

DIC

# PROJECT 1 : PROBLEM 5

## TWITTER DATA ANALYSIS USING R SHINY

NAME : ANIKET THIGALE

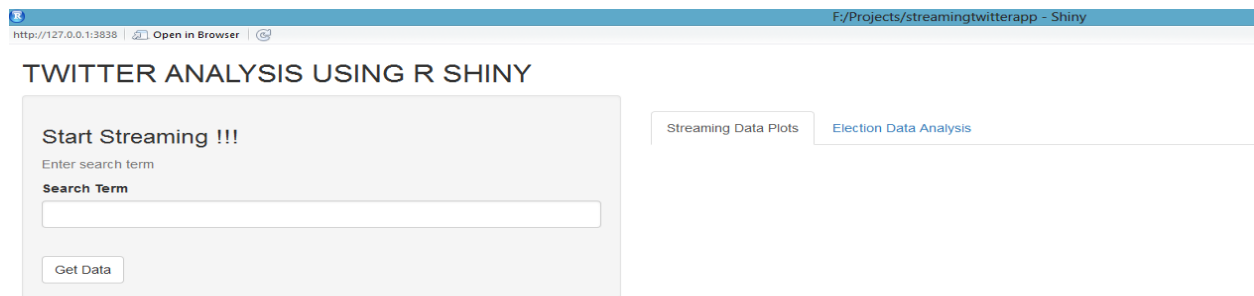
UB NO : 50168090

---

### PROBLEM DESCRIPTION

Write a R script to process and analyze streaming Twitter data about ongoing elections using R Shiny. Also collect tweets for a week from various regions of the country and plot the daily and weekly trend and summarize the stats.

### UI



The UI consists of a search field box for entering the term that you want to query twitter for, an action button when you want to start streaming and tabPanels for displaying some real time plots. It also has a tabPanel for displaying summary and some plots like trends for the data that we have initially collected for a week.

There are two components of this Problem:

