

# CHAPTER 1

## INTRODUCTION

### 1.1 Introduction to Twitter Data Analytics

Social media like Twitter is a continuous source of real time information. The data present on social media can be used as knowledgeable information to serve various domains like depicting natural calamity (earthquakes etc.), finding the thoughts of people about any person or a particular news/event. Tweet text is a short text message limited by 140 characters in length posted by users on Twitter. People often post messages about interesting news, their daily activities, thoughts, feelings, as well as health conditions. Twitter<sup>®1</sup> is a massive social networking site tuned towards fast communication. More than 140 million active users publish over 400 million 140-character “Tweets” every day<sup>2</sup>. Twitter’s speed and ease of publication have made it an important communication medium for people from all walks of life. Twitter has played a prominent role in socio-political events, such as the Arab Spring<sup>3</sup> and the Occupy Wall Street movement<sup>4</sup>. Twitter has also been used to post damage reports and disaster preparedness information during large natural disasters, such as the Hurricane Sandy.

Twitter is a “micro-blogging” social networking website that has a large and rapidly growing user base. Thus, the website provides a rich bank of data in the form of “tweets,” which are short status updates and musings from Twitter’s users that must be written in 140 characters or less to tell others what they are doing, what they are thinking, or what is happening around them. According to the latest Twitter entry in Wikipedia, the number of Twitter users has climbed to 190 million and the number of tweets published on Twitter every day is over 65 million .As an increasingly popular platform for conveying opinions and thoughts, it seems natural to mine Twitter for potentially interesting trends regarding prominent topics in the news or popular culture.

Tweets and status updates range from important events to inane comments. Most messages contain little informational value but the aggregation of millions of messages can generate important knowledge. Several Twitter studies have demonstrated that aggregating millions of messages can provide valuable insights into a population. Barbosa and Feng (2010) classified tweets by sentiment, a first step towards measuring public opinion, such as political sentiment, which has been shown to track public political opinion and predict election results (Tumasjan et al. 2010; O’Connor et al. 2010). Eisenstein et al. (2010) studied lexical variations across geographic areas directly from tweets. Others have monitored the spread of news (Lerman and Ghosh 2010), detected the first mention of news events (Petrovic, Osborne, and Lavrenko 2010), and monitored earthquakes (Sakaki, Okazaki, and Matsuo 2010). Copyright © 2011, Association for the Advancement of Artificial Intelligence ([www.aaai.org](http://www.aaai.org)). All rights reserved. Twitter users often publicly express personal information; messages like “I got da flu” and “sick with this flu it’s taking over my body ughhhh” are common. Knowing that a specific user has the flu may not be interesting, but millions of such messages can be revealing, such as tracking the influenza rate in the United Kingdom and United States (Lampos and Cristianini 2010; Culotta 2010b). Furthermore, tweets are not isolated events: they occur with specific times, locations, languages and users. Aggregating over millions of users could provide new tools for research.

## 1.2 Rapid Miner

Rapid Miner is the most powerful, easy to use and intuitive graphical user interface for the design of analytic processes. Hundreds of data loading, data transformation, data modelling, and data visualization methods with access to a comprehensive list of data sources including Excel, Access, Oracle, IBM DB2, Microsoft SQL, Netezza, Teradata, MySQL, Postgres, SPSS, Salesforce.com, and hundreds more! Easily integrate your own specialized algorithms into Rapid Miner by leveraging its powerful and open extension APIs. Rapid Miner Studio runs on every major platform and operating system. With more and more employees bringing their own devices, you need flexibility in how and where your analytics are processed. Rapid Miner Studio breaks away from the limitations of traditional data analysis tools and allows you to work with large data sources.

RapidMiner (formerly YALE)

is the world-wide leading open-source data mining solution due to the combination of its leading-edge technologies and its functional range.

Applications of RapidMiner cover a wide range of realworld data mining tasks. The modular operator concept of RapidMiner (formerly YALE) allows the design of complex nested operator chains for a huge number of learning problems in a very fast and efficient way (rapid prototyping). The data handling is transparent to the operators. They do not have to cope with the actual data format or different data views - the RapidMiner core takes care of all necessary transformations.

## 1.3 Twitter 4J

Twitter4J is an unofficial Java library for the Twitter API. With Twitter4J, you can easily integrate your Java application with the Twitter service. Twitter4J is an unofficial library.

Twitter4J is featuring:

- ✓ 100% Pure Java - works on any Java Platform version 5 or later
- ✓ Android platform and Google App Engine ready
- ✓ Zero dependency : No additional jars required
- ✓ Built-in OAuth support
- ✓ Out-of-the-box gzip support
- ✓ 100% Twitter API 1.1 compatible

## 1.4 JxBrowser

**JxBrowser** is a cross-platform Java library that allows integrating Google Chromium-based web browser component into the Java Swing/AWT/JavaFX application. With JxBrowser we can embed a lightweight Swing/JavaFX component into our Java application to display modern web pages, supporting the latest web standards such as HTML5, CSS3, and JavaScript etc. JxBrowserPanel is a component that currently has the sole purpose of "host" for the JxBrowser. It is a class that inherits from JPanel and only has the BrowserView the JxBrowser inside. To use this component we need only jxbrowser.

## JXBROWSER FEATURES:

- **Rich and easy-to-use API.**
- **Great documentation** with many examples.
- **You're in control of your app.** Compared to other similar products, JxBrowser does not execute native code in your Java application process. Forget about memory usage issues caused by native web browser engine that allocates too much memory of your Java process. All the native code runs in separate native processes. In the unlikely event of a browser crash (e.g. because of a plugin problem), your app will continue working and you can restore the browser back.

**Great support.** Responsive support team will answer your request within 24 hours.

## 1.5 Google API

**Google APIs** is a set of APIs developed by Google which allow communication with Google Services and their integration to other services. Examples of these include Search, Gmail, Translate or Google Maps. Third-party apps can use these APIs to take advantage of or extend the functionality of the existing services.

The APIs provide functionality like analytics, machine learning as a service (the Prediction API) or access to user data (when permission to read the data is given). Another important example is an embedded Google map on a website, which can be achieved using the Static maps API, Places API or Google Earth API.

---

### Common uses

- **User registration** is commonly done via Google+ sign in, which allows users to securely log in to 3rd party services with their Google+ account using the Google+ API. This is currently available from within Android, iOS or JavaScript. It is popular to include a “Sign in with Google” button in Android apps, as typing login credentials manually is time-consuming due to limited screen size. As the user is usually signed into their Google account on their mobile device, signing-in/signing-up for a new service with a Google is usually a matter of a few button clicks.
- **Drive apps** are various web applications (often third party) which work within Google Drive using the Drive API. Users can integrate these apps into their Drive from the Chrome Web Store which allows them to do work entirely in the cloud.<sup>[9]</sup> There are many apps available for collaborative document editing (Google Docs, Sheets), picture/video editing, work management or for sketching diagrams and workflows.
- **Custom Search** allows web developers to provide a search of their own website by embedding a custom search box and using the Custom Search API. They can customize the search results and make money off the ads shown using AdSense for Search.
- **App Engine apps** are web apps that run on the Google App Engine, a platform-as-a-service (PaaS) cloud computing platform which allows web developers to run their

websites in Google datacentres. These web apps often take advantage of APIs to manipulate services such as TaskQueue (a distributed queue), BigQuery (a scalable database based on Dremel) or DataStore.

- **Gadgets** are mini-applications built in HTML, JavaScript, Flash and Silverlight that can be embedded in webpages and other apps. They can run on multiple sites and products (even writing them once allow users to run them in multiple places).

## 1.6 Frequency Visualization

Cirrus is a word cloud displaying the frequency of words appearing in a corpus. Words occurring more frequently appear larger. Cirrus is a freely available visualization tool that generates a word cloud from a document or group of documents. The initial word cloud contains all of the most frequent words from the text, laid out graphically to fit together within an elliptical shape. Words may be oriented either vertically or horizontally, and they are rendered in different colours for ease of viewing. The user should be aware that word orientation and colour are strictly decorative elements, they do not indicate anything meaningful about the word. The size of the word, however, indicates the frequency with which it appears in the document. The larger the font size, the more frequent the word. The user can mouse over a word to see the precise number of times that it appears in the document. One of the problems with word frequency visualization is the prevalence of conjunctions and articles in most texts. These less semantically interesting words often drown out the adjectives, adverbs, nouns, and verbs that are more likely to be of interest to the researcher. Cirrus offers an easy solution to this problem. By clicking the options icon in the upper right hand menu, a user can apply a list of stop words to the visualization. These stop words will be treated as noise, and stripped from the cloud, allowing words of greater interest to surface. Cirrus offers two pre-built lists of stop words, one for common English words, and one for common French words. The user can view the lists, so she know which words will be removed before applying them.

## 1.7 Maptive

Maptive is an online tool Transform raw location data into a beautiful, customized Google map in a matter of seconds.

### FEATURES:

- Large or small data sets, 10 or 100,000 locations.
- Safe & Secure
- Cloud-Based
- Get access to all the customization and editing tools you'll ever need to create beautiful, customized maps.
- Create territories that make sense and improve the profitability of your business.
- Precision level control makes it easy to filter data and show only what you want to see on the map you create and customize.
- Optimize Routes / Directions
- Fully customized and configured to meet the unique needs of any business or organization—large or small.
- Share your map privately with specific individuals, publish publicly to the web, or embed within a webpage or blog using simple HTML code.
- The maps you create and share can be easily accessed and viewed on any smartphone or tablet devices.
- Track Your Location with GPS

## 1.8 Plotly

**Plotly**, also known by its URL, [Plot.ly](#), <sup>[1]</sup> is an online analytics and data visualization tool, headquartered in Montreal, Quebec. Plotly provides online graphing, analytics, and stats tools for individuals and collaboration, as well as scientific graphing libraries for Python, MATLAB, Perl Julia, Arduino, and REST

**Plotly** has a graphical user interface for importing and analysing data into a grid and using stats tools. Graphs can be embedded or downloaded. Mainly used to make creating graphs faster and more efficient.

Plotly is the easiest way to graph and share your data. It standardizes the graphing interface across scientific computing languages (Python, R, MATLAB, etc.) while giving rich interactivity and web share ability that has not been possible before with matplotlib, ggplot, MATLAB, etc.

On the Plotly website, you can style your graphs with a GUI, so you don't have to spend hours writing code that simply changes the legend opacity.

Plotly does this all while backing up your graphs on the cloud, so that years later, you can find data that may have otherwise been on a harddrive in a landfill. If you make your data public, other people can also find your graphs and data. The best practice that we have today for saving and sharing research data is to entomb it as a thesis in the engineering library basement. All that is changing.

## **CHAPTER 2**

### **TWITTER DATA ANALYTICS**

#### **2.1 INTRODUCTION TO TWITTER DATA ANALYTICS**

Social media sites like Twitter, Facebook helps to communicate and share information in mass .As a result, large and tons of data is present which can be used to discover information related to any news/event, person or any health issue. This information can thus be used to serve many purposes like:

- ▶ Developing surveillance system
- ▶ Mining information etc.

