

San Francisco Bay Guardian Paper Text Mining

Nelson Corrocher

August 21, 2016

The following section prepare the corpus of the text mining. In this exercise, the folder texts/ is being used so the text object must be put in that folder.

```
## [1] "C:/Users/Nelson Corrocher/WorkSpace/R/Week7/texts"
```

```
## [1] "SFBG.txt"
```

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

The following section pre-process the words (removes punctuation and number, convert to lowercase, remove stopwords, extract the word stem and remove whitespaces).

```
docs <- tm_map(docs, removePunctuation) # *Removing punctuation:*
docs <- tm_map(docs, removeNumbers)     # *Removing numbers:*
docs <- tm_map(docs, tolower)           # *Converting to lowercase:*
docs <- tm_map(docs, removeWords, stopwords("english")) # *Removing "stopwords"
library(SnowballC)
docs <- tm_map(docs, stemDocument)      # *Removing common word endings* (e.g., "ing", "es")
docs <- tm_map(docs, stripWhitespace)   # *Stripping whitespace
docs <- tm_map(docs, PlainTextDocument)
```

Staging the Data

```
dtm <- DocumentTermMatrix(docs)
# tdm <- TermDocumentMatrix(docs) This line is commented because it is not used
```

The following section is used for some data familiarization and to catch evident outliers:

```
freq <- colSums(as.matrix(dtm))
length(freq)
```

```
## [1] 7744
```

```
ord <- order(freq)
m <- as.matrix(dtm)
dim(m)
```

```
## [1] 1 7744
```

```
write.csv(m, file="DocumentTermMatrix.csv") # The excel file in the working directory now holds word co
```

Some cleaning up of the data below. The output shows two rows of numbers. The top number is the frequency with which words appear and the bottom number reflects how many words appear that frequently:

```
# Start by removing sparse terms:
dtms <- removeSparseTerms(dtm, 0.1) # This makes a matrix that is 10% empty space, maximum.
head(table(freq), 20) ### Word Frequency
```

```
## freq
##      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15
## 5107 1218  469  257  156  113   82   65   44   28   33   31   15    9    5
##      16     17     18     19     20
##      14     10     10      4      5
```

```
tail(table(freq), 20)
```

```
## freq
##  35  36  37  38  40  41  42  43  47  51  52  54  56  61  62  64  70  73
##   2   2   2   1   1   2   1   1   1   4   1   1   1   1   3   1   1   1
## 112 134
##   2   1
```

Next, let's consider only biggest frequencies (in this example, 42):

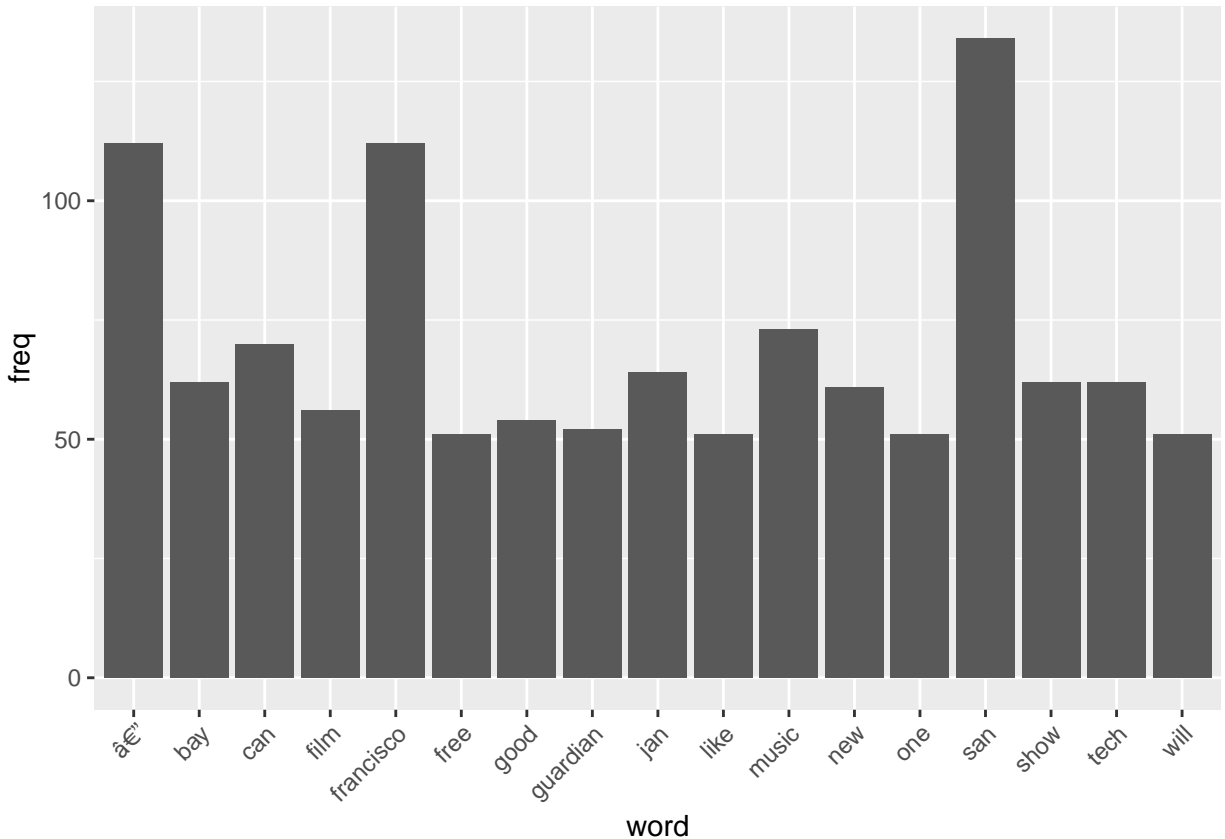
```
freq <- colSums(as.matrix(dtms)) # Matrix is too big to be shown on this document
f_list <- findFreqTerms(dtm, lowfreq=30) # <<<< Change it to the most appropriate for the data
f_list
```

```
## [1] "â&U+0080><U+0094>"      "also"      "arts"      "bay"      "call"
## [6] "can"      "city"      "classifi"  "cultur"   "drink"
## [11] "feb"      "film"      "first"     "food"     "francisco"
## [16] "fre"      "free"      "get"       "good"     "guardian"
## [21] "housing"  "ing"       "jan"       "just"     "last"
## [26] "like"     "music"     "new"       "news"     "now"
## [31] "one"      "open"      "people"    "san"      "sat"
## [36] "selector" "sfbg"      "show"      "sun"      "tech"
## [41] "time"     "will"
```

Plotting Word Frequencies (Only words that appear at least 50 times):

```
##
## Attaching package: 'ggplot2'

## The following object is masked from 'package:NLP':
##
##      annotate
```



Next, we can try to find the correlation between terms that occurs in the corpus. In this context, correlation is a quantitative measure of the co-occurrence of words in multiple documents. Since we have only one document in this corpus, it will always return 0. This, it was set not to show anything.

```
findAssocs(dtm,f_list ,corlimit = 0.0) # <<<< Adjust the corlimit to the desired correlation level. Obs
```

Here we create a word cloud based on the frequency of the words.

```
## Loading required package: RColorBrewer
```