1. Election Data Analysis
  2. Streaming Data to R Shiny
-

---

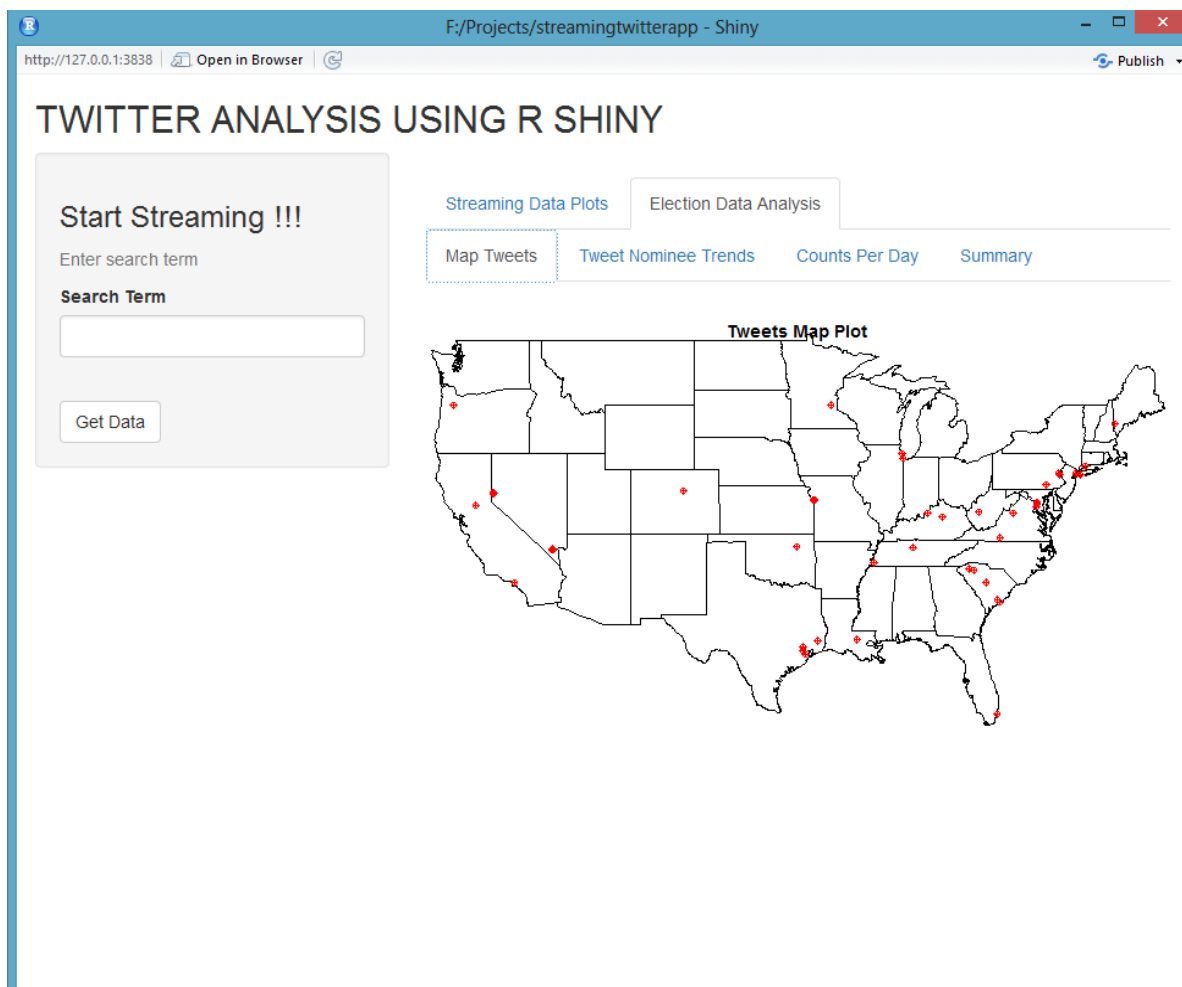
# ELECTION DATA ANALYSIS

## DATA

We use the election twitter data collected in problem 1 for analysis. This data has been collected over a week from February 20th 2016 to February 27th 2016 using various keywords like election2016, clinton, trump, cruz, bush etc.

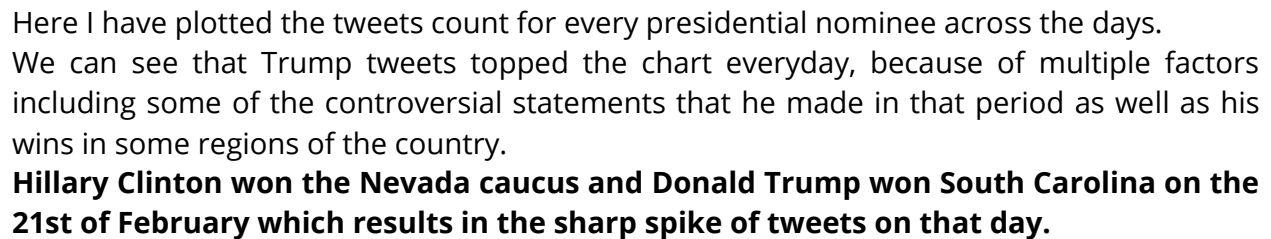
## PLOTS

### 1. MAP



We filter tweets which have latitudes and longitudes to display them on the map. According to this, there are more tweets coming from the East Coast, but since only a small number of tweets have latitude and longitude, this is very inaccurate.

Streaming Data Plots	Election Data Analysis	
Map Tweets	Tweet Nominee Trends	Counts Per Day
		Summary

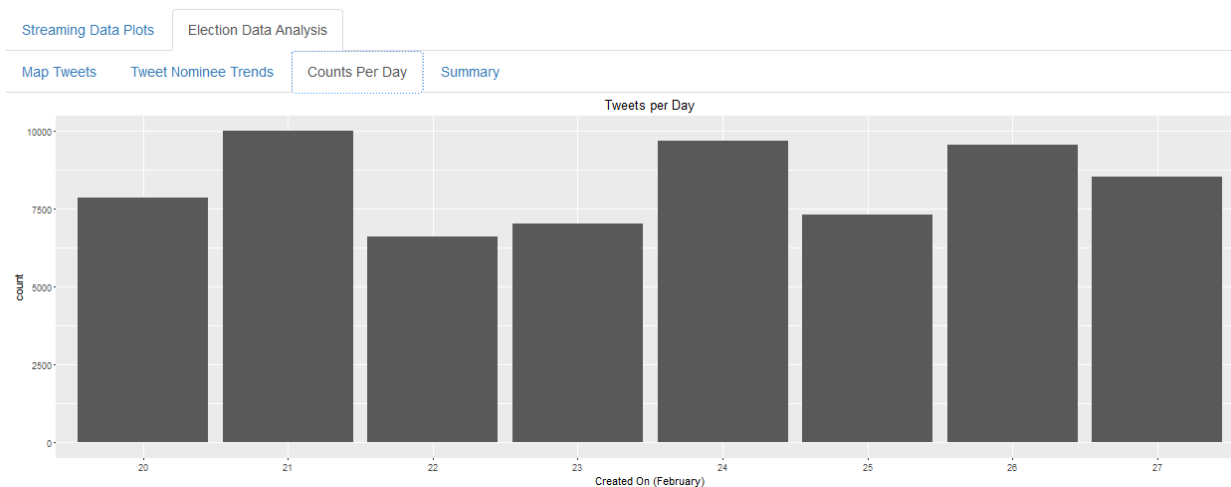


---

I created a wordcloud for all tweets in which Trump was mentioned by extracting text, cleaning it by removing punctuations, stop words and certain other text, stemming words, creating a corpus, constructing Term Document Frequency Matrix using the **tm**, **SnowballC** and **wordcloud** package available in R.