Twitter Data Analytics is divided into three categories:

- News Analysis
- Celebrity Analysis
- Health Analysis

#### **2.1.1 INTRODUCTION TO NEWS ANALYSIS**

The news media, and more specifically print media, serve as valuable sources of information and powerful modes of communication. This power controls much of what people understand of events that occur around the world on a daily basis. The way information is transferred to its recipients comes through various forms of communication, all of which is framed to meet the goals of the providing source. The aim of this project is to find out the thoughts of people about any news or event occurred. This is done using Rapid Miner, a java based software for extracting the tweets, processing it and analysing the feeling of the people by categorizing that into three categories: positive, negative and neutral

#### **2.1.2 INTRODUCTION TO CELEBRITY ANALYSIS**

Social media sites like Twitter, Facebook helps to communicate and share information in Mass .As a result, large and tons of data is present which can be used to discover Biomedical and health related information. This information can thus be used to serve Many purposes like:

- To analyse behaviour of any person about him/her
  - To analyse behaviour of other persons towards any him/her
  - To analyse follower's tweets.
  - Geomapping of follower's geolocation
  - Analyses of two twitter data analysing tools which are Tweetchup and Twitnomy.
- Our aim to analyse the data present on social media. People often post messages about Interesting news, their daily activities, thoughts and feelings. This may help to analyse Their behaviour and behaviour of other people's towards him/her. It makes easy to analyse Positive and negative attitude from their tweets. Celebrity's publicity or how

famous they are, what people think about them and how positive and negative image they have.

This can be analysing from tweets. Using their follower's tweets and geolocation it can be Easily analyse their behaviour and how their popularity varying country wise or any Specific area wise. It will also help to find terror area connection by getting number of Followers lying in that area. Twitter message, when used in combination, helps increase The accuracy of our prediction. The data is also verified by analysing twitter data using Twitter data analysing tools which are twitonomy&tweetchup.

The basic need for social networking behaviour analyzing is to analyse the data posted into

Twitter and to transform those tweet related data into related information for tweet Analysis. Also to fetch follower's information and any particular person's information and

Relate it to the required information in order to make relations among these info and to Mine the behaviour related data so as to analyze positive and negative attitude from their Tweets. Celebrities publicity or how famous they are, what people think about them and How positive and negative image they have, all this can be analyze from tweets. Twitter Message, when used in combination, helps increase the accuracy of our prediction.

### **2.1.3 INTRODUCTION TO HEALTH ANALYSIS**

Twitter is a great data resource to monitor these health issues prevailing among general population. This project includes the extraction of tweets using Twitter API and then analysing them in many possible ways. Geo mapping is one of the important features explored in this project. Static Geo-Mapping and Dynamic Geo-Mapping are the two types which have been discussed. Apart from this, a tool named Cirrus has been used to take the analysing process to another step. Our aim is to handle the tweets on Twitter relating to health related issues and then extract the required information for analysis.

## **2.2 NEED FOR TWITTER DATA ANALYTICS**

The basic need for twitter data analytics is to analyse the data posted into Twitter and to transform those tweet related data into related information for tweet analysis. Also to fetch follower's information and relate it to the required information in order to make relations among these info and to mine the related data.

## **2.3 OBJECTIVE OF TWITTER DATA ANALYTICS**

Our aim is to handle the tweets on Twitter relating to any news, person or health related issues and then extract the required information for analysis. This project includes the extraction of tweets using Twitter API and Rapid Miner and then analysing them in many possible ways. Geo mapping is one of the important features explored in this project. Static Geo-Mapping and Dynamic Geo-Mapping are the two types which have been discussed. Apart from this, a tool named Cirrus has been used to take the analysing process to another step.

# TWITTER DATA ANALYTICS

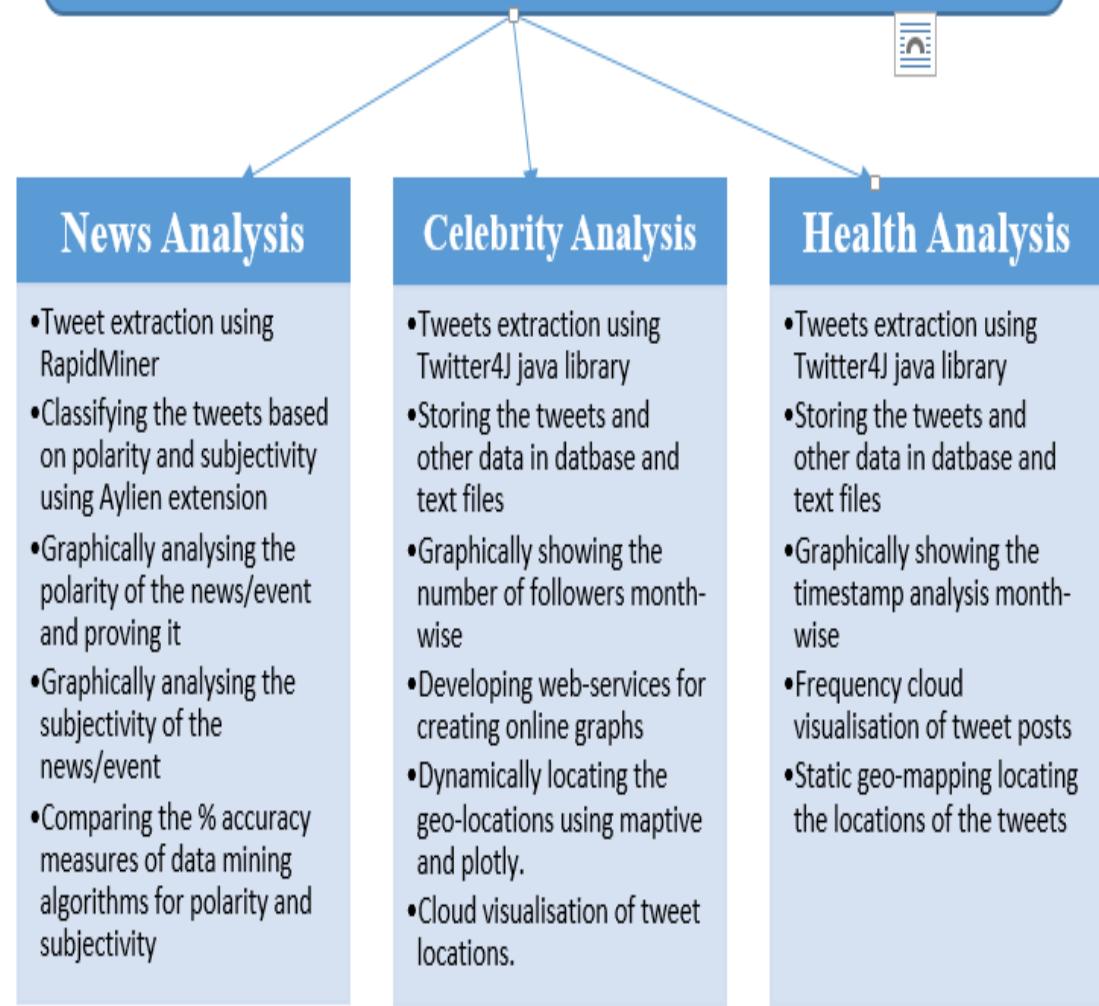


Figure 1: Flowchart showing the Twitter Data Analytics

## **CHAPTER 3**

# **DESIGNING AND IMPLEMENTATION OF TWITTER DATA ANALYTICS**

In this project we present a more general approach that discovers associations from tweets. This project first aims at extracting tweets using two different approaches. Various tweets' attributes like username, language of the tweet, month, year and location from where it was posted are extracted along with the tweet messages. These attributes, along with the original tweet message are then stored in various formats. Using this data from database, location frequency is calculated for every country, which is then used by JxBrowser and Google API service of maps marker for Geo-Mapping. Cirrus, another tool used for analysing the extracted data, is a word cloud displaying the frequency of words appearing in a corpus. This tool not only allows us to remove the stop words, but it also helps us to analyse the tweets by calculating the frequency of important and relevant words.

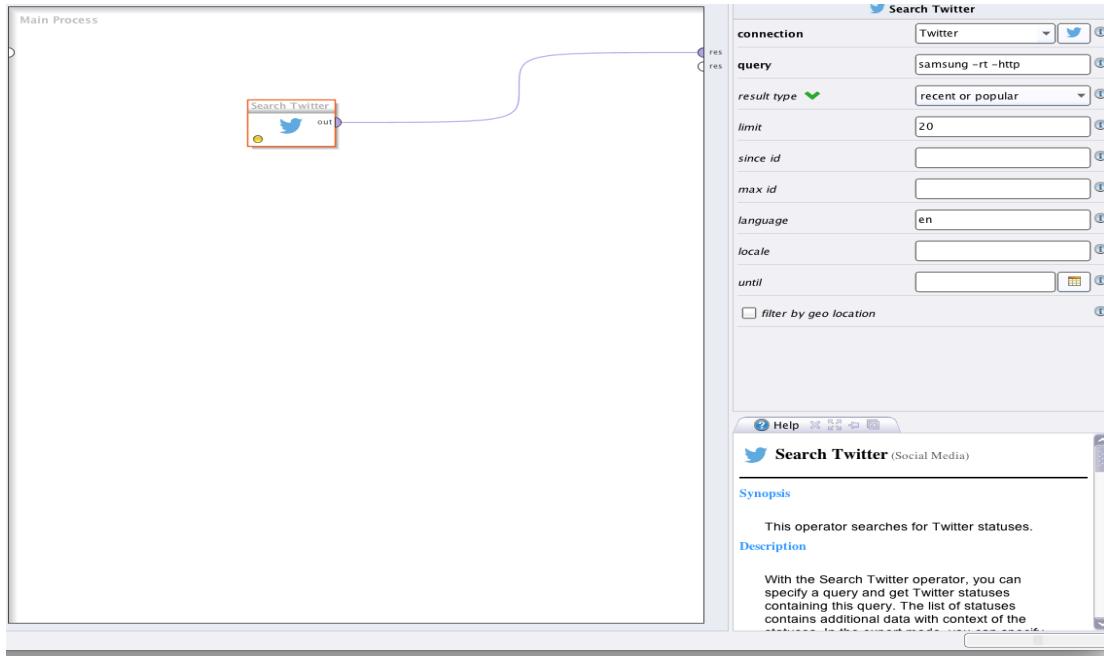
### **3.1 Tweets Extraction and Collection**

This phase is further categorised as per the three categories of our project.

#### **3.1.1 Tweets Extraction and Collection in News Analysis**

Using Search Twitter operator in Rapid Miner, tweets can be extracted and collected in table form. The table consists of columns as: from user, to user, text, relocations, source, language, rewet count.

Figure 2 shows an example on how to gather tweets in Rapid Miner. Below we're searching for tweets containing the keyword "Samsung" using Search Twitter Operator. We've cleaned up our search a little by removing rewets (-rat) and links (-http). We've also restricted the number of tweets to collect to 20 and decided we only want to see English tweets by adding "en" in the language parameter. We've also indicated that we want only recent or popular tweets to be returned using the Result type parameter.



**Figure2: Extraction of tweets**

After running the process as shown in figure1, the output screen is shown in figure 3 which is in table form consisting of multiple columns consisting of username, user-id, text,source,retweet count etc. proving various kinds of data to be analysed further.

