# Data Mining with R

## Text Mining

Hugh Murrell

## reference books

These slides are based on a book by Yanchang Zhao:

**R and Data Mining: Examples and Case Studies**.

*http://www.RDataMining.com*

for further background try Andrew Moore's slides from:
*http://www.autonlab.org/tutorials*

and as always, **wikipedia** is a useful source of information.

# text mining

This lecture presents examples of text mining with R.

We extract text from the BBC's webpages on Alastair Cook's letters from America. The extracted text is then transformed to build a term-document matrix.

Frequent words and associations are found from the matrix. A word cloud is used to present frequently occuring words in documents.

Words and transcripts are clustered to find groups of words and groups of transcripts.

In this lecture, *transcript* and *document* will be used interchangeably, so are *word* and *term*.

# text mining packages

many new packages are introduced in this lecture:

- ▶ **tm**: [Feinerer, 2012] provides functions for text mining,
- ▶ **wordcloud** [Fellows, 2012] visualizes results.
- ▶ **fpc** [Christian Hennig, 2005] flexible procedures for clustering.
- ▶ **igraph** [Gabor Csardi , 2012] a library and R package for network analysis.

# retrieving text from the BBC website

This work is part of the Rozanne Els PhD project

She has written a script to download transcripts direct from the website
http://www.bbc.co.uk/programmes/b00f6hbp/broadcasts/1947/01

The results are stored in a local directory, ACdatedfiles, on this apple mac.

# loading the corpus from disc

Now we are in a position to load the transcripts directly from our hard drive and perform corpus cleaning using the `tm` package.

```
> library(tm)
> corpus <- Corpus(
+     DirSource("./ACdatedfiles",
+       encoding = "UTF-8"),
+       readerControl = list(language = "en")
+     )
```

## cleaning the corpus

now we use regular expressions to remove at-tags and urls
from the remaining documents

```
> # get rid of html tags
> pattern <- "</?\\w+((\\s+\\w+(\\s*=\\s*(?:\".*?\"|'
> rmHTML <- function(x)
+   gsub(pattern, "", x)
> corpus <- tm_map(corpus, rmHTML)
> writeCorpus(corpus,path="./ac",
+             filenames = paste("d",
+                 seq_along(corpus),
+                 ".txt", sep = ""))
> corpus <- Corpus(
+     DirSource("./ac",
+       encoding = "UTF-8"),
+       readerControl = list(language = "en")
+     )
```

# further cleaning

now we use text cleaning transformations:

```
> # make each letter lowercase
> corpus <- tm_map(corpus, tolower)
> # remove white space
> corpus <- tm_map(corpus, stripWhitespace)
> # remove punctuation
> corpus <- tm_map(corpus, removePunctuation)
> # remove generic and custom stopwords
> my_stopwords <- c(stopwords('english'),
+  c('dont','didnt','arent','cant',
+    'one','two','three'))
> corpus <- tm_map(corpus,
+               removeWords, my_stopwords)
```

# stemming words

In many applications, words need to be stemmed to retrieve their radicals, so that various forms derived from a stem would be taken as the same when counting word frequency.

For instance, words update, updated and updating should all be stemmed to updat.

Sometimes stemming is counter productive so I chose not to do it here.

```
> # to carry out stemming
> # corpus <- tm_map(corpus, stemDocument,
> #                   language = "english")
```

# building a term-document matrix

A term-document matrix represents the relationship between terms and documents, where each row stands for a term and each column for a document, and an entry is the number of occurrences of the term in the document.

```
> (tdm <- TermDocumentMatrix(corpus))

A term-document matrix (54690 terms, 911 documents)
Non-/sparse entries: 545261/49277329
Sparsity           : 99%
Maximal term length: 33
Weighting          : term frequency (tf)
```

# frequent terms

Now we can have a look at the popular words in the term-document matrix,

```
> (tt <- findFreqTerms(tdm, lowfreq=1500))
 [1] "ago"       "american"  "called"
 [4] "came"      "can"       "country"
 [7] "day"       "every"     "first"
[10] "going"     "house"     "just"
[13] "know"      "last"      "like"
[16] "long"      "man"       "many"
[19] "much"      "never"     "new"
[22] "now"       "old"       "people"
[25] "president" "said"      "say"
[28] "states"    "think"     "time"
[31] "united"    "war"       "way"
[34] "well"      "will"      "year"
[37] "years"
```
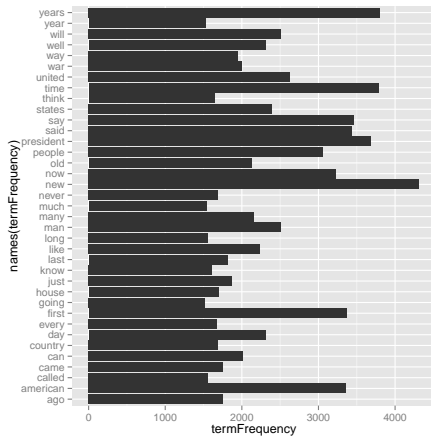
# frequent terms

