

# Analysis of Twitter Activity and Sentiment on Starbucks

*Zihuan Qiao Teammate: CJ Xiang*

## I. Introduction

In this generation of social media, almost everyone has a social media account to interact with each other almost everyday. People express their ideas freely on any topics, from politics to what they eat today. There are also a series of functions enabling different kinds of interaction among users, like replying, forwarding, that enlarges the impact of social media on users. For its huge amount of users and highly interactive property, social media data plays an increasingly significant role in business analysis and financial analysis nowadays.

Twitter, one of the most popular social medias now, has its API available to the public which makes data analysis taking advantage of social media data easier than ever. By using twitter API, one can connect to twitter database to get data of a specific location during a specific time on a specific topic.

In this work, we try to use Starbucks twitter data to see the twitter activity and sentiment on Starbucks. Starbucks is an American coffee company and coffeehouse chain that has very high popularity. There are over 13,107 Starbucks in the United States. It is meaningful to see whether this popular and important coffeehouse chain company has good reputation on the social media. This can be helpful for many business purposes.

This project will be arranged as following: we will first look at the number of twitter on Starbucks compared with Dunkin' Donuts and their distribution in the US on a map to see which one is popular as a topic on social media. Then, we want to show what are people talking about Starbucks when mentioning it on Twitter by using word clouds. Further analysis is related to sentiment analysis. We use ordinal logistic regression to see the relationship between location and sentiment to Starbucks (positive, neutral or negative). Taking different level of influence of each tweets has on other users, we also conduct analysis on the variables that can reflect the influence level. We use kernel density estimation to get plots of density of the variables and use bootstrap to estimate their means.

## II. Method

### 1. Data

In addition to Starbucks, we also collect Dunkin' Donuts twitter data for comparison. Dunkin' Donuts is another American coffeehouse chain. The company has grown to become one of the largest coffee and baked goods chains in the world. In order to compare their popularity on Twitter, we set all the other parameters to be the same except for the searching topic.

```
## Collect data from Twitter

### Corresponding roauth R code and stream R code can be found
###in the separate R code files: roauth.R, stream_Starbucks.R, stream_Dunkin' Donuts.R
### Here starts from reading the saved RDS data

dunkinUS.df <- readRDS("Dunkin' Donuts US Data.RDS")
starbucksUS.df <- readRDS("Starbucks US Data.RDS")
dim(dunkinUS.df)

## [1] 6 42
```

```
dim(starbucksUS.df)
```

```
## [1] 210 42
```

In the same amount of time, 6000 seconds, we collect 210 tweets about Starbucks in the US and only 6 tweets about Dunkin' Donuts in the US. Each dataset has 42 variables. Huge amount of difference in number of tweets indicate that Starbucks is much more popular as a topic on twitter.

```
## draw maps of twitter activity distribution
```

```
library(ggplot2)
```

```
#draw map of Starbucks tweets in the USA domain
```

```
map.data <- map_data("state")
```

```
points <- data.frame(x=as.numeric(starbucksUS.df$place_lon),  
                    y=as.numeric(starbucksUS.df$place_lat))
```

```
points <- points[points$y>25,]
```

```
ggplot(map.data)+
```

```
  geom_map(aes(map_id = region),
```

```
            map=map.data,
```

```
            fill="white",
```

```
            color="grey20",size=0.25)+
```

```
  expand_limits(x = map.data$long, y = map.data$lat)+
```

```
  theme(axis.line = element_blank(),
```

```
        axis.text = element_blank(),
```

```
        axis.ticks = element_blank(),
```

```
        axis.title = element_blank(),
```

```
        panel.background = element_blank(),
```

```
        panel.border = element_blank(),
```

```
        panel.grid.major = element_blank(),
```

```
        plot.background = element_blank(),
```

```
        plot.margin = unit(0 * c(-1.5, -1.5, -1.5, -1.5), "lines"))+ 
```