ExampleSet (20 examples, 1 special attribute, 11 regular attributes)					
Row No.	Id	Created-At	From-User	Text	From
1	653285841	Oct 11, 201	Vala Afshar	Digital trucks save lives —Samsung trucks shows drivers traffic ahead http://t.co/TFXQFDyw9t	2597
2	653806661	Oct 13, 201	Jimmy Kimm	Tonight on #Kimmel @KirstenDunst #Fargo, @NathanFielder #NathanForYou & @Purity_Ring on the Samsung outdoor st:	3403
3	650686811	Oct 4, 2015	Samsung Mc	No need for new machines. Samsung Pay works where your card works. http://t.co/p0SpNjZID	2971
4	653967291	Oct 13, 201	VéroniqueG	My new Samsung galaxy s6 edge 64gb my girlfriend got me for my birthday :) http://t.co/AIxEPCvif	2165
5	653967201	Oct 13, 201	DunceTrista	#défence #Samsung Techwin K9 Thunder SPT howitzer is competing for a 800 M US \$ tender for 100 units in #India.	3865
6	653967177	Oct 13, 201	Heidi Kenne	@VentureBeat: LG file for G Pay trademark to take on Apple, Google, and Samsung... http://t.co/AX3flUmsr	3296
7	653967101	Oct 13, 201	Cool Geekz	Samsung tv fireplace mounted. Wires concealed in wall to prep for built in shelving to be installed.... https://t.co/LsbW	2946
8	653967091	Oct 13, 201	cubenstein	@FlareZephyr check the chip, see if he/she got tsmc or samsung	2504
9	653967081	Oct 13, 201	DunceTrista	The samsung m1 vodafone 360 incorporates touchscreen technical skill: juXBQWtY	1240
10	653967001	Oct 13, 201	blue neighb	i hate samsung and i hate this stupid fucking country	3218
11	653966871	Oct 13, 201	meem pls	Aww my new iPhone has the shitty Samsung chip instead of the good TSMC one. #chipgate I'm sad	2894
12	653966831	Oct 13, 201	D Vipa	#Fitbit #Samsung #Note4 Why do I have a notification that can't go away notifying me my notifications are running? http://tcc	3387
13	653966781	Oct 13, 201	VéroniqueG	#défence via #Jane's #India's MoD shortlisted #Samsung Techwin #K9 Thunder 155 mm/52-calibre self-propelled tra	3865
14	653966751	Oct 13, 201	Joaquín Qeu	Samsung Galaxy Tab Pro SM-T900 32GB, Wi-Fi, 12.2" - Black (Latest Model) #TV - Bid Now! Only \$280.01 http://tcc	3860
15	653966721	Oct 13, 201	blue neighb	@troyesivan @Beats1 @zanelowe HI. IM A SAMSUNG USER. HOW DO I LISTEN HELP	3218
16	653966701	Oct 13, 201	Duncan Coo	Hello sexy followers OMG I must be off my head I have just brought a new Samsung S6 edge £600.00 oh well it's onh	2984
17	653966681	Oct 13, 201	Brian'Mullig	The Samsung Galaxy S6 Edge Giveaway   SegmentNext Deals https://t.co/Xx3fbXRf9 via @SegmentNext	5995
18	653966531	Oct 13, 201	HRH James	@samsungireland S4 mini. 2nd time to happen to this device. Third time to happen with Samsung. Good job its nobrad.	2159
19	653966391	Oct 13, 201	Amanda De	Samsung is literally the worst phone. My note 3's camera is garbage...	4547
20	653966321	Oct 13, 201	Oktar Akin	We are shortlisted at LIA Non-Traditional category for 2 times with Samsung Hearing Hands. Yep! @leoburnett @LIAa	7074

**Figure 3: Example Set of gathered tweets**

### 3.1.2 Tweets Extraction and Collection in Celebrity Analysis

The developer registers on twitter to access APIs. The Twitter then issues customer token and secret. Then the application directs user to Twitter to verify user credentials where the user enters the required credentials. These credentials are verified by Twitter and if correct, then Twitter issues a OAuth verifier. After this, the developer requests access token through the application and Twitter issues him the access token and the secret key so that the user can get access to the requested contents and information.

Figure 4 represents that Twitter has issued an authentication URL to verify credentials of user and Figure 5 shows that user has been verified by authorizing app.

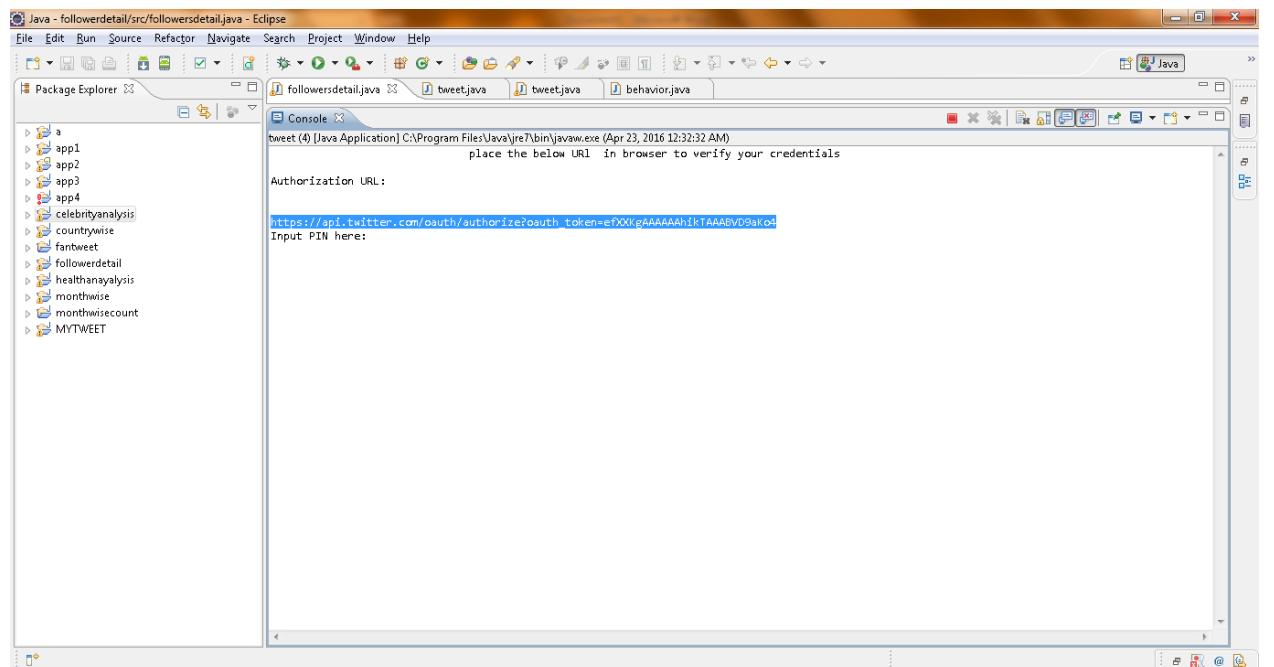


Figure 4: Authentication URL

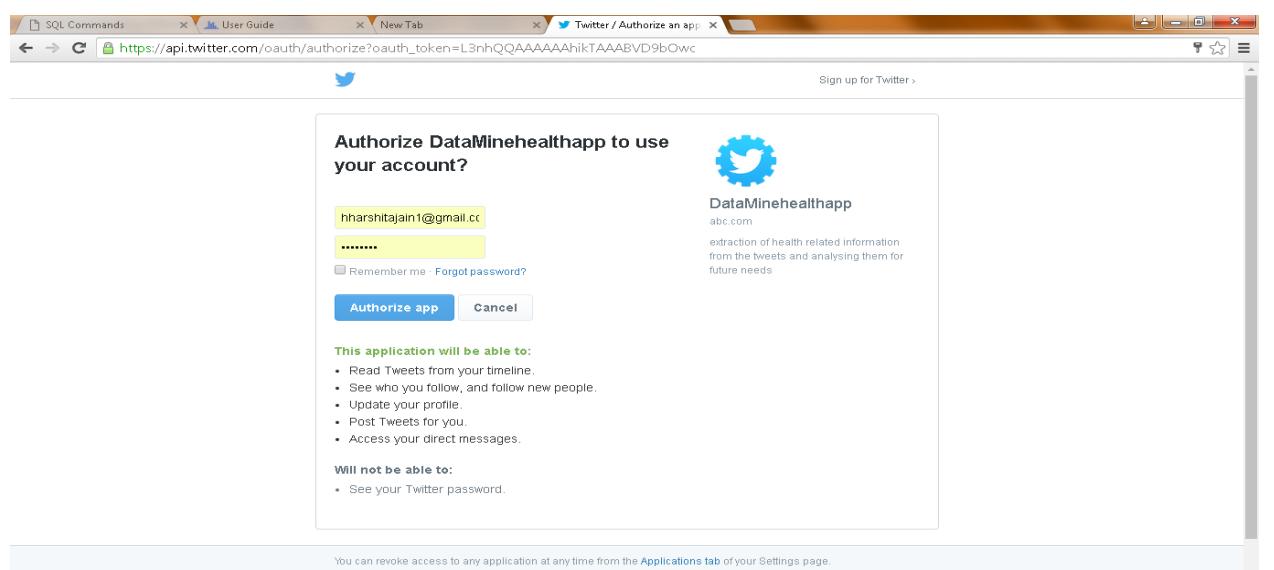
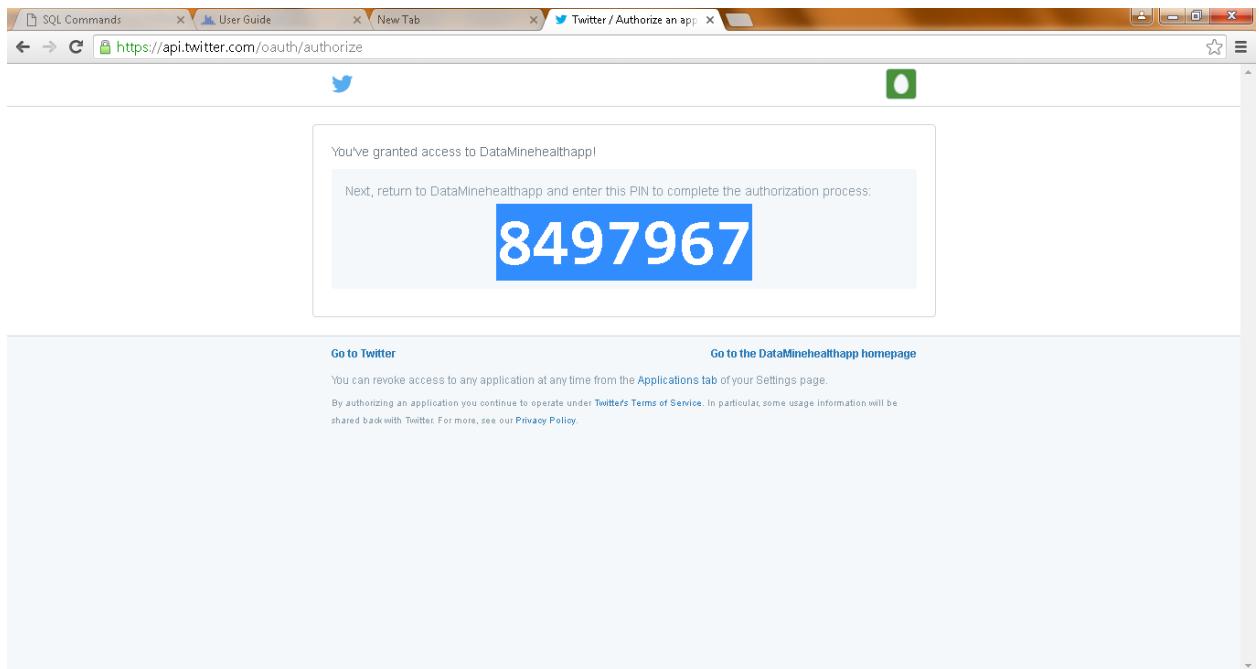


