

SENTIMENT ANALYSIS OF SOCIAL MEDIA CONTENTS

HIMANSHU RAWAT

INTRODUCTION

- In the past one decade, there has been an exponential surge in the online activity of people across the globe.
- The volume of posts that are made on the web every second runs into millions. To add to this, the rise of social media platforms has led to flooding to content on the internet.
- Social media is not just a platform where people talk to each other, but it has become very vast and serves many more purposes. It has become a medium where people
 - Express their interests.
 - Share their views.
 - Share their displeasures.
 - Compliment companies for good and poor services.
- In this project, we will be analyzing the sentiments that are reflected from the post of people on social media networks.
- For our study, we have chosen “Twitter” as our social media platform. And the coding language that we used to perform the analysis is R.

TWITTER

- An online social networking service that enables users to send and read short 140-character messages called “tweets” (Wikipedia).
- Over 320 million monthly active users (as of 2016).
- Creating over 500 million tweets per day



Techniques and Tools

Techniques

- Text mining
- Topic modelling
- Sentiment analysis
- Social network analysis

Tools

- Twitter API
- R and its packages:
 - twitterR
 - GGLOT2
 - ROAuth
 - sentiment
 - Word Cloud

PROCESS FLOW

Extract tweets and followers from the Twitter website with R and the twitterR package

With the tm package, clean text by removing punctuations, numbers, hyperlinks and stop words, followed by stemming and stem completion

Build a term-document matrix

Analysis of most frequently used words

Analysis of sentiment with the sentiment package

Generation of Word Cloud

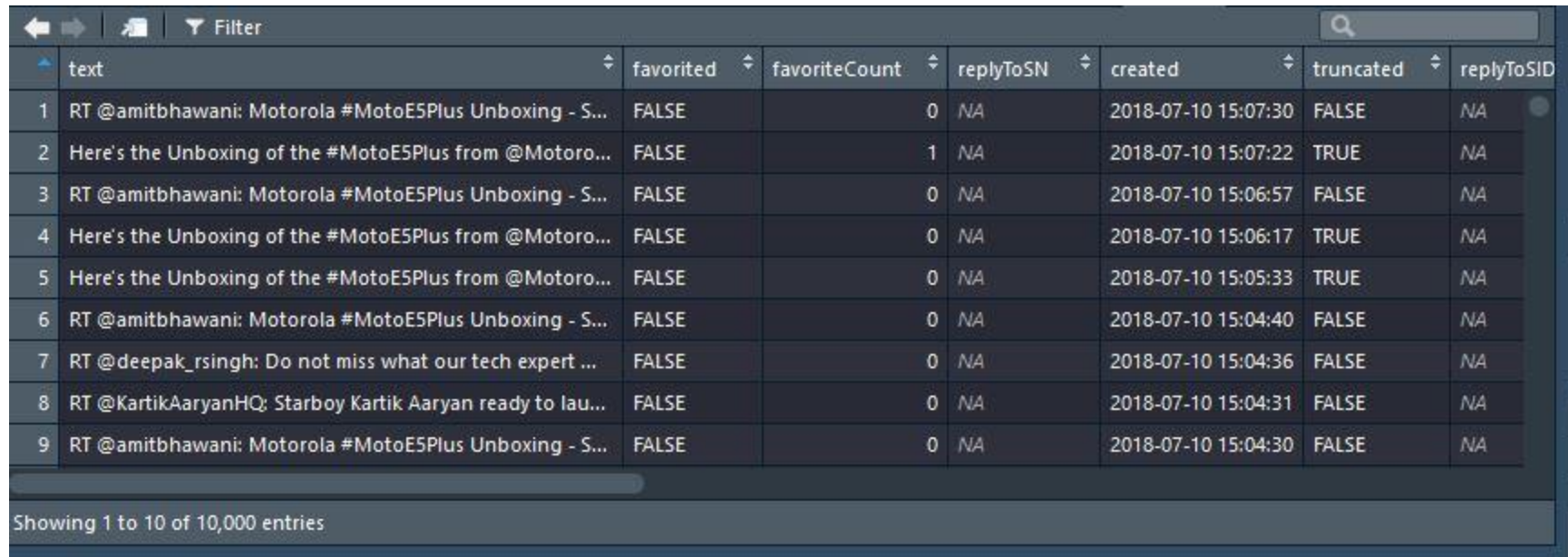


Functions & Generated Results

(Case Study on Moto E 5 Plus Launch)