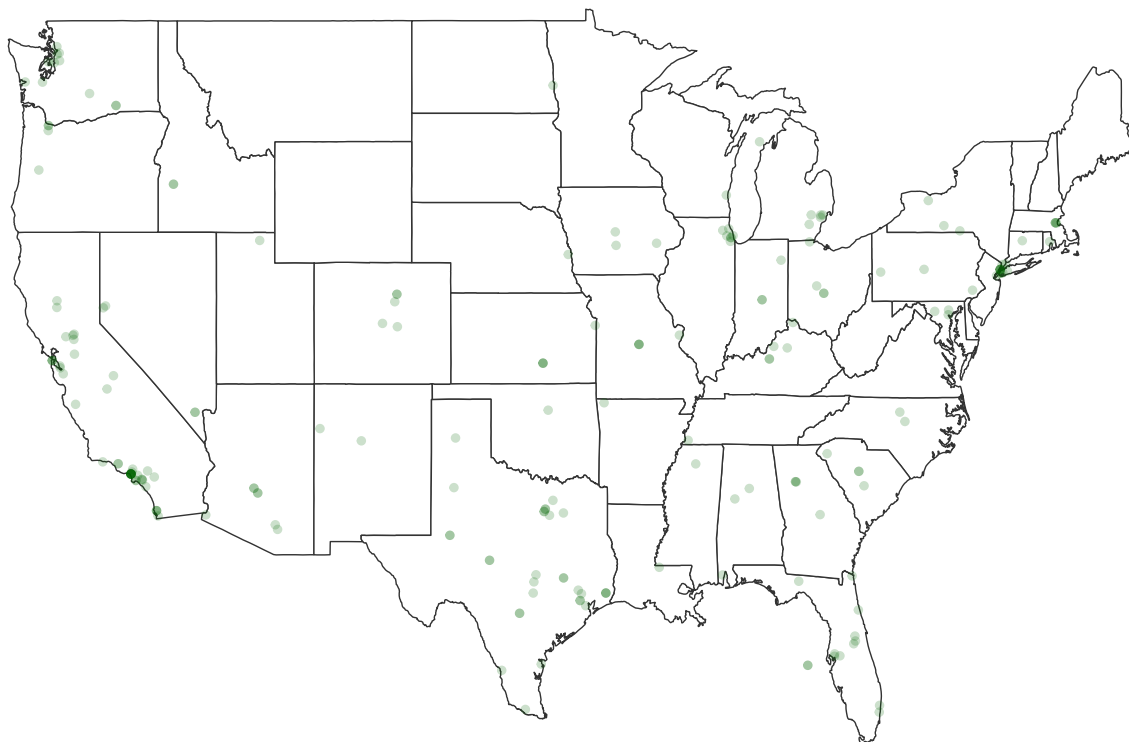
```
  geom_point(data = points,
```

```
             aes(x = x, y = y), size = 1,
```

```
             alpha = 1/5, color = "darkgreen")+ 
```

```
  ggtitle("Tweets mentioning Starbucks in the U.S.")
```