Figure 5: User credential are verified



**Figure 6: Pin for authentication**

The user/developer can get access to his application through the customer token and secret keys. Using these keys he access twitter to grant him the authentication pin which here is 8497967(figure 6) so as to verify the application on twitter .once the user verifies the application he/she has to enter the keyword on the console..

```

Java - followerdetail/src/followerdetail.java - Eclipse
File Edit Run Source Refactor Navigate Search Project Window Help
Package Explorer D3
followersdetail.java tweet.java tweet.java behavior.java
Console
tweet (4) [Java Application] C:\Program Files\Java\jre7\bin\javaw.exe (Apr 23, 2016 12:34:42 AM)
the following are your token details
Access Token: 3431346011$PGo7P7aNpiaD1vB79qTykh5I6HSiSISN53waQI
Access Token Secret: DpNaJuN8XOLjb2XoKlav3LdvAr39DtiYkSLUCHr1mCbjqK

Enter the string to which the related posts are to be displayed
BeingSalmanKhan
| wait your tweets are being fetching up !!!!  

| your tweets are
doone
$ID@723581143816110081$@LANGUAGE@$und$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Sat Apr 23 00:06:02 IST 2016$@TWEET$@https://t.co/dCacHgpRNG
$@EOTWEET$@
$ID@72339621705344205$@LANGUAGE@$en$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Fri Apr 22 11:51:12 IST 2016$@TWEET$@$support @powerminds by riding ur
$@EOTWEET$@
$ID@72321356548673536$@LANGUAGE@$en$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Thu Apr 21 23:45:23 IST 2016$@TWEET$@Great initiative by BMC for waste
$@EOTWEET$@
$ID@72321356773898240$@LANGUAGE@$en$@GEO@$NUNBAI$@TIMESTAMP@$Thu Apr 21 23:44:38 IST 2016$@TWEET$@My Galaxy is Garbage free !!! Who
$@EOTWEET$@
$ID@721665226553954816$@LANGUAGE@$en$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Sun Apr 17 17:12:52 IST 2016$@TWEET$@Lkg fwd 2 this film TRAFFIC http
$@EOTWEET$@
$ID@721411448170110976$@LANGUAGE@$und$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Sun Apr 17 00:24:27 IST 2016$@TWEET$@https://t.co/Z7SpvePIwb
$@EOTWEET$@
$ID@9774614151258112$@LANGUAGE@$et$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Tue Apr 12 12:00:15 IST 2016$@TWEET$@Aur yeh raha teaser #SultanTeaser
http://t.co/cobofUJu9
$@EOTWEET$@
$ID@9710499388528103424$@LANGUAGE@$in$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Mon Apr 11 17:46:36 IST 2016$@TWEET$@$Sultan ka Pehta Daav #SultanPoste
$@EOTWEET$@
$ID@971098979993378819$@LANGUAGE@$en$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Sun Apr 10 21:52:53 IST 2016$@TWEET$@So happy to see Anant Ambani, lots
$@EOTWEET$@
$ID@9717000747120865280$@LANGUAGE@$en$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Mon Apr 04 20:17:54 IST 2016$@TWEET$@Chalo r u ready for Bigg Boss 10
$@EOTWEET$@
$ID@716221811948707840$@LANGUAGE@$und$@PLACE@$null$@GEO@$NUNBAI$@TIMESTAMP@$Sat Apr 02 16:42:41 IST 2016$@TWEET$@https://t.co/c5laOthMZO
$@EOTWEET$@

```

**Figure 7: keyword"BeingSalmanKhan" entered**

After the authorization and authentication process the user/developer has to specify the keyword or twitter user id in the console output (as shown in figure 7) so that only those tweets and its related information can be collected in which the keyword or from that user id is present. Then once the specification of keyword is done then tweets get extracted on

the console itself and other information is stored in database.

### 3.1.3 Tweets Extraction and Collection in Health Analysis

A tweet has tweet text, user name, time-stamp, and location. The data pre-processor module reads the raw tweet data from the database, extracts the tweet text, tweet timestamp and user location, and stores the data back into the database.

A language constraint was applied while extracting tweets in order to get tweets of only one particular language, well known to us for better analysis of tweets. Once this constraint was successfully applied, the information attributes such as the username, location, timestamp and the tweet message posted were stored in a database. This was done as a step to make a move further for our analysis and research work.

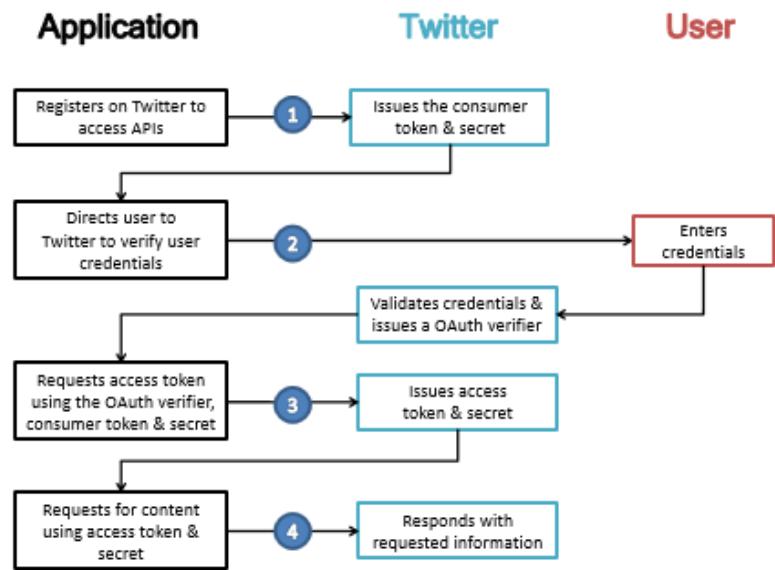
The data collector module continuously downloads the tweet data using Twitter Streaming API. Twitter Streaming API allows high-throughput near real-time access to global stream of public tweets that matches pre-specified filter predicates. Multiple parameters may be specified in a single connection to the streaming API to determine what tweets will be delivered on the stream. Firstly the tweets were extracted using twitter4j java library. Twitter4J is an unofficial Java library for the Twitter API. With Twitter4J, we easily integrated our Java application with the Twitter service.

Figure 8 shows the steps for extraction of tweet related data and collecting the information in database. Firstly the tweets were filtered on the basis on language constraint in order to extract tweets where language of the follower is English using “en” constraint where “en” stands for English language. After the application of language constraints, the tweet related information such as username, his id number, language, timestamp, geolocation and his/her posted tweet were extracted. Then the tweets were saved in text file and other information that i.e. username, his id number, language, timestamp, geolocation were stored in database so that sql queries can be applied for analysis.



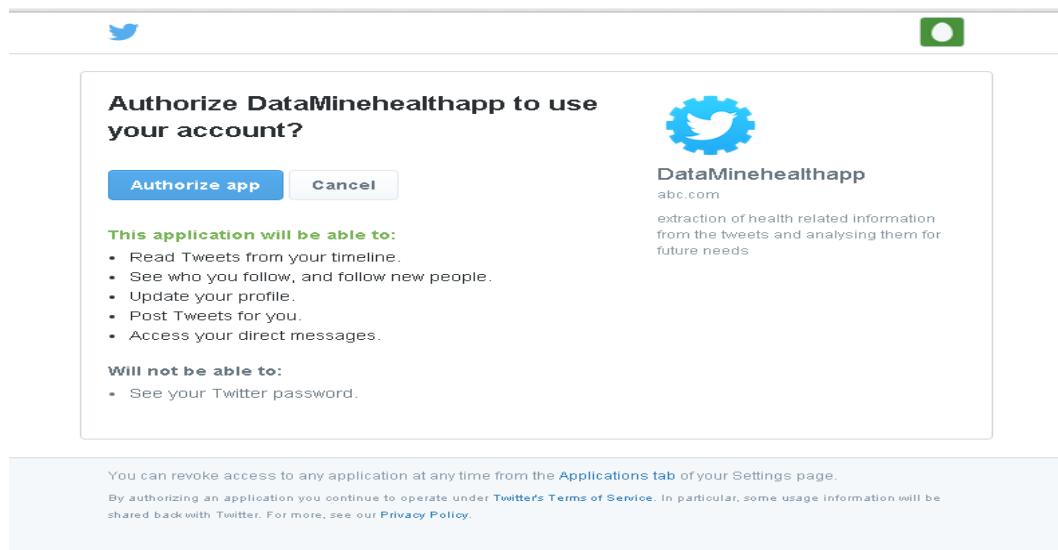
**Figure 8: Steps for tweet extraction and collection**

Figure 9 shows relationship amongst the twitter application, twitter and the user. With the help of application, the developer registers on twitter to access APIs. The twitter then issues customer token and secret then the application directs user to twitter to verify user credentials where the user enters the required credentials, these credentials are verified by twitter and if correct then twitter issues a OAuth verifier. After this the developer requests access token through the application and twitter issues him the access token and the secret key so that the user can get access to the requested contents and information.

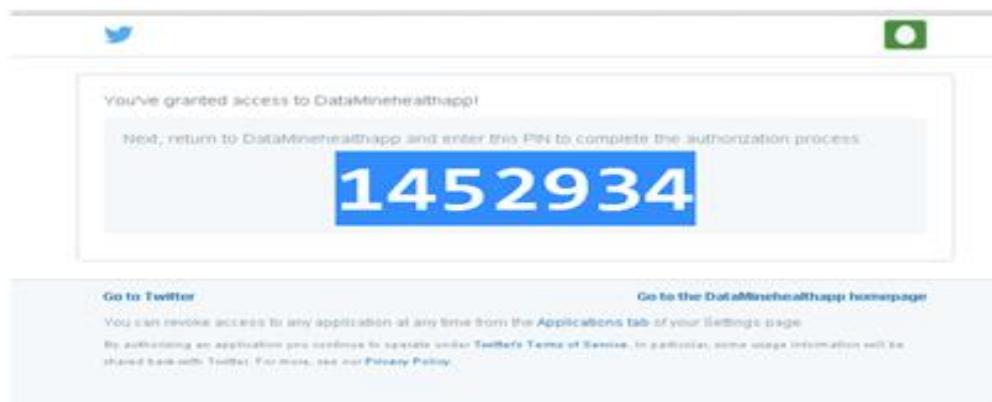


**Figure 9: Integration of Java Application with Twitter Services**

The user/developer can get access to his application through the customer token and secret keys. Using these keys he access twitter to grant him the authentication pin which here is 1452934(figure 11) so as to verify the application on twitter .once the user verifies the application he/she has to enter the keyword on the console..