There were a large number of tweets about Trump's win in the South Carolina and Nevada elections and hence we can see the words "win", "southprimari", "nevada" as some of the top ones. Other keywords associated with Trump include "nevertrump", a hashtag which was trending, "racism", "absent" which refers to the GOP debate he didn't attend and so on.

## 4. TWEET VOLUME



Just gives the volume of tweets recorded for every day in the range of 20th February till 27th February.

## 5. SUMMARY

Streaming Data Plots		Election Data Analysis				
Map Tweets	Tweet Nominee Trends	Counts Per Day	Summary			
text	favorited	favoriteCount	created	truncated	id	statusSource
Length:66566	Mode :logical	Min. : 0.0000	Length:66566	Mode :logical	Length:66566	Length:66566
Class :character	FALSE:66566	1st Qu.: 0.0000	Class :character	FALSE:66566	Class :character	Class :character
Mode :character	NA's :0	Median : 0.0000	Mode :character	NA's :0	Mode :character	Mode :character
		Mean : 0.5982				
		3rd Qu.: 0.0000				
		Max. :1370.0000				
screenName	retweetCount	isRetweet	retweeted	replyToSN	replyToUID	replyToSID
Length:66566	Min. : 0.0	Mode :logical	Mode :logical	Length:66566	Length:66566	Length:66566
Class :character	1st Qu.: 0.0	FALSE:35738	FALSE:66566	Class :character	Class :character	Class :character
Mode :character	Median : 1.0	TRUE :30828	NA's :0	Mode :character	Mode :character	Mode :character
	Mean : 114.9	NA's :0				
	3rd Qu.: 21.0					
	Max. :1569.0					
longitude	latitude	created_day	tweet_date			
Length:66566	Length:66566	Min. :2016-02-20 00:00:01	Length:66566			
Class :character	Class :character	1st Qu.:2016-02-21 20:44:17	Class :character			
Mode :character	Mode :character	Median :2016-02-24 04:07:04	Mode :character			
		Mean :2016-02-24 02:16:16				
		3rd Qu.:2016-02-26 02:52:18				
		Max. :2016-02-27 23:59:56				

Gives us some more information about the tweet fields like count, mean and type of field.

## STREAMING DATA

We use the **streamR** package available for R to continuously stream tweets and write an observer block to detect changes.

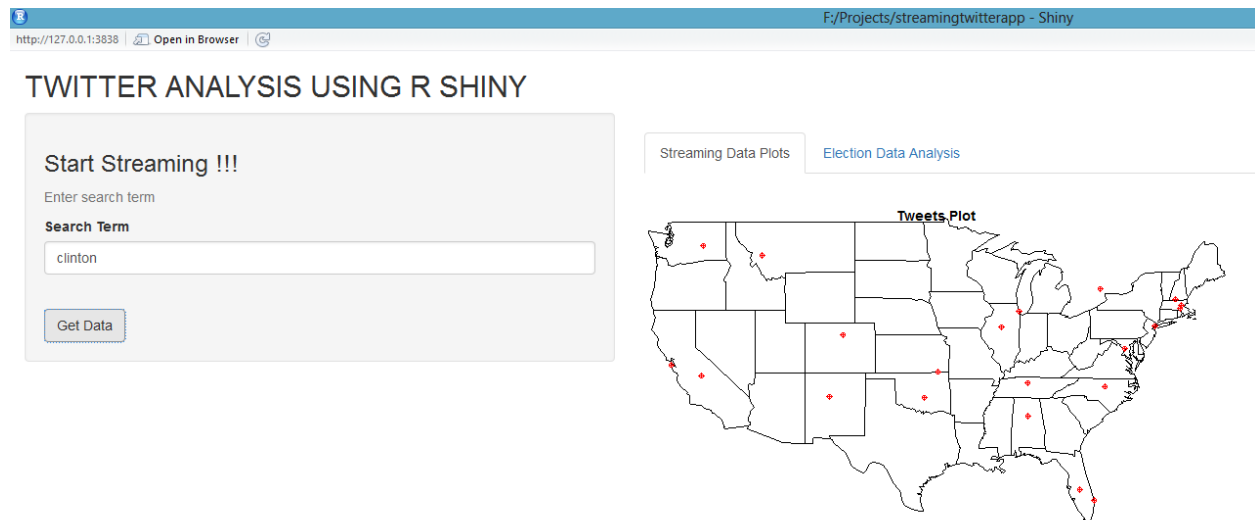
We use the `filterStream()` function in `streamR` to stream tweets which get written to a file in JSON format and use the `parseTweets()` function to retrieve the tweets from the JSON file into a dataframe in R.