Launched: 10 July-18

Retrieve Tweets



The screenshot shows a web-based interface for viewing a list of tweets. At the top, there is a navigation bar with a back arrow, a search icon, and a 'Filter' button. Below this is a table with 8 columns: 'text', 'favorited', 'favoriteCount', 'replyToSN', 'created', 'truncated', and 'replyToSID'. The table contains 10 rows of data, each representing a tweet. The first row is a retweet of a tweet about Motorola MotoE5Plus unboxing. The second row is the original tweet. The third row is another retweet. The fourth row is the original tweet again. The fifth row is another retweet. The sixth row is the original tweet. The seventh row is a tweet from @deepak_rsingh. The eighth row is a tweet from @KartikAaryanHQ. The ninth row is a retweet of the Motorola tweet. The bottom of the interface shows a status bar indicating 'Showing 1 to 10 of 10,000 entries'.

	text	favorited	favoriteCount	replyToSN	created	truncated	replyToSID
1	RT @amitbhawani: Motorola #MotoE5Plus Unboxing - S...	FALSE	0	NA	2018-07-10 15:07:30	FALSE	NA
2	Here's the Unboxing of the #MotoE5Plus from @Moto...	FALSE	1	NA	2018-07-10 15:07:22	TRUE	NA
3	RT @amitbhawani: Motorola #MotoE5Plus Unboxing - S...	FALSE	0	NA	2018-07-10 15:06:57	FALSE	NA
4	Here's the Unboxing of the #MotoE5Plus from @Moto...	FALSE	0	NA	2018-07-10 15:06:17	TRUE	NA
5	Here's the Unboxing of the #MotoE5Plus from @Moto...	FALSE	0	NA	2018-07-10 15:05:33	TRUE	NA
6	RT @amitbhawani: Motorola #MotoE5Plus Unboxing - S...	FALSE	0	NA	2018-07-10 15:04:40	FALSE	NA
7	RT @deepak_rsingh: Do not miss what our tech expert ...	FALSE	0	NA	2018-07-10 15:04:36	FALSE	NA
8	RT @KartikAaryanHQ: Starboy Kartik Aaryan ready to lau...	FALSE	0	NA	2018-07-10 15:04:31	FALSE	NA
9	RT @amitbhawani: Motorola #MotoE5Plus Unboxing - S...	FALSE	0	NA	2018-07-10 15:04:30	FALSE	NA