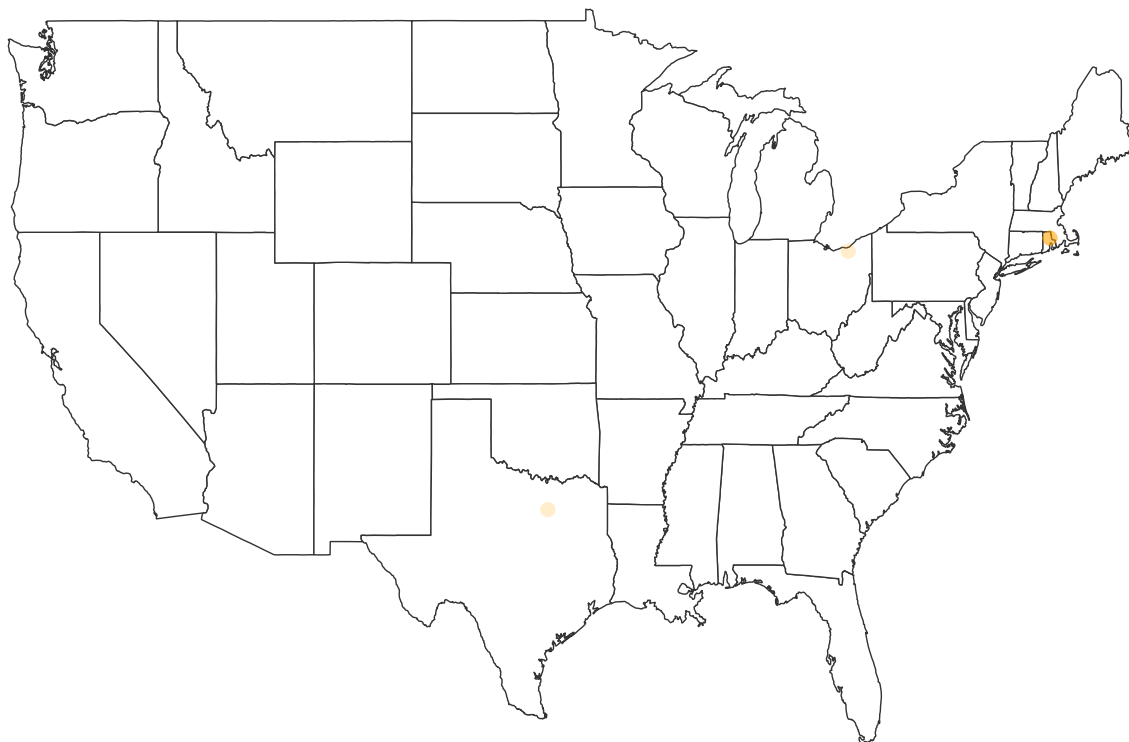
## Tweets mentioning Starbucks in the U.S.



```
#draw map of Dunkin' Donuts tweets in the USA domain
map.data <- map_data("state")
points <- data.frame(x=as.numeric(dunkinUS.df$place_lon),
                     y=as.numeric(dunkinUS.df$place_lat))

points <- points[points$y>25,]
ggplot(map.data)+
  geom_map(aes(map_id=region),
           map=map.data,
           fill="white",
           color="grey20",size=0.25)+
  expand_limits(x=map.data$long,y=map.data$lat)+
  theme(axis.line=element_blank(),
        axis.text=element_blank(),
        axis.ticks=element_blank(),
        axis.title=element_blank(),
        panel.background=element_blank(),
        panel.border=element_blank(),
        panel.grid.major=element_blank(),
        plot.background=element_blank(),
        plot.margin=unit(0*c(-1.5,-1.5,-1.5,-1.5),"lines"))+
  geom_point(data=points,
            aes(x=x,y=y),size=2,
            alpha=1/5,color="orange")+
  ggtitle("Tweets Mentioning Dunkin' Donuts in USA")
```

## Tweets Mentioning Dunkin' Donuts in USA



Two maps above show the twitter activity distribution in the country. Each point represents a tweet. We can see that tweets on Starbucks distribute almost evenly in this US, except for the Midwest region. There are barely points in this region which indicates topic Starbucks is not hot there. While in terms of Dunkin' donuts, since there are only six data in the US, so the points on the map are also very sparse. Tweets on Dunkin' Donuts are far less active than Starbucks.

## 2. Variable

There are altogether 42 variables in the original data from Twitter, including text, followers count, favourites count, name, tweet time, country, place longitude, place latitude and so on. But we are not going to use all of them. In this work, we are only going to focus on six of them, they are text, followers count, favourites count, full name, place latitude and place longitude.

Table 1: Names and descriptions of Variables

Variable Name	Description and Variable Labels
text	tweets text
followers count	number of followers of user
favourites count	number of favourites of tweet
full name	user location full name, including city name and state abbreviation
place latitude	user location latitude
place longitude	user location longitude

### 3. Text Mining

#### a. Word Frequency

In order to find what twitter users are talking about when mentioning Starbucks, we use text mining technique to deal with text variable in the data which are tweets contents. In R, the main package to perform text mining is tm package. In addition to that, we also use wordcloud package along with RColorBrewer package to visualize the word frequency. Here we draw two word clouds, the first one corresponds to the Starbucks tweets content, the second one is about the hashtag frequency.

```
## text analysis: wordclouds
```

```
# Let's import necessary packages needed for generating a wordcloud  
library(tm)
```

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

```
##  
## Attaching package: 'NLP'
```

```
## The following object is masked from 'package:ggplot2':  
##  
##      annotate
```

```
library(wordcloud)
```

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

```
library(RColorBrewer)  
library(stringr)  
  
# remove Emoji and wiered characters  
Star_text <- sapply(starbucksUS.df$text, function(row) iconv(row, "latin1", "ASCII", sub=""))  
# create a corpus  
Star_corpus = Corpus(VectorSource(Star_text))  
# create document term matrix applying some transformations  
tdm = TermDocumentMatrix(Star_corpus,  
  control = list(removePunctuation = TRUE,  
  stopwords = c("Starbucks", stopwords("english")),  
  removeNumbers = TRUE, tolower = TRUE))  
  
# define tdm as matrix  
m = as.matrix(tdm)  
# get word counts in decreasing order  
word_freqs = sort(rowSums(m), decreasing=TRUE)  
# create a data frame with words and their frequencies  
dm = data.frame(word=names(word_freqs), freq=word_freqs)  
# plot wordcloud  
wordcloud(dm$word, dm$freq, random.order=FALSE, colors=brewer.pal(8, "Dark2"))
```



From the second word cloud, which is the word cloud of hashtags, it's even clearer that hashtags like jobs, career, hiring still are among the tops in frequency. It again indicates that jobs at Starbucks is a hot topic in twitter to some extent.