**Figure 10: Authorization**



**Figure 11:Pin for authentication**

After the authorization and authentication process the user/developer has to specify the keyword in the console output (as shown in figure 12) so that only those tweets and its related information can be collected in which the keyword is present. Then once the specification of keyword is done then tweets get extracted on the console itself (figure 13) and other information is stored in database

```
Java - app1/src/tweet.java - Eclipse
File Edit Run Source Refactor Navigate Search Project Window Help
Console
tweet [Java Application] C:\Program Files\Java\jre7\bin\javaw.exe (Sep 16, 2015 10:55:10 AM)
        place the below URL in browser to verify your credentials

Authorization URL:

https://api.twitter.com/oauth/authorize?oauth_token=exXlZwAAAAAAwhkTAAABT9sdizN
Input PIN Here: 1452934
        the following are your token details

Access Token: 3431346011-2PGoP7aMpiaDlV879gTykhS1GHS1SN5v3wQzI
Access Token Secret: DpNu0u8XOL3b2KOkIav3ldVAr390tIyksLchRimCJbqk

Enter the string to which the related posts are to be displayed
debug
```

**Figure 12:** Keyword specified like dengue

**Figure 13:**tweets extracted for Health analysis

## 3.2STORING THE EXTRACTED DATA

The extracted data from Twitter is stored into excel files (in case of news analysis) and in database and text files (in case of celebrity and health analysis).

### 3.2.1STORING THE TWITTER DATA IN EXCEL FILES

In news analysis, after extracting Tweets and other data like usernames,retweet count,id,location,source are stored in excel files for further analysis and visualisation.

As shown in figure 14, the tweet posts and other data like usernames,retweet count,id,location,source are stored in excel files for further analysis and visualisation of news.

1	Created From-User	From-User Id	To-User	To-User Id	Langua	Source	Text	Geo-Lo	Geo-Lc	Retwee Id	
2	NDTV	37034483.0		-10.en	cahrel	"Freedom 251launched as 'world's cheapest smartphone' at Rs 251https://t.co/v44JNqJr75		198.0	###		
3	TIMESSHOW	240649814.0		-10.en	cahrel	"Indian company Ringing Bells launches the world's cheapest smartphone Freedom 251 priced at ₹ 251https://t.co/LLr1Gm5oz8		104.0	###		
4	Gadgets 360	4307446.0		-10.en	cahrel	"Here's an exclusive first look at Freedom 251, the Rs. 251 smartphone that's making the headlines today. #Freedom 251https://t.co/z2nN67H9yx		186.0	###		
5	Aleks Moiss	17339370.0		-10.en	cahrel	"Freedom 251, a smartphone with a price tag of 251. But will it remain the cheapest option available? https://t.co/JakzJohY95		.0	###		
6	Adami J	2829746594.0		-10.en	cahrel	"India #NewDeals 45 Ringing Bells 251 smartphone goes on sale in India, website promptly pushes @FollowNewsNow		.0	###		
7	Times of India People	4000000.0		-10.en	cahrel	"Freedom 251 launch: Indian company Ringing Bells launches the world's cheapest smartphone Freedom 251 (via VideoPing) Be... https://t.co/nv3nJmLM#gadgets		.0	###		
8	Police4Jobs	2893730625.0		-10.en	cahrel	"#Breaking Freedom 251 website crashes after it received 6 lac hits per second Read More: https://t.co/dgJgNgodu		.0	###		
9	police4jobs	2893730625.0		-10.en	cahrel	"#Breaking Freedom 251 website crashes after it received 6 lac hits per second Read More: https://t.co/dgJgNgodu		.0	###		
10	ajaykayna	2520438893.0	singhanu	330345112.0.en	cahrel	"@singhanu3 pingutne freedom 251 order kya kya?"		.0	###		
11	Adharsh Bharadwaj	136253203.0	ajuns	39535322.0.en	cahrel	"@ajuns @nercitizen get freedom 251?"		.0	###		
12	Jyotirmoy Roy	27620030943.0	RingingB	4768566844.0.en	cahrel	"@RingingB why don't you give that freedom 251to sell Amazon, or any online shopping site ????		.0	###		
13	Raul Rubio™	103683684.0		-10.es	cahrel	"Conoces de Freedom 251 un smartphone con especificaciones "decentes" a un precio de escándalo. https://t.co/wHmEykqbQD		.0	###		
14	Vijay Chavsey	4200000.0		-10.en	cahrel	"Freedom 251 or freedom 420		.0	###		
15	Soham Bhattacharj	15222000.0	RingingB	4768566844.0.en	cahrel	"@RingingB I am waiting for that when will the freedom 251 will be back in stock."		.0	###		
16	Soham Bhattacharj	6953958.0		-10.en	cahrel	"@RingingB I am waiting for that when will the freedom 251 will be back in stock."		4.0	###		
17	Crystal Manandhar	54905982.0		-10.en	cahrel	"@freedom 251 website akash banaya ja raha h janai ka jana Kya samble freedom 251 sir farj aye aah h		.0	###		
18	Shivshakti Mihra Shriv	4237611732.0	narendra	18839785.0.hi	cahrel	"@narendra modi woh sab akash banaya ja raha h janai ka jana Kya samble freedom 251 sir farj aye aah h		.0	###		
19	Jiendra Singh	74422349.0		-10.en	cahrel	"@RingingB I am waiting for that when will the freedom 251 will be back in stock."		.0	###		
20	srujana jinku	184349958.0		-10.en	cahrel	"@RingingB I am waiting for that when will the freedom 251 will be back in stock."		.0	###		
21	Bhanwar kanwar	*****		-10.en	cahrel	"I want to booking freedom 251 please book name bhanwar kanwar address 262 nico housing board bhwara		.0	###		
22	Sheetal Bhandarpur	140739795.0		-10.en	cahrel	"Isn't credibility the importance compared to Freedom 251?"		.0	###		
23	Green Beaver	4850400.0		-10.en	cahrel	"Freedom 251 has a sexy cover that doesn't slug."		.0	###		
24	sunil shukla	320056362.0		-10.en	cahrel	"Freedom 251 is a very good smartphone."		.0	###		
25	Prat-Earthman	387039863.0		-10.en	cahrel	"Technology is great but #India Freedom 251 needs help me https://t.co/77QXNjPv		.0	###		
26	Anku Chauhan	10142474.0		-10.en	cahrel	"I still don't know anyone who managed to buy a book freedom 251"		.0	###		
27	mayakrishnan	3220598362.0		-10.en	cahrel	"How to get freedom 251 smartphonehttps://t.co/4qGFD2Fh"		.0	###		
28	vishal	18763397810.0	MKBHD	29873662.0.en	cahrel	"@MKBHD Check out Freedom 251"		.0	###		
29	ayush kataria	2747806438.0		-10.en	cahrel	"I think that freedom 251 is full too (price) (fraud) #Freedom251"		.0	###		
30	ahmed jo	123036637.0		-10.en	cahrel	"Reduced quantity and price too reduced #Freedom 251 looks like a scam https://t.co/llubJbJEIDZ"		.0	###		
31	Angus Sharma	469203637.0		-10.en	cahrel	"What's up with the price freedom 251?"		.0	###		
32	Manoj-Sharma	489258267.0		-10.en	cahrel	"@vishal shukla what's up with freedom 251?"		.0	###		
33	Nathani	764345390.0		-10.en	cahrel	"Did any one booked freedom 251?"		2.0	###		
34	Prachi	72823253.0		-10.en	cahrel	"Walls into the temple, breaks coconut, zipz coconut water.... receives Freedom 251 in return."		.0	###		

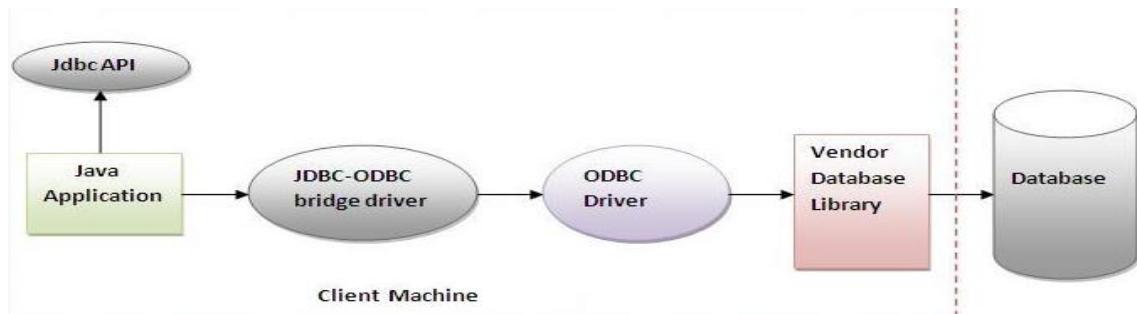
Figure 14: Excel file storage of news tweets

### 3.2.2STORING THE TWITTER DATA IN DATABASE AND TEXT FILES

After storing the extracted attributes in the database for celebrity and health, appropriate queries were used to calculate the frequency of tweets and other issues required for analysis.

For storing the data extracted, Jdbc odbc has been implemented using oracle 10g express edition.

Figure 15 explains that the java application has been integrated through JDBC ODBC bridge driver and a vendor library as (ODBC14.JAR).java application has been integrated through odbc driver with oracle 10 g express edition database.



**Figure 15: Client Machine**

### **Storing the extracted data of celebrity**

In this project we present a more general approach that analyze behaviour of any person from his/her tweets and behaviour of his/her followers towards him/her from tweets. This project first aims at extracting tweets from Twitter API. Various tweets' attributes like username, language of the tweet, month, year and location from where it was posted are extracted along with the tweet messages. These attributes, along with the original tweet message are then stored in a database. Two text files are stored which contains positive and negative words list. By comparing tweets of followers and other peoples ,tweets are analyzed to know behaviour about that person. For analyzing, Cirrus, another tool used for

analyzing the extracted data, is a word cloud displaying the frequency of words appearing in a corpus. This tool not only allows us to remove the stop words, but it also helps us to analyze the tweets frequency month wise to know activeness on twitter is increasing or decreasing .Using follower's geolocation, geomapping is done to find how much followers are lying in terror area so that it will help to find connection of that persons contacts and it will also help to analyse popularity of any celebrity varying country wise or any specific area wise. Two twitter data analysis tools are also analysed which are twitonomy and tweetchup to verify data which we get after java coding and by tool.

### **Tweets are extracted for celebrity “ salman khan” and stored in DataBase**

Firstly the tweets were extracted usingtwitter4j java library. Twitter4J is an unofficial Java library for the Twitter API. WithTwitter4J, we easily integrated our Java application with the Twitter service.

Firstly the tweets were filtered on the basis on language constraint in order to extract tweets where language of the follower is English using “en”constraint where “en” stands for English language. After the application of language constraints, the tweet related information such as username, his id number, language,timestamp, geolocation and his/her posted tweet were extracted. Then the tweets were saved in text file and other information that i.e. username, his id number, language,timestamp, geolocation were stored in database so that sql queries can be applied for analysis.

Figure 16represents the collection of tweets extracted of the celebrity “Salman Khan” With the id “BeingSalmanKhan”. Figure 17 represents the description of the table “celebrity”. the screenname column represents the screenname or username of the follower of salman khan. language stores the language and location stores the location from where the user has tweeted and time represents their time of tweeting.