Note that the frequent terms are ordered alphabetically, instead of by frequency or popularity. To show the top frequent words visually, we make a barplot of them.

```
> termFrequency <-
+   rowSums(as.matrix(tdm[tt,]))
> library(ggplot2)
> qplot(names(termFrequency),
+       termFrequency, geom="bar")  +
+   coord_flip()
```
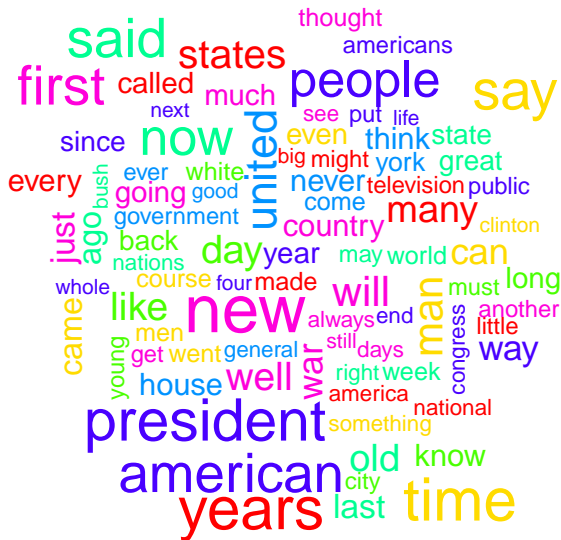
# frequent term bar chart

# wordclouds

We can show the importance of words pictorally with a *wordcloud* [Fellows, 2012]. In the code below, we first convert the term-document matrix to a normal matrix, and then calculate word frequencies. After that we use wordcloud to make a pictorial.

```
> tdmat = as.matrix(tdm)
> # calculate the frequency of words
> v = sort(rowSums(tdmat), decreasing=TRUE)
> d = data.frame(word=names(v), freq=v)
> # generate the wordcloud
> library(wordcloud)
> wordcloud(d$word, d$freq, min.freq=900,
+           random.color=TRUE,colors=rainbow(7))
```

# wordcloud pictorial

## clustering the words

We now try to find clusters of words with hierarchical clustering.

Sparse terms are removed, so that the plot of clustering will not be crowded with words.

Then the distances between terms are calculated with dist() after scaling.

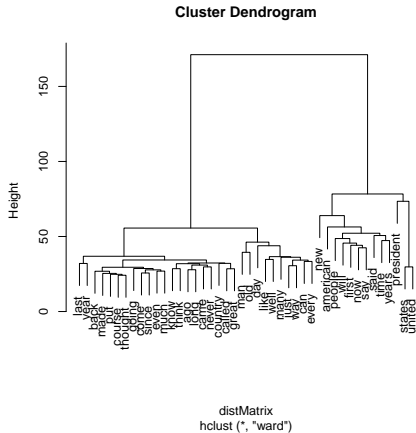After that, the terms are clustered with hclust() and the dendrogram is cut into 10 clusters.

The agglomeration method is set to ward, which denotes the increase in variance when two clusters are merged.

Some other options are single linkage, complete linkage, average linkage, median and centroid.

## word clustering code

```
> # remove sparse terms
> tdmat <- as.matrix(
+    removeSparseTerms(tdm, sparse=0.3)
+ )
> # compute distances
> distMatrix <- dist(scale(tdmat))
> fit <- hclust(distMatrix, method="ward")
> plot(fit)
```

# word clustering dendogram



**Cluster Dendrogram**

distMatrix
hclust (*, "ward")

# clustering documents with k-medoids

We now try *k-medoids* clustering with the `Partitioning Around Medoids` algorithm.

In the following example, we use function `pamk()` from package `fpc` [Hennig, 2010], which calls the function `pam()` with the number of clusters estimated by optimum average silhouette.
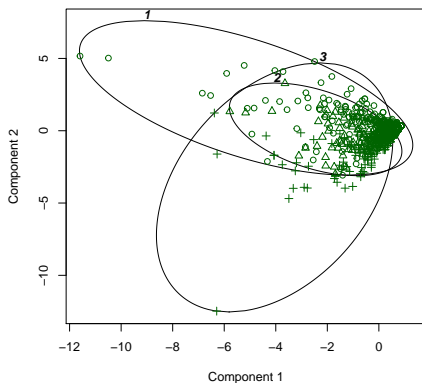
Note that because we are now clustering documents rather than words we must first transpose the term-document matrix to a document-term matrix.

# k-medoid clustering

```
> # first select terms corresponding to presidents
> pn <- c("nixon","carter","reagan","clinton","roose
+         "truman","bush","ford")
> # and transpose the reduced term-document matrix
> dtm <- t(as.matrix(tdm[pn,]))
> # find clusters
> library(fpc)
> pamResult <- pamk(dtm, krange=2:6,
+                     metric = "manhattan")
> # number of clusters identified
> (k <- pamResult$nc)
[1] 3
> # extract a pam object and plot it
> pamObj <- pamResult$pamobject
> # produce cluster plot
> clusplot(pamObj,labels=4, col.clus=1)
```

# document clustering plot



clusplot(pam(x = sdata, k = k, diss = diss, metric = "manhattan'

Component 1
These two components explain 32.11 % of the point variability.

## generate cluster wordclouds

```
> layout(matrix(c(1,2,3),1,3))
> for(k in 1:3){
+ cl <- which( pamObj$clustering == k )
+ tdmk <- t(dtm[cl,])
+ v = sort(rowSums(tdmk), decreasing=TRUE)
+ d = data.frame(word=names(v), freq=v)
+ # generate the wordcloud
+ wordcloud(d$word, d$freq, min.freq=5,
+           random.color=TRUE,colors="black")
+ }
> layout(matrix(1))
```

# cluster wordclouds

# exercises

No exercises this week!

make sure you submit all your outstanding assignments and
check you assignment marks on student central

## before 6h00 Monday 12th May