## b. Sentiment Analysis

Sentiment analysis is another important part of conducting text mining. Inspired by the results from word frequency part that some words related to service quality are also frequently mentioned in tweets, and that many of these words are positive words, we further perform sentiment analysis on the text variable. In this part, we try to find twitter users' attitude towards Starbucks.

In this project, we perform sentiment analysis in R using `syuzhet` package. To get the sentiment of each tweet, we apply the `get_nrc_sentiment` command. The `get_nrc_sentiment` implements Saif Mohammad's NRC Emotion lexicon. According to Mohammad, "the NRC emotion lexicon is a list of words and their associations with eight emotions (anger, fear, anticipation, trust, surprise, sadness, joy, and disgust) and two sentiments (negative and positive)".

```
## sentiment analysis

# required packages
library(dplyr)
library(ggplot2)
library(syuzhet)
library(plotrix)

# extract text
Star_text <- starbucksUS.df$text

# clean text
# We try to get rid of Emoji and wieried characters
Star_text <- sapply(starbucksUS.df$text,
  function(row) iconv(row, "latin1", "ASCII", sub=""))

# remove retweet entities
Star_text = gsub("(RT|via)((?:\\b\\W*@[\\w+)+)", "", Star_text)
# remove at people
Star_text = gsub("@\\w+", "", Star_text)
# remove punctuation
Star_text = gsub("[[:punct:]]", "", Star_text)
# remove numbers
Star_text = gsub("[[:digit:]]", "", Star_text)
# remove html links
Star_text = gsub("http\\w+", "", Star_text)
# remove unnecessary spaces
Star_text = gsub("[ \\t]{2,}", "", Star_text)
Star_text = gsub("^\\s+|\\s+$", "", Star_text)

# define "tolower error handling" function
try.error = function(x)
{
  # create missing value
  y = NA
  # tryCatch error
  try_error = tryCatch(tolower(x), error=function(e) e)
  # if not an error
  if (!inherits(try_error, "error"))
    y = tolower(x)
  # result
```

```

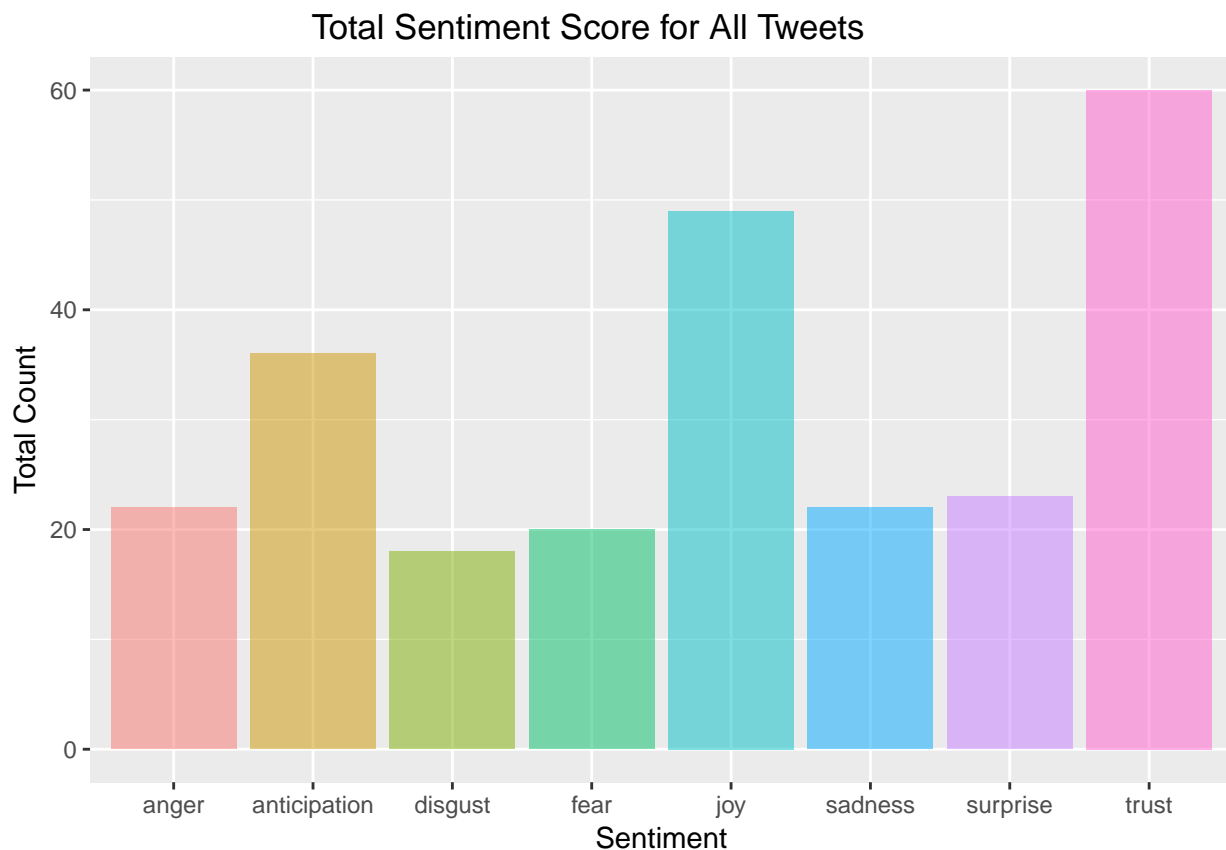
    return(y)
}

# lower case using try.error with sapply
Star_text = sapply(Star_text, try.error)

#extract sentiment
mySentiment <- get_nrc_sentiment(Star_text)

# plot sentiment
sentimentTotals <- data.frame(colSums(mySentiment[,c(1:8)]))
names(sentimentTotals) <- "count"
sentimentTotals <- cbind("sentiment" = rownames(sentimentTotals), sentimentTotals)
rownames(sentimentTotals) <- NULL
ggplot(data = sentimentTotals, aes(x = sentiment, y = count)) +
  geom_bar(aes(fill = sentiment), stat = "identity", alpha=0.5) +
  theme(legend.position = "none") +
  xlab("Sentiment") + ylab("Total Count") +
  ggtitle("Total Sentiment Score for All Tweets")