"you're not a robot. You can't just conjure up motivation when you don't have it." <https://t.co/NaMKxdu1iv>  
"Your life is tetris. Stop Playing It Like Chess." <https://t.co/VYC5mF4CJ>

"Semantics" <https://t.co/0qkpeec21o>

RT @hadip: US vote this week will erode my freedom, make me a 2nd-class citizen. [Please help me: https://t.co/kdsw..](https://t.co/TBQo73Merd)  
"It's not about the altitude, it's the attitude."  
"The day I became a millionaire" by @dh <https://t.co/xpgvTlIpj8>  
"Shift Your Mindset By Saying Less of These Four Things" <https://t.co/BAb7urT5pk>

#siteX conference in Dubai. Robert Scoble @scobleizer on stage #siteX2015 [http://t.co/wcoUrxvc](https://t.co/wcoUrxvc)  
on excusing us for making us "creative". Source: <https://t.co/ndaoCMrTna> [http://t.co/SC2acTsfwn](https://t.co/SC2acTsfwn)

If I look at the 80+ countries we've visited so far, at the moment, Iran is the ideal place to be an entrepreneur" <https://t.co/jalcAHHiq>

RT @anidash: Lots of takes on business insider sale; I see it as yet another blog that launched on Movable Type & sold for \$400M+. <http://t.co/..>  
6 years vs. 6 years <https://t.co/4kpargrnxm> than washington post. (via @zerohedge) <http://t.co/AbT1gxggsc>

"In establishing the rule of law, the first 5 centuries are always the hardest" (road to prosperity is hard & long) <http://t.co/2HbfQaqnlh>

"I want to deliver a product that our customers want, not one that our investors want" by @jasonfried <https://t.co/nqo8zkcuhk>

"Why Blacksmiths Are Better At Startups Than You" <https://t.co/UTmbqfj0z2>

RT @ewy: Launched The Million Dollar Homepage 10 years ago today... how time flies! [http://t.co/oMhpqzwxDLk](https://t.co/oMhpqzwxDLk)

"You're not a robot. You can't just conjure up motivation when you don't have it." <https://t.co/NaMKxdu1iv>

"Your life is tetris. Stop Playing It Like Chess." <https://t.co/VYC5mF4CJ>

RT @hadip: US vote this week will erode my freedom, make me a 2nd-class citizen. [Please help me: https://t.co/kdsw..](https://t.co/TBQo73Merd)  
"It's not about the altitude, it's the attitude."  
"The day I became a millionaire" by @dh <https://t.co/xpgvTlIpj8>  
"Shift Your Mindset By Saying Less of These Four Things" <https://t.co/BAb7urT5pk>

#siteX conference in Dubai. Robert Scoble @scobleizer on stage #siteX2015 [http://t.co/wcoUrxvc](https://t.co/wcoUrxvc)  
on excusing us for making us "creative". Source: <https://t.co/ndaoCMrTna> [http://t.co/SC2acTsfwn](https://t.co/SC2acTsfwn)

If I look at the 80+ countries we've visited so far, at the moment, Iran is the ideal place to be an entrepreneur" <https://t.co/jalcAHHiq>

RT @anidash: Lots of takes on business insider sale; I see it as yet another blog that launched on Movable Type & sold for \$400M+. <http://t.co/..>  
6 years vs. 6 years <https://t.co/4kpargrnxm> than washington post. (via @zerohedge) <http://t.co/AbT1gxggsc>

"In establishing the rule of law, the first 5 centuries are always the hardest" (road to prosperity is hard & long) <http://t.co/2HbfQaqnlh>

"I want to deliver a product that our customers want, not one that our investors want" by @jasonfried <https://t.co/nqo8zkcuhk>

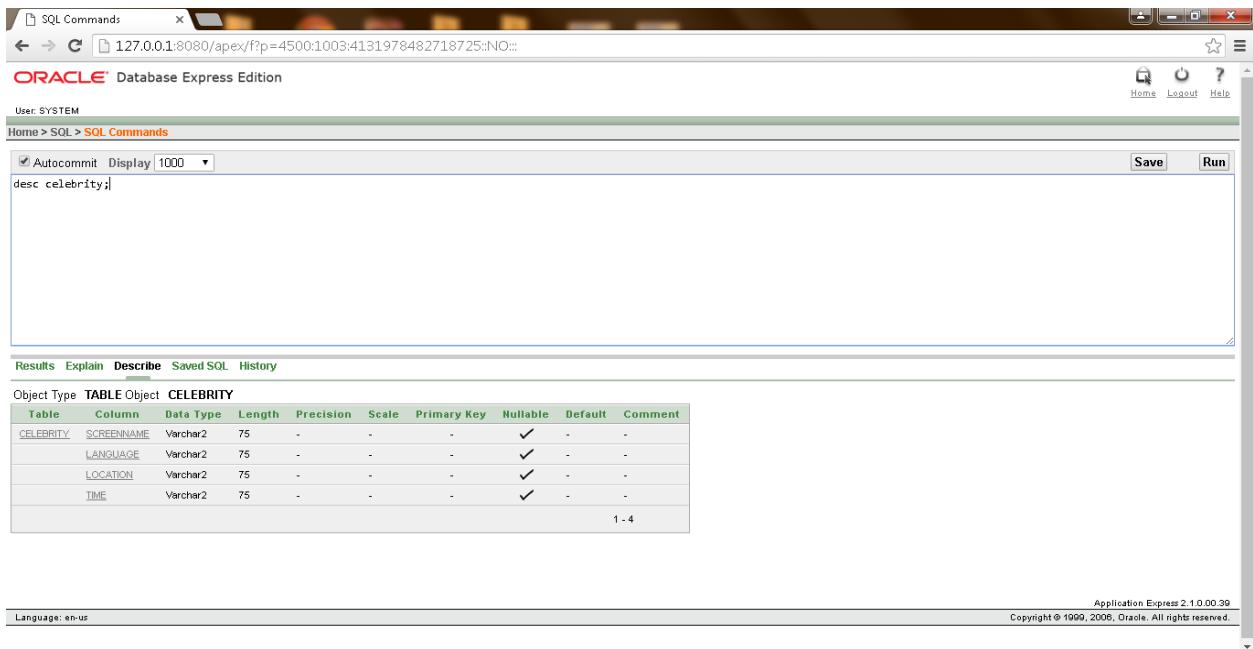
"Why Blacksmiths Are Better At Startups Than You" <https://t.co/UTmbqfj0z2>

RT @ewy: Launched The Million Dollar Homepage 10 years ago today... how time flies! [http://t.co/oMhpqzwxDLk](https://t.co/oMhpqzwxDLk)

Surat kemo ch! I am coming to ur city tomorrow for Da- bang! Chalo jaldi tickets book karlo! Will see u all soon .  
watching rocky 1 now... taaa taaa taaa taaa taaa Taaaaaa! thestlystalloneis adm1 se seekho yeh hal ho <https://t.co/2Twf1Zmu0>  
samai mal aayi <https://t.co/RuvPuvP903>

: ) <https://t.co/qHnHQ70tdx>  
Celebrate with us with me tomorrow! watch the world Television Premiere of Prem Ratan Dhan Payo on Star GOLD at 1 PM!  
#salmankhan.txt - Notepad  
File Edit Format View Help  
"you're not a robot. You can't just conjure up motivation when you don't have it." <https://t.co/NaMKxdu1iv>  
"Your life is tetris. Stop Playing It Like Chess." <https://t.co/VYC5mF4CJ>

**Figure 16:Tweets extracted for celebrity**



**Figure 17:Tweets Storage in Database for celebrity**

Categorization of tweets extracted into “positive” and “negative” on the basis of ontology specified

The tweets extracted are further analysed using the java code into positive and negative tweets on the basis of ontology as specified by us. firstly there are two text files containing positive and negative words which are fed as input to java code with the tweets of salman khan and as result various text files are created showing which the positive tweets are and why or due to which word it has been stated as positive and same with the negative case.

Figure 18 represent the various negative word in each line of text file for example “2 faced”, “abort”. Figure 19 represent the various positive word in each line of text file for example “happy”, “a+”.

```
negative.txt - Notepad
File Edit Format View Help
2-faced
2-faces
abnormal
abolish
abominable
abominably
abominate
abomination
abort
aborted
aborts
abrade
abrasive
abrupt
abruptly
abscond
absence
absent-minded
absentminded
absurd
absurdity
absurdly
absurdness
abuse
abused
abuses
abusive
abyssal
abyssally
abysses
accidental
accost
accursed
accusation
accusations
accuse
accuses
accusing
accusingly
acerbate
acerbic
acerbically
ache
ached
aches
aching
acrid
acridly
acridness
acrimonious
```

Figure 18:ontology of negative words

```
positive.txt - Notepad
File Edit Format View Help
happy
delicious
good
sweet
yes
yeah
a+
abound
abounds
abundance
abundant
accessible
accessible
acclaim
acclaimed
acclimation
accolade
accolades
accommodative
accommodative
accommodating
accomplished
accomplishment
accomplishments
accurate
accurately
achievable
achievement
achievements
achievable
adaptable
adaptive
adequate
adjustable
admirable
admirably
admiration
admire
admirer
adoring
adoring
adorably
adorable
adore
adored
adorer
adoring
adoringly
adroit
adroitly
adulate
```

Figure 19:ontology of positive word

Figure 20: represents the categorization of positive tweets into text file from the text file containing both positive and negative tweets i.e salmankhan.txt and Figure 21: represents the positive words which have been used from the ontology during categorization. for example the tweet "First round of legal battle in #Salman case in SC won by Attorney General who convinces court that appeal needs to be exa..." is positive due to the word "won" stored in ontology.

**Figure 20:positive tweets extracted for keyword "BeingSalmanKhan"**

**Figure 21:**positive words used from the ontology

Fig 22: represents the categorization of negative tweets into text file from the text file containing both positive and negative tweets i.e salmankhan.txt and Fig 23: represents the negative words which have been used from the ontology during categorization. for example the tweet "SC sends notice to Salman on plea against acquittal" is negative due to the word "plea" stored in ontology.

**Figure 22:negative tweets extracted for keyword”BeingSalmanKhan”**

**Figure 23:negative words used from the ontology**

**Follower details for salman khan are extracted and stored in text file**

The followers of salman khan on twitter are extracted and their details are stored in database.

Fig.24 contains two column screenname representing the follower name with their location.

