## Project\_3

## Robert Sellers March 15, 2016

Libraries used.

```
library(data.table)
## Warning: package 'data.table' was built under R version 3.2.4
library(dplyr)
## Attaching package: 'dplyr'
## The following objects are masked from 'package:data.table':
##
       between, last
##
## The following objects are masked from 'package:stats':
##
##
       filter, lag
## The following objects are masked from 'package:base':
##
##
       intersect, setdiff, setequal, union
library(RCurl)
## Loading required package: bitops
library(stringr)
library(twitteR)
## Warning: package 'twitteR' was built under R version 3.2.4
##
## Attaching package: 'twitteR'
## The following objects are masked from 'package:dplyr':
##
       id, location
require(wordcloud)
## Loading required package: wordcloud
```

```
## Warning: package 'wordcloud' was built under R version 3.2.4

## Loading required package: RColorBrewer

require(tm)

## Loading required package: tm

## Loading required package: NLP
```

Twitter OAuth. Will not be evaluated. This is only required when generating new datasets. Both included, only one needed.

```
#Robert's Credentials
#consumer_key <- "DaAA9z8QvnxsdLOSpIr1oYwvP"
#consumer_secret <- "bCfsuODQyoYKMxPoHhZy2LxvvVqBvSM1LemzBtqm6YFeylWKUE"
#access_token <- "558596891-tDxN7T34cyVJJaBc4ExGTAq6wRfFBBlyHb2IzQvM"
#access_secret <-"nHqi5sVT2XRoCoSuOdYQnboCg1h35w5hRvtg657t8ROX8"
#setup_twitter_oauth(consumer_key, consumer_secret, access_token, access_secret)

#Chirag's credentials
#consumer_key <- "DaAA9z8QvnxsdLOSpIr1oYwvP"
#consumer_secret <- "bCfsuODQyoYKMxPoHhZy2LxvvVqBvSM1LemzBtqm6YFeylWKUE"
#access_token <- "558596891-tDxN7T34cyVJJaBc4ExGTAq6wRfFBBlyHb2IzQvM"
#access_secret <-"nHqi5sVT2XRoCoSuOdYQnboCg1h35w5hRvtg657t8ROX8"
```

Sample data mining code. The following code was used to generate the .csv files that you will find inside: https://github.com/RobertSellers/SlackProjects/tree/master/data

```
twitter_results_Feb_March_19<-searchTwitter("#datascience", n=10000)
twitter_results_Feb_March_19 <- Map(as.data.frame, twitter_results_Feb_March_19 )
twitter_results_Feb_March_19 <- rbindlist(twitter_results_Feb_March_19 )
write.csv(twitter_results_Feb_March_19 , file = "C:/Users/Robert/Desktop/CUNY/GitHub/R/data/twitter_results_Feb_March_19 )</pre>
```

Loading the data sources. Currently dating from 3/16, 3/18, 3/19. Will run only 3/16 for this.

```
twitter_results_march_16<-read.csv(file="https://raw.githubusercontent.com/RobertSellers/SlackProjects/states/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slackProjects/slac
```

twitter\_results\_march\_19<-read.csv(file="https://raw.githubusercontent.com/RobertSellers/SlackProjects/

Word Cloud function.

To Do: We may want a second "stopWords" variable input.

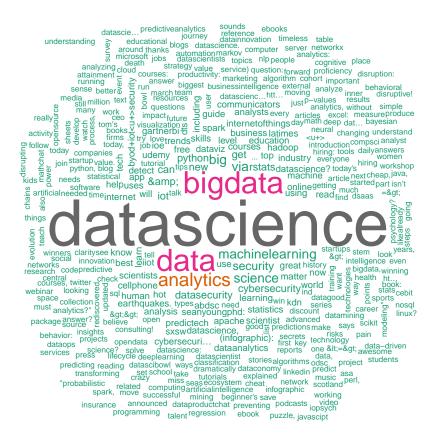
To Do: The results from this word cloud (that are relevant) should be added to the skills.csv.

To Do: Find a way to export a frequency table from this data?

```
dataScienceWordCloud<-function(twitterData){</pre>
  text<-twitterData$text
  corpus<- Corpus(VectorSource(text))</pre>
  corpus <- tm_map(corpus, stripWhitespace)</pre>
  corpus <- tm_map(corpus, content_transformer(tolower))</pre>
  corpus <- tm map(corpus, removeNumbers)</pre>
  toSpace <- content_transformer(function (x , pattern ) gsub(pattern, " ", x))
  corpus <- tm map(corpus, toSpace, "http\\S+\\s*")</pre>
  corpus <- tm_map(corpus, toSpace, "http\\S+\\s*") #twice????</pre>
  corpus <- tm_map(corpus, toSpace, "#")</pre>
  corpus <- tm_map(corpus, toSpace, "http\\w+")</pre>
  corpus <- tm_map(corpus, toSpace, "@//w+")
  corpus <- tm_map(corpus, toSpace, "https\\w+")</pre>
  corpus <- tm_map(corpus, toSpace, "uselect")</pre>
  corpus <-tm_map(corpus, removeWords, c(stopwords("english"),"#"))</pre>
  wordcloud(corpus, random.order=F, min.freq=1, max.words=400, colors=brewer.pal(8, "Dark2"))
}
```

## Running the word cloud on March 16

dataScienceWordCloud(twitter\_results\_march\_16)



Lookup table function.

ToDo: This uses the "skill.csv" as a look up table that is also applied to the other team's URL method. We ought to add new values to that csv based on what we find in the twitter dataset.

```
lutSkills<-read.csv(file= "https://raw.githubusercontent.com/RobertSellers/SlackProjects/master/data/sk
lookupFrequencies <-function(twitterData,lookupTable) {
   lookupTable<-as.data.frame(lookupTable)
   lookupTable$Skill<-paste0("\\<",lookupTable$Skill,"\\>")
   lookupTable["counts"]<-NA
   lookupTable$Skill[7] <- "xxxxxxxxxxx" #C++ not working as a keyword
   i<-1
   for(i in 1:nrow(lookupTable)) {
        lookupTable$counts[i]<-length(grep(lookupTable$Skill[i], twitterData$text))
   }
   colfunc<-colorRampPalette(c("red","yellow","springgreen","royalblue"))
   colnamesbarplot <- as.character(lookupTable$Skill)
   barplot(lookupTable$counts,main="Word Counts",horiz=TRUE,col=(colfunc(50)),axes=TRUE, names.arg=colnargrid(nx=NULL, ny=NA,col="black")
box()
}</pre>
```

Running the function on March 16th

To Do - Fix plot & refine the function. Sort the data and continue to update the skills lookup table to get better data.

To Do - We ultimately want to run this on each date to ensure that we have consistency between dates and then ultimately to combine all of the data.

```
lookupFrequencies(twitter_results_march_16,lutSkills)
```

## **Word Counts**