Showing 1 to 10 of 10,000 entries

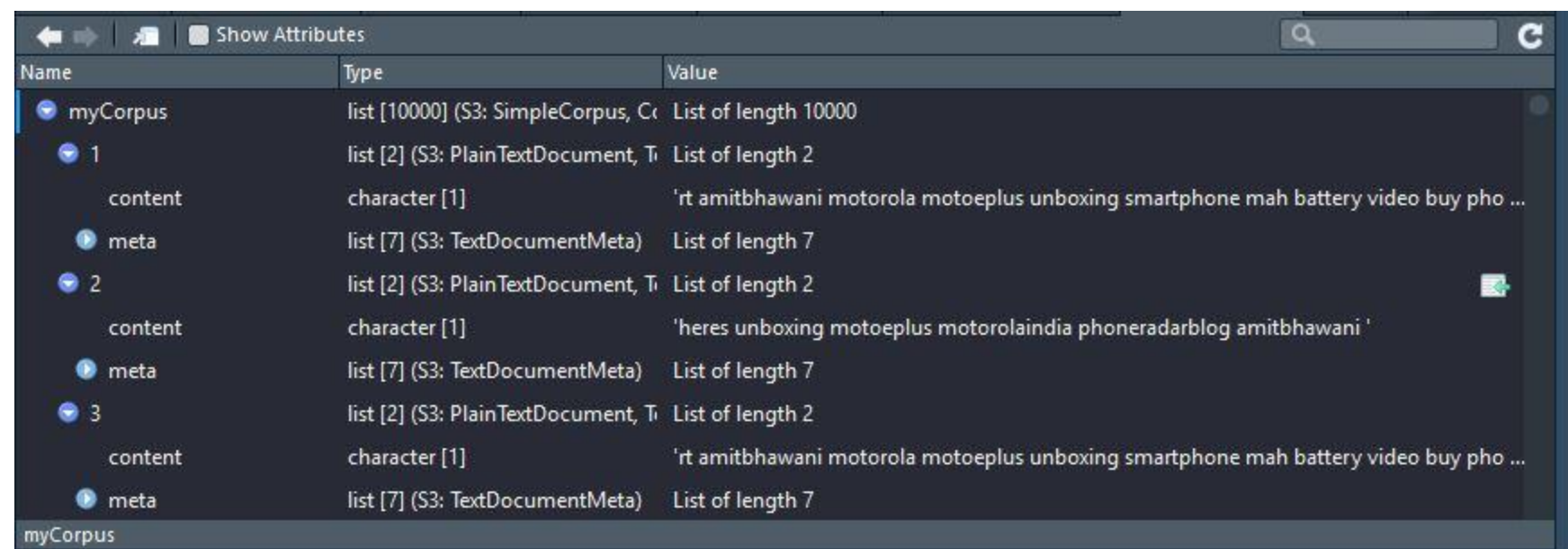
Following functions have been used:

- **twitteR**-for connecting with twitter rest api
- **ROAuth**- for authorization with twitter

*Tweets were and extracted and then converted to a data-frame as shown here.

*No of Tweets extracted =10,000

Text Cleaning



Name	Type	Value
myCorpus	list [10000] (S3: SimpleCorpus, Co	List of length 10000
1	list [2] (S3: PlainTextDocument, Ti	List of length 2
content	character [1]	'rt amitbhawani motorola motoeplus unboxing smartphone mah battery video buy pho ...
meta	list [7] (S3: TextDocumentMeta)	List of length 7
2	list [2] (S3: PlainTextDocument, Ti	List of length 2
content	character [1]	'heres unboxing motoeplus motorolaindia phoneradarblog amitbhawani '
meta	list [7] (S3: TextDocumentMeta)	List of length 7
3	list [2] (S3: PlainTextDocument, Ti	List of length 2
content	character [1]	'rt amitbhawani motorola motoeplus unboxing smartphone mah battery video buy pho ...
meta	list [7] (S3: TextDocumentMeta)	List of length 7

myCorpus

“tm” function is used to pre-process(cleaning, removing hashtags ,hyperlinks etc)

The following operations were performed in a step wise manner

Tweets were converted to lowercase

Removal of URL's

Removal of number punctuations

Removal of Stop Words

Removal of white spaces.

Stemming

```
42 # Stem & Stemming
43 myCorpus <- tm_map(myCorpus, stemDocument) # stem words
44 writeLines(strwrap(myCorpus[[190]]$content, 60))
45 ## r refer card data mine now provid link packag cran packag
46 ## mapred hadoop ad
47 stemCompletion2 <- function(x, dictionary) {
48   x <- unlist(strsplit(as.character(x), " "))
49   x <- x[x != ""]
50   x <- stemCompletion(x, dictionary=dictionary)
51   x <- paste(x, sep="", collapse=" ")
52   PlainTextDocument(stripwhitespace(x))
53 }
54 myCorpus <- lapply(myCorpus, stemCompletion2, dictionary=myCorpusCopy)
55 myCorpus <- Corpus(VectorSource(myCorpus))
56 writeLines(strwrap(myCorpus[[190]]$content, 60))
57
```