The screenshot shows the Oracle Database Express Edition interface. In the SQL Commands window, the command `desc followerdata;` is run, displaying the table structure:

Table	Column	Data Type	Length	Precision	Scale	Primary Key	Nullable	Default	Comment
FOLLOWERDATA	SCREENNAME	Varchar2	75	-	-	-	✓	-	-
	LOCATION	Varchar2	75	-	-	-	✓	-	-

Below the table definition, the results of the `select * from followerdata;` query are shown:

SCREENNAME	LOCATION
shreyashingala	India
BinaChoudhary	-
raj444patil	-
UmeshPatil17454	-
kathirmenu	-
4fac730676a467	-
HrMathurc18	-
Khushi90206029	-
mtukulap	-
prakashbaghel14	-
1234_rama	delhi,india
AmaadChmanzoor	-
Javed455648	-
sahidsaind18	-
IshaqulSyeda	-
Brajakishore17	-
esa749c8b8345c	-
LadDev	-
kamalpr04733232	-
fareenif15	-

Figure 24:database for followerData

The location extracted is then analysed using geolocation showing which country has more followers and terror stated like khazakistan does has any follower indicating the popularity of the celebrity region wise

The screenshot shows the Oracle Database Express Edition interface. In the SQL Commands window, the command `select * from followerdata;` is run, displaying the results of the query:

SCREENNAME	LOCATION
shreyashingala	India
BinaChoudhary	-
raj444patil	-
UmeshPatil17454	-
kathirmenu	-
4fac730676a467	-
HrMathurc18	-
Khushi90206029	-
mtukulap	-
prakashbaghel14	-
1234_rama	delhi,india
AmaadChmanzoor	-
Javed455648	-
sahidsaind18	-
IshaqulSyeda	-
Brajakishore17	-
esa749c8b8345c	-
LadDev	-
kamalpr04733232	-
fareenif15	-

Figure 25:follower detail with their location

Fig.25 shows the relative position of various followers from various parts of country.for example “shreyashingala is from”india”

This data can now be exported to csv file and then later on fed to the Maptive to generate geolocations showing the user and the relative frequency of the location.further. it can be categorized into country wise mapping which indicates which country is active

## Follower views about salman khan are extracted and analysed monthwise

The views of the salman khan followers are extracted using java code month wise by paging.the count of relative month has been collected which can further be used to indicate as to which month did the salman khan was popular among people or their followers and by what percentage

Fig26,27,28:shows the collection of follower views in various months like September(sep.txt) December January(jan.txt) and which are further counted using java code

```

jan.txt - Notepad
File Edit Format View Help
Watch the launch episode of Comedy Nights Live at 10pm. All the best and do whatever u want to do man, but don't trouble your mother.
Throw back London dreams . https://t.co/SAJNec7smp
Thanku for being on Bigg boss @sauravindiatv @anjanaomkashyap https://t.co/X5h2j821AV
dolby k1 doll Mouth publicity se collections aaJ have picked up. Glad u guys hv liked the film n spoken abt it.wah yaar kamal karte ho aap
#BeingHumanClothing now available on the south African website http://t.co/P9R1c3xvck
Thank you for being on Bigg Boss. I am Bin Laden - dead or alive https://t.co/vizbzbdyta
dolly k1 doll Mouth publicity se collections aaJ have picked up. Glad u guys hv liked the film n spoken abt it.wah yaar kamal karte ho aap
want to find our k1 aakhir day hame kiski? May be this will give you some clue : http://t.co/kMr51EU1q3
see it , its too much fun like i said
The whole baaghi bhaijaan team watched the movie in a theatre close to us abt 35 km frm us n v loved it . Ma kasaam
Dilwale , script the word is superb . u guys wid love it
Pulkit
Abhishek dogra as director
Amitabh Bachchan as producer
Much
Courier service chahiye.. 4 days wait karne padega.. http://t.co/4kn5zrjfics
well done fans guys .. well done
Thank u to all the fans who have seen Jai Ho n liked it . I totally appreciate it guys . u been wonderful , god bless
@lalitkumar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/weslfifk700
@BeingGaurav This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/MCXL90QNVg
@mragaginiya This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/r42vZkaph
@lalitkumar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0V0L0d
@gala_pankti This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/peV72ykoGg
@yassar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/201hdmZis
@lalitkumar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0Xf3u3iB
@iamrahulraz This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/nqpx21bzFz
@hereforsalmankhan This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/Bwcp9ZHT0g
@pawan_kr This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0X
@greenian This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/jyj3wvkwrt
@MMIndia This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/rm4tP8nHrn
@S14 This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/TmV3emj
@S14Amer This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/lnf0T39
@lap20in This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/u401otd4kc
@Laugh_Riot This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/2YvdW63L
@Shahrukhkhan This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0Q
@mrnameet This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/uaq/cIOPA6
@sandeepmungar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/F5ywupxpz
@shahrukh This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0
@zuhra_afshan This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/mmTz2a9v1n
@deherajksingh08 This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/4cosAKo1s
@pawankalyan This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0up0xit
@bhayakrm This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/l8dhkR7a
@SNElazar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/1wejazwXOp
@bhit_lion_lover This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/mosadised
@saurabhik This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0
@saurabhik This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/gn2easdprN
@sandeepmungar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/m3awor5f2
@snehal_khatri This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/0
@SNElazar This Republic Day, remember, you do not need to wear a uniform to serve the country. Jai Ho http://t.co/txFdpotk0m

```

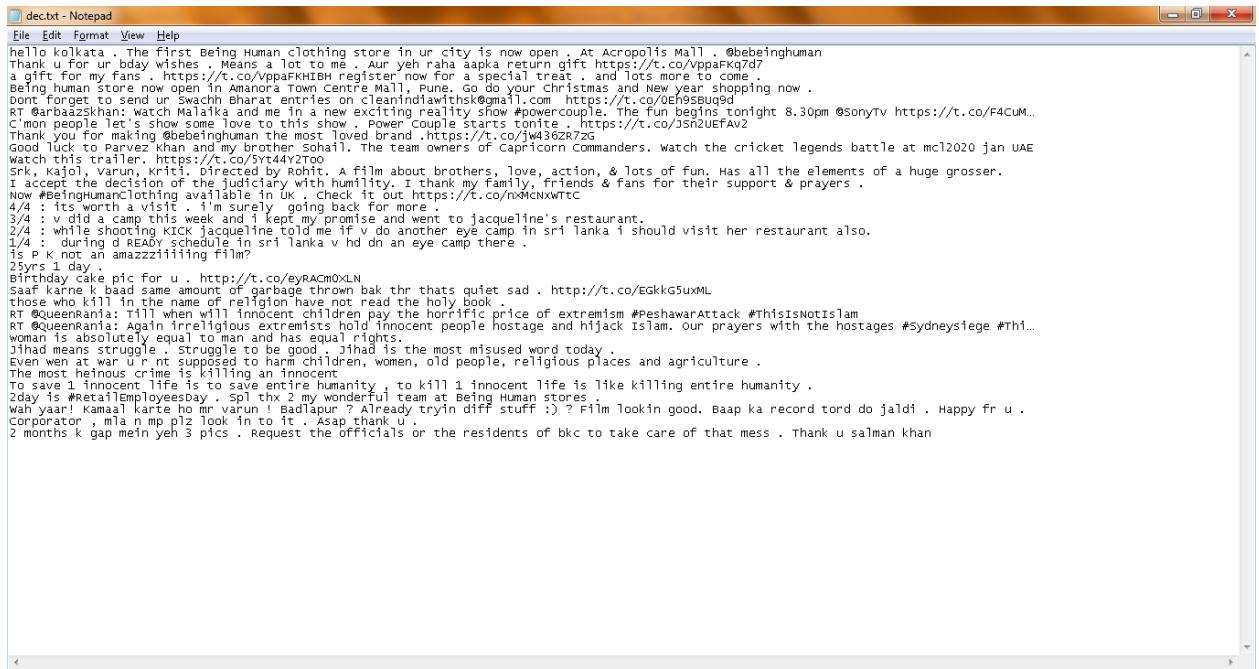
Figure 26: follower views in January

```

sep.txt - Notepad
File Edit Format View Help
all the best PC . do watch it guys . #quantico @priyankachopra
Oh acccha on 3rd oct .
Wen is pc's show releasing in India guys ?
Good afternoon everyone ! Suraj Singh's Prem again #PremRatanDhanPayo @rajshri @Foxstarhindi http://t.co/QA9IAa02k4
Create releasing on 25th Nov @theslystallone Looks Amazing as always https://t.co/h2xDQjljb
Some ppl click pics with me and then misuse them. This is not okay.
A fake Facebook page claims that I am casting for a film. Beware of fakes and rumors. Neither me nor my managers are casting for any project
RT @rajcheerfull: #BREAKING To all @biggboss & @BeingSalmanKhan fans here's the first look of #biggboss9 http://t.co/p6Qmmd0x0b Coming thi...
Thank you. Appreciate it. Retweeting my other fans. https://t.co/04qIKla57x
Wishing #teamIndia all the best for the #WC2015 . Heard the girls r playing very well @lumsoccer http://t.co/isNvY783ck
Signed up for tickets of #teamIndia vs #Ireland . Will sign more tomorrow. tag #HerokaTicketer http://t.co/yjfPNfgxr3
Plz vote for my friends Dimitri vegas and Like Mike as 1Artist/Group on http://t.co/ea0q7ycwOy
Thank u for all the compliments on suraj baba and athiya for songs, dances n action of hero. Appreciate it.
collections superb , whole team of hero super happy vit em. wd not have been possible vit out u going to see it. Thank u.
Thanks to the people who hv gone to see Hero n welcomed suraj n athiya in the film world n all those who r going to see it . Lots of respect
and if you don't know where the post office is , tol ask your mom.
This is the original way to mail. this is how it was before emails and messaging and all.
Send your ticket(s) to P. O. Box 960, Bandra (W), Mumbai - 400051. Apna Sam, address our twitter handle bhi likhna.
Watch Hero. Send me your ticket. I will sign 100 tickets. One of them could be yours.
Welcome them into the film world vit open arms vit enjoying them, supporting, encouraging them in cinemas vit cetees n taali's #herotomorrow
Friends families , family friends hai family. http://t.co/vwF0la74tL
Bigg boss season 9 is fine #biggboss9 promo shoot . @rajcheerfull @ColorSTV http://t.co/krg6f1xf2u
RT @rajcheerfull: colors welcomes #Bhajrangibhaijaan @BeingSalmanKhan to #biggboss9 Presented by @snapdeal Powered by @oppomobileindia thi...
Watch #herotomorrow at next track #khudaoutTomorrow at 12 pm. http://t.co/JRUf385cOE
Remembering London dreams with the roadies rockstar @annivijaysingha watch out for movie #3AM in cinemas tomorrow . http://t.co/jSYrCdkbYB
All d best @vinayvirmani24 & @ajayvirmani1 @DrCabbie
Hello Canada, Dr Cabbie releases today . njoy. http://t.co/juogymYkr
Check this out.http://t.co/ETBTU08x3

```

Figure 27: follower views in September



**Figure 28: follower views in December**

Fig 29:indicated the relative count of tweets of followers in various months.it indicates that tweet count is highest in January(646) and lowest in May(9).

The count of these months can now be plotted on plotly

```

Java - followerdetail/src/followersdetail.java - Eclipse
File Edit Run Source Refactor Navigate Search Project Window Help
Console
[terminated] count [Java Application] C:\Program Files\Java\jre7\bin\javaw.exe (Apr 23, 2016 12:08:49 AM)
Total number of lines in january : 646
Total number of lines in february: 10
Total number of lines in march: 29
Total number of lines in april: 13
Total number of lines in may : 9
Total number of lines in june : 57
Total number of lines in july : 61
Total number of lines in august: 44
Total number of lines in september : 33
Total number of lines in october: 50
Total number of lines in november : 28
Total number of lines in december: 34
done

```

**Figure 29:count of follower views in various months.**

## Storing the health related data in database and text files

Figure 30 describes the table storing the attributes of tweets extracted. The table is named twittergeo storing the user id, language in which tweets has been extracted, which in this case has been restricted to en (English), this merely acts a constraint. The next column is of geo which specifies the location of user from where it has been posted and year describes timestamp of the tweets extracted. Figure 31 is the magnified view to the twittergeo table shown in figure 30

The screenshot shows the Oracle Database Express Edition interface. The title bar says "ORACLE Database Express Edition". The top menu bar includes "Home", "Logout", and "Help". Below the menu is a toolbar with icons for search, refresh, and help. The main area shows a SQL command window with the text "desc twittergeo2". The results pane below displays the table structure:

Table	Column	Data Type	Length	Precision	Scale	Primary Key	Nullable	Default	Comment
TWITTERGEO2	<u>ID</u>	Varchar2	50	-	-	✓	-	-	
	<u>LANGUAGE</u>	Varchar2	20	-	-	✓	-	-	
	<u>GEO</u>	Varchar2	100	-	-	✓	-	-	
	<u>MONTH</u>	Varchar2	20	-	-	✓	-	-	
	<u>YEAR</u>	Varchar2	20	-	-	✓	-	-	
	<u>KEYWORD</u>	Varchar2	75	-	-	✓	-	-	

At the bottom right of the results pane, it says "1 - 6". The footer of the page includes "Results Explain Describe Saved SQL History" and "Object Type TABLE Object: TWITTERGEO2". At the very bottom, it says "Language: en" and "Application Express 2.1.0.0.39 Copyright © 1999, 2008, Oracle. All rights reserved."