```



```

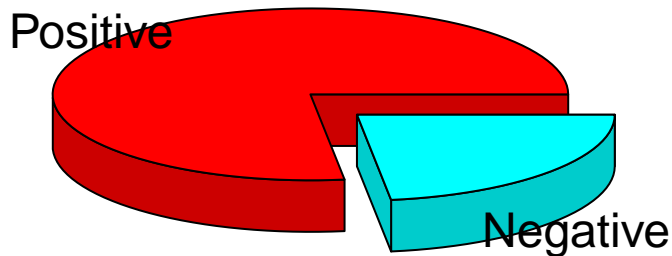
# Pie Chart with Percentages
pos <- sum(mySentiment$positive)
neg <- sum(mySentiment$negative)
slices <- c(pos, neg)
lbls <- c("Positive", "Negative")

```



```
pie3D(slices,labels=lbls,explode=0.12,
      main="Pie Chart of Postive and Negative Tweets")
```

## Pie Chart of Postive and Negative Tweets



```
# Combine clean data with sentiment data
starfull <- cbind(starbucksUS.df,mySentiment[,9:10])
```

‘Total Sentiment Score for all Tweets’ shows difference in number of eight emotions based on all the tweets text in the US. Trust, joy and anticipation are three highest frequency emotions. Especially bars for trust and joy are significantly higher than the other bars. Top three emotions are all positive emotions. Level of other five emotions including anger, disgust, fear, sadness and surprise are close to each other. And four out of these five emotions are negative emotions. We also notice that the difference between sadness and anticipation is only 5, which not very big actually. To conclude, there are more tweets user that hold a positive attitude towards Starbucks.

‘Pie Chart of Positive and Negative Tweets’ summarize the proportion of positive tweets and negative tweets. This pie chart clearly shows that positive tweets take up more than three quateriles of the total tweets. Thus, we can draw the same conclusion as from the bar chart.