# Wk2_SLP1_Nikki Fitzherbert

Nikki Fitzherbert 13848336

05 March 2021

Note: The code was copied from https://www.analytics-tuts.com/zipfs-law-introduction-text-analytics/

## 1. Install and Load Packages

```
library(tm)
library(SnowballC)
library(wordcloud)
library(RColorBrewer)
```

## 2. Load, Transform and Clean Data

```
filePath <- "https://archive.org/stream/AnneFrankTheDiaryOfAYoungGirl_201606/
Anne-Frank-The-Diary-Of-A-Young-Girl_djvu.txt"

text <- readLines(filePath)
docs <- Corpus(VectorSource(text))

toSpace <- content_transformer(function (x , pattern ) gsub(pattern, " ", x))
docs <- tm_map(docs, toSpace, "/")
docs <- tm_map(docs, toSpace, "@")
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, stripWhitespace)
```
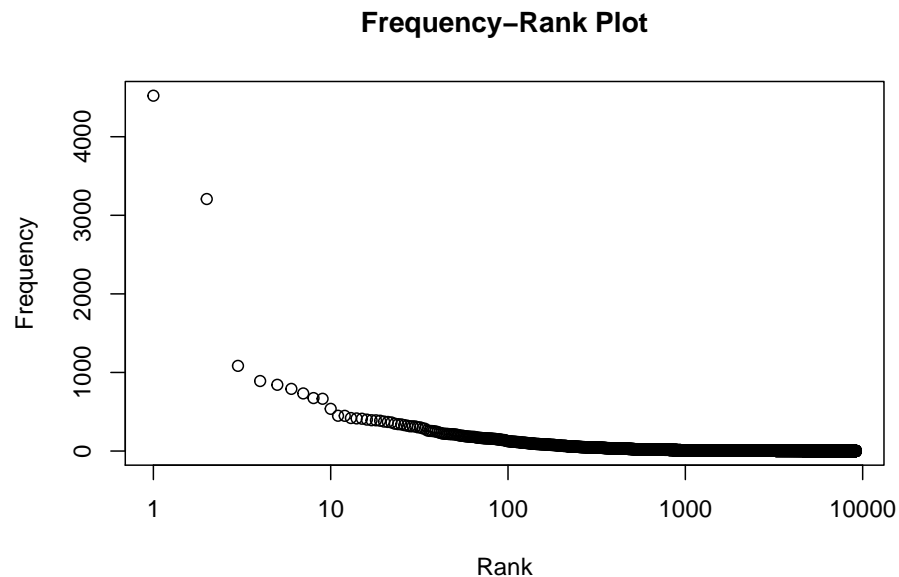
```
dtm <- TermDocumentMatrix(docs)

m <- as.matrix(dtm)
v <- sort(rowSums(m),decreasing=TRUE)
d <- data.frame(word = names(v),freq=v)

zipf <- cbind(d, Rank = 1:nrow(d), per=100 * d$freq / sum(d$freq))
```
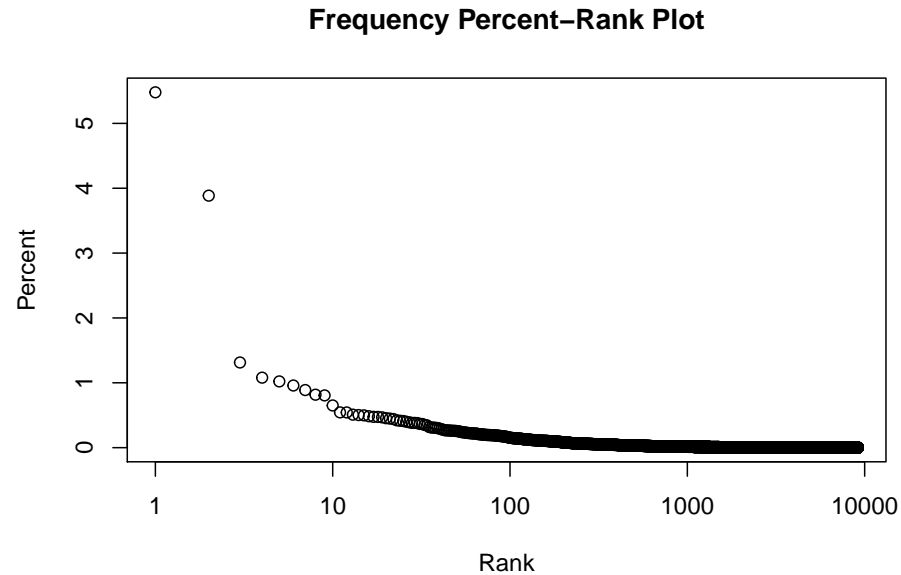
### 3. Visualising Zipf's Law

```
plot(zipf$Rank, zipf$freq,
     xlab = "Rank", ylab = "Frequency", main = "Frequency-Rank Plot",
     log = "x")
```

**Frequency–Rank Plot**

```
plot(zipf$Rank, zipf$per,
     xlab = "Rank", ylab = "Percent", main = "Frequency Percent-Rank Plot",
     log = "x")
```

**Frequency Percent–Rank Plot**



```
wordcloud(words = zipf$word, freq = d$freq, min.freq = 1, max.words=200,
random.order=FALSE, rot.per=0.35, colors=brewer.pal(8, "Dark2"))
```