Usage:

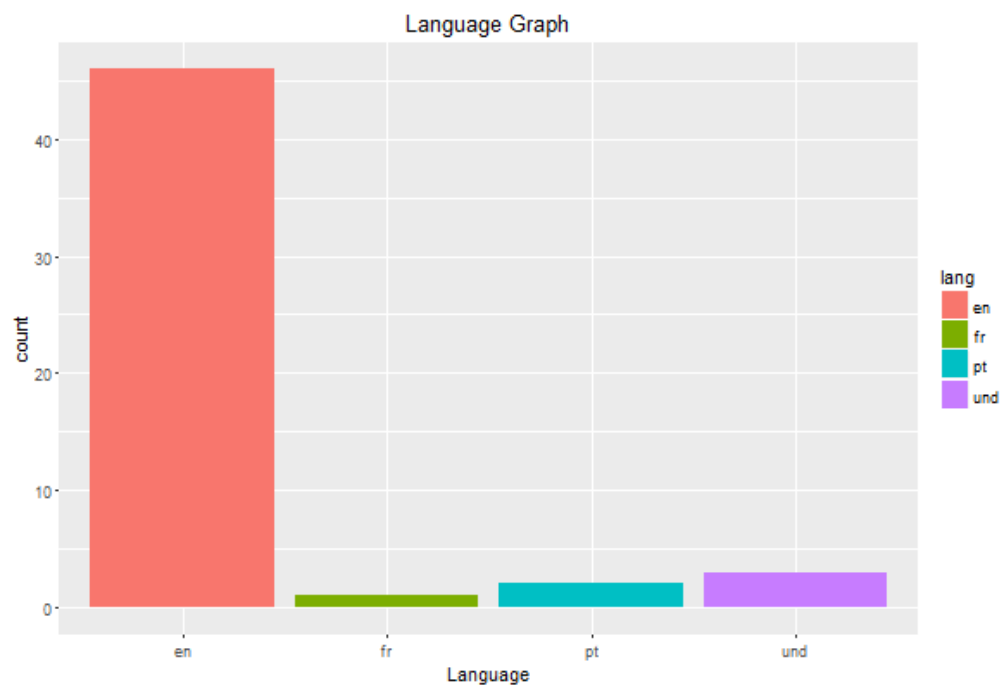
```
filterStream(file.name = "tweets.json", timeout=5, track = data, oauth = my_oauth)  
parseTweets("tweets.json", verbose = FALSE)
```

We get 43 fields for each tweet using the above API and we make use of some of them for generating the real time plots.

## PLOTS (All get updated in real time)

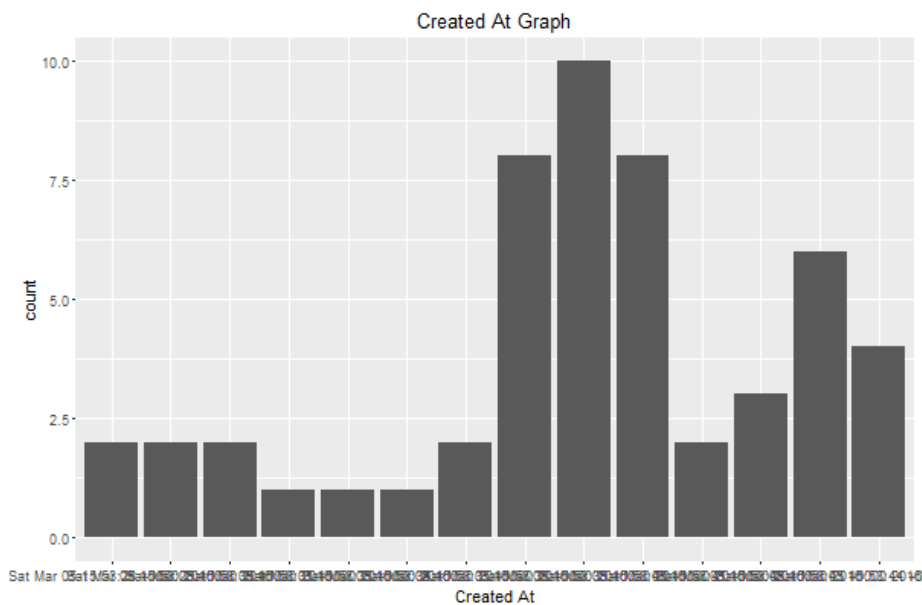


I used a map to display the locations from the streaming tweets. The red points indicate the locations which get updated in real time.



---

This language graph displays the tweets collected for each language which gets updated constantly.



This graph keeps track of count of new tweets against their created time.

## CONCLUSION

We have learnt to use many statistical, NLP and graph tools in R and use them to derive useful conclusions and display them using R Shiny.