map(reuters, content_transformer(removeWords), stopwords("english"))
reuters <- tm_map(reuters, stemDocument)
reuters <- tm_map(reuters, content_transformer(stripWhitespace))
```

```
as.character(reuters[[1]])
```

```
## [1] "diamond shamrock corp said effect today cut contract price crude oil dlrs barrel reduct bring post pr"
```

An incredibly useful representation of the corpus is as a *document term matrix*. This is a *td* matrix, where *t* is a vector of all terms used in the corpus, and *d* is a vector of all documents in the corpus. Each row then represents the frequency of a term in each document, and each row shows a documents constituent terms and their frequency.

```
dtm <- DocumentTermMatrix(reuters)
inspect(dtm[1:5,1:20])
```

```
## <<DocumentTermMatrix (documents: 5, terms: 20)>>
## Non-/sparse entries: 6/94
## Sparsity           : 94%
## Maximal term length: 10
## Weighting          : term frequency (tf)
##
##      Terms
## Docs  abdulaziz abil abl abroad accept accord across act activity add
## 127      0    0  0      0      0      0      0  0      0  0
## 144      0    2  0      0      0      0      0  0  0      0  0
## 191      0    0  0      0      0      0      0  0  0      0  0
## 194      0    0  0      0      0      0      0  0  0      0  0
## 211      0    0  0      0      0      0      0  0  0      0  0
```

```
##      Terms
## Docs  added address adher advantag advisers agenc agr agre agreement
##  127    0      0      0      0      0      0  0  0  0      0
##  144    1      4      0      1      0      0  1  0      1
##  191    0      0      0      0      0      0  0  0      0
##  194    0      0      0      0      0      0  0  0      0
##  211    0      0      0      0      0      0  0  0      0
##      Terms
## Docs  agricultur
##  127      0
##  144      0
##  191      0
##  194      0
##  211      0
```

Some useful operations can be applied to the term document matrix, such as finding the most frequent terms, or terms that are highly correlated.

```
findFreqTerms(dtm, lowfreq = 15)
```

```
## [1] "barrel" "bpd"    "crude"  "dlrs"   "kuwait" "last"   "market"
## [8] "mln"    "offici" "oil"    "one"    "opec"   "price"  "reuter"
## [15] "said"   "saudi"  "will"
```

```
findAssocs(dtm, "opec", corlimit = 0.85)
```

```
## $opec
##      emerg      meet      analyst production      oil
##      0.92      0.92      0.91      0.87      0.86
```

## Narrative

Now that I have summarised the basic text mining functionality provided by *tm*, I'll move on to the additional functionality written in *Narrative*. I will use the `Narrative::function` syntax to distinguish between functions written in *Narrative* or other packages.

Word and character counts are often useful for normalising statistics generated from your corpus. We can generate vectors of each using the `characterCount` and `wordCount` functions. `characterCount` accepts the corpus as an argument, whereas `wordCount` uses a document term matrix.

```
wc <- Narrative::wordCount(dtm)
cc <- Narrative::characterCount(reuters)
```

Each function returns a numeric vector giving the word or character count for each document in our corpus. These vectors can be added to our corpus metadata using the `addToMetaData` function.

```
reuters <- Narrative::addToMetaData(corpus = reuters, vector = wc, tag = "word_count")
reuters <- Narrative::addToMetaData(corpus = reuters, vector = cc, tag = "character_count")
meta(reuters[[1]])
```

```
## author      : character(0)
## timestamp   : 1987-02-26 17:00:56
## description :
## heading     : DIAMOND SHAMROCK (DIA) CUTS CRUDE PRICES
## id          : 127
## language    : en
## origin      : Reuters-21578 XML
## topics      : YES
## lewissplit  : TRAIN
## cgisplit    : TRAINING-SET
## oldid       : 5670
## topics_cat  : crude
## places      : usa
## people      : character(0)
## orgs        : character(0)
## exchanges   : character(0)
## word_count  : 57
## character_count: 353
```

**IO** Narrative provides a couple of IO functions that make handling meta data associated with a corpus easier. To demonstrate, I'll first write our *reuters* corpus to disk using the `saveCorpus` function. This function takes a corpus, and writes it to the specified directory, along with an `.Rdata` file containing all associated metadata.

```
library(Narrative)
```

```
## Narrative v0.2
```

```
Narrative::saveCorpus(reuters,"C:\\Users\\324240\\Desktop\\test_files", save_metadata = T)
```

```
## metadata saved to " C:\Users\324240\Desktop\test_files " as " reuters_metadata.RDS "
```

We can then read this data back in to R, including the metadata, using the `readSeparateText` function.

```
reuters2 <- Narrative::readSeparateText("C:\\Users\\324240\\Desktop\\test_files",load_metadata = T,meta
```

## Term Frequency Analysis