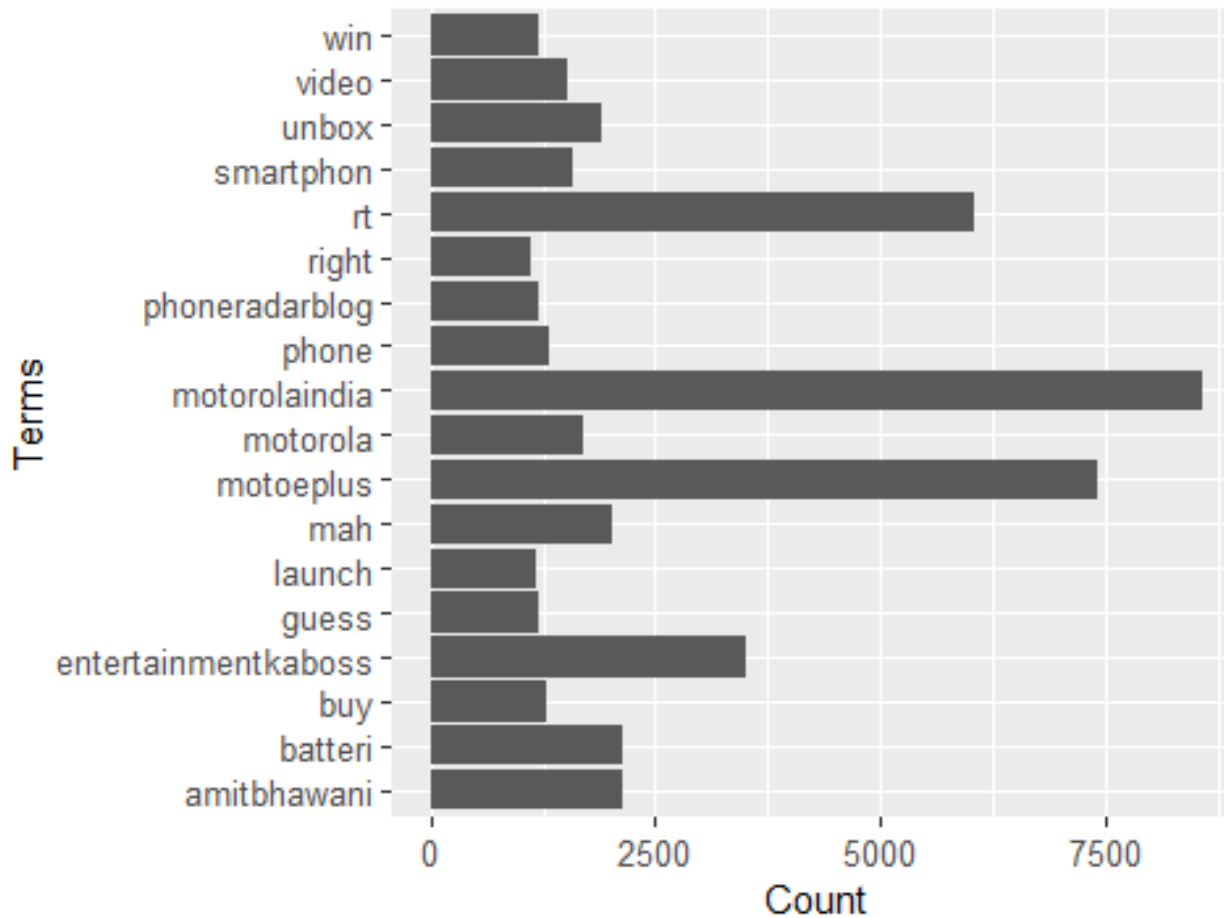
This is also a part of “tm” function and was used for stemming so that root words were consolidated and arbitrary words removed.

Term Doc Matrix

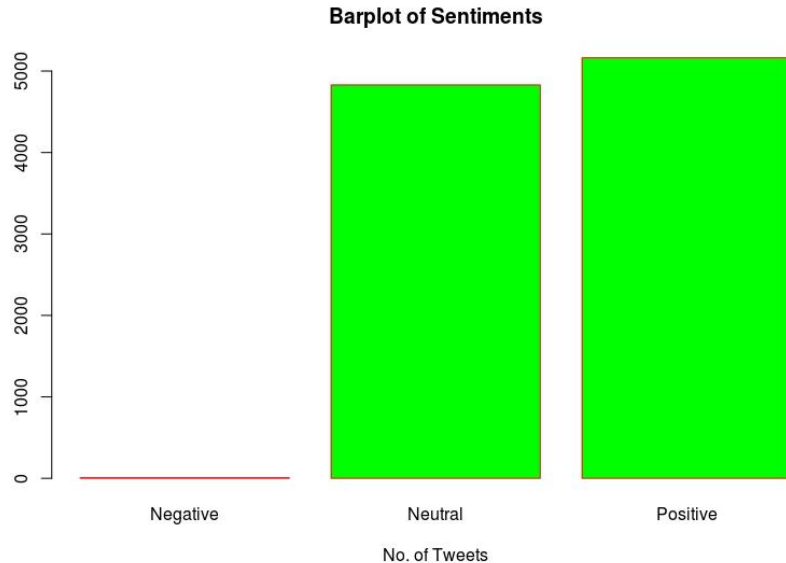
```
78 ##Build Term Document Matrix
79
80 tdm <- TermDocumentMatrix(myCorpus,
81                           control = list(wordLengths = c(1, Inf)))
82 tdm
83
84 idx <- which(dimnames(tdm)$Terms %in% c("r", "data", "mining"))
85 as.matrix(tdm[idx, 21:30])
86
```

Term doc matrix was created to obtain the frequency of words in the collection of tweets. This matrix is then used to generate the words frequency plot. The matrix is further employed to generate the sentimental plots.

Words frequency Plot



Sentiment Analysis- Maximum Entropy Classifier

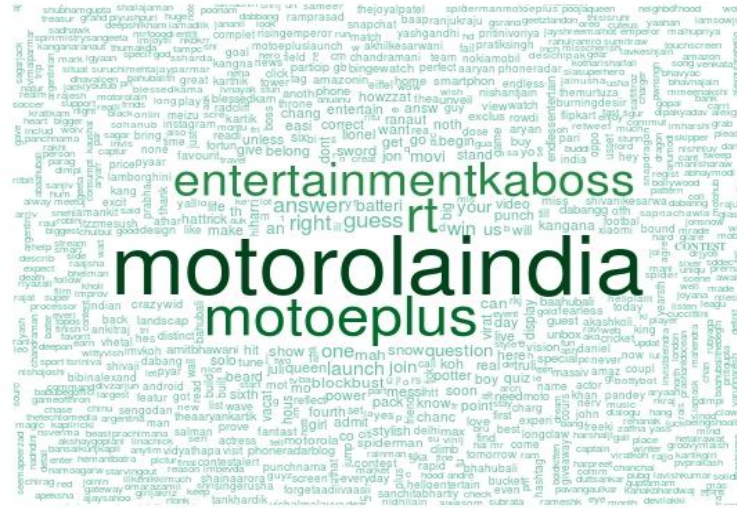


The independent words extracted from each tweet are matched with predefined lexicon. The lexicon is generated containing positive, negative & neutral sentiments using Maximum Entropy Classifier Approach. The cumulative score of each word gives the tweet a value(-1,0,+1), on the basis of which the tweet is defined as negative, neutral or positive. The count of each individual tweet is then taken to determine whether the sentiment is negative, neutral or positive for a particular subject (MotoE5 plus in our study).

STATISTICS:

Total Tweets	: 10,000	Hashtag	: Moto E5 Plus
Positive	:5163		
Negative	:7		
Neutral	:4830		

Word Cloud



A Word cloud is generated basis the frequency of words appearing in the tweets. The most frequently occurring words are displayed more prominently in the cloud of All the words present in the tweets. This is achieved by using “Word Cloud “ function .