Figure 30 :description of table storing the tweets the data type of entities.

Table	Column	Data Type	Length	Precision	Scale	Primary Key	Nullable	Default	Comment
TWITTERGEO2	<u>ID</u>	Varchar2	50	-	-	-	✓	-	-
	<u>LANGUAGE</u>	Varchar2	20	-	-	-	✓	-	-
	<u>GEO</u>	Varchar2	100	-	-	-	✓	-	-
	<u>MONTH</u>	Varchar2	20	-	-	-	✓	-	-
	<u>YEAR</u>	Varchar2	20	-	-	-	✓	-	-
	<u>KEYWORD</u>	Varchar2	75	-	-	-	✓	-	-

1 - 6

Figure 31: detailed description

Figure 32 reflects the data stored in twitter geo by selecting the rows from the table. For instance the first row describes the id and the language as en, geo as i=United Kingdom month as October year as 2015 and keyword cancer. Figure 33 is the magnified view of the result obtained by selecting the rows of twittergeo. Figure 34 shows the text file storage of tweets for cloud visualisation analysis.

ORACLE® Database Express Edition

User SYSTEM

Home > SQL > SQL Commands

Autocommit Display 5000 ▾

select \* from twittergeo2

Save Run

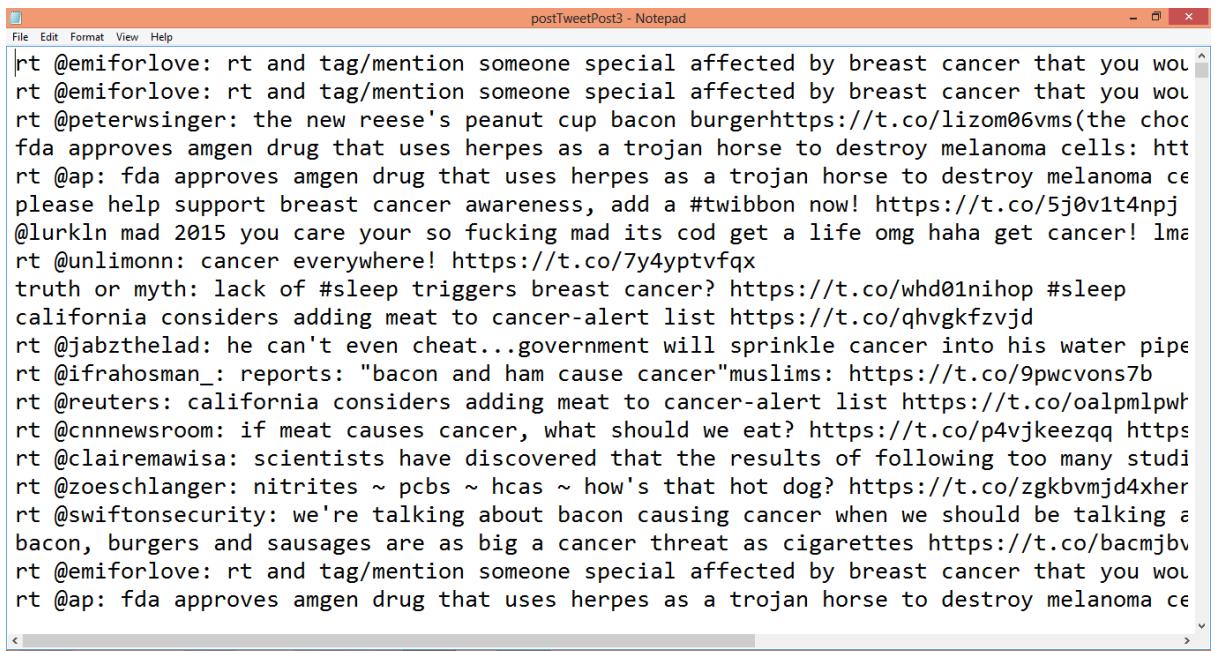
Results Explain Describe Saved SQL History

ID	LANGUAGE	GEO	MONTH	YEAR	KEYWORD
meanpeenmachine	en	-	October	2015	cancer
mattrees_	en	united kingdom	October	2015	cancer
xxmoonoofmylife	en	fuera de este mundo	October	2015	cancer
bifdusregularis	en	qarth	October	2015	cancer
thelaurrenj	en	united states	October	2015	cancer
jochurchill4	en	-	October	2015	cancer
_pisstance	en	-	October	2015	cancer
food_news_aus	en	-	October	2015	cancer
freedomforthwin	en	on my way to heaven	October	2015	cancer
cancer_advice	en	nyc	October	2015	cancer
_nikkias	en	-	October	2015	cancer
boldskyliving	en	bengaluru	October	2015	cancer

Figure 32: stored data in table twittergeo

ID	LANGUAGE	GEO	MONTH	YEAR	KEYWORD
meanpeenmachine	en	-	October	2015	cancer
mattrees_	en	united kingdom	October	2015	cancer
xxmoonoofmylife	en	fuera de este mundo	October	2015	cancer
bifdusregularis	en	qarth	October	2015	cancer
thelaurrenj	en	united states	October	2015	cancer
jochurchill4	en	-	October	2015	cancer
_pisstance	en	-	October	2015	cancer
food_news_aus	en	-	October	2015	cancer
freedomforthwin	en	on my way to heaven	October	2015	cancer
cancer_advice	en	nyc	October	2015	cancer
_nikkias	en	-	October	2015	cancer
boldskyliving	en	bengaluru	October	2015	cancer
amieynnn	en	peninsula → borneo	October	2015	cancer
homaihass	en	johor darul takzim	October	2015	cancer
tenaj_gerundio	en	cebu	October	2015	cancer
mrkeithreb	en	-	October	2015	cancer
navosnavos478	en	-	October	2015	cancer
so_tsuda	en	glasgow, scotland	October	2015	cancer
wilfreymorena	en	-	October	2015	cancer

Figure 33: detailed data stored in twittergeo



A screenshot of a Microsoft Notepad window titled "postTweetPost3 - Notepad". The window contains a large block of text representing a list of tweets. The tweets are in a conversational format, often starting with "rt @username:" followed by a message. Many tweets include URLs and hashtags. The text is in black font on a white background with standard Windows-style scroll bars.

```
rt @emiforlove: rt and tag/mention someone special affected by breast cancer that you wou
rt @emiforlove: rt and tag/mention someone special affected by breast cancer that you wou
rt @peterwsinger: the new reese's peanut cup bacon burgerhttps://t.co/lizom06vms(the choc
fda approves amgen drug that uses herpes as a trojan horse to destroy melanoma cells: htt
rt @ap: fda approves amgen drug that uses herpes as a trojan horse to destroy melanoma ce
please help support breast cancer awareness, add a #twibbon now! https://t.co/5j0v1t4npj
@lurkln mad 2015 you care your so fucking mad its cod get a life omg haha get cancer! lmao
rt @unlimonnn: cancer everywhere! https://t.co/7y4yptvfqx
truth or myth: lack of #sleep triggers breast cancer? https://t.co/whd01nihop #sleep
california considers adding meat to cancer-alert list https://t.co/qhvgkfzvjd
rt @jabzthelad: he can't even cheat...government will sprinkle cancer into his water pipe
rt @ifrahosman_: reports: "bacon and ham cause cancer"muslims: https://t.co/9pwcvons7b
rt @reuters: california considers adding meat to cancer-alert list https://t.co/oalpmplpwk
rt @cnnnewsroom: if meat causes cancer, what should we eat? https://t.co/p4vjkeezqq https://t
rt @clairemawisa: scientists have discovered that the results of following too many studi
rt @zoeschlanger: nitrates ~ pcbs ~ hcacs ~ how's that hot dog? https://t.co/zgkbvmjd4xher
rt @swiftonsecurity: we're talking about bacon causing cancer when we should be talking about
bacon, burgers and sausages are as big a cancer threat as cigarettes https://t.co/bacmjbv
rt @emiforlove: rt and tag/mention someone special affected by breast cancer that you wou
rt @ap: fda approves amgen drug that uses herpes as a trojan horse to destroy melanoma ce
```

Figure 34: text file storage for tweets

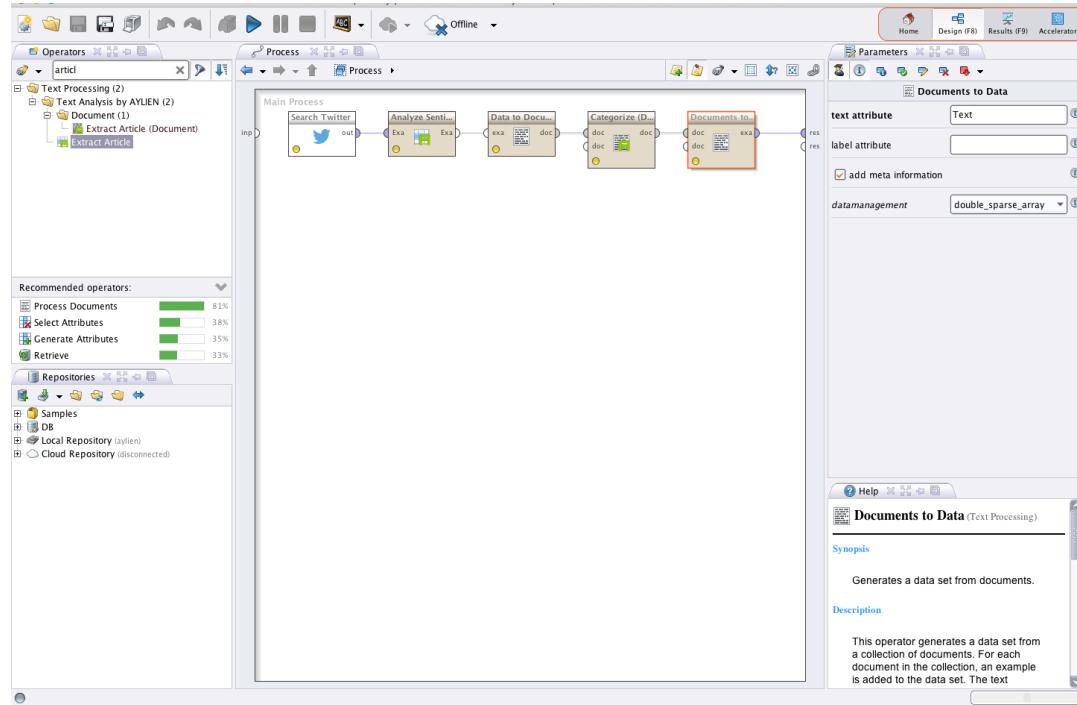
# CHAPTER 4

## REPRESENTATION AND ANALYSIS TOOLS

### 4.1 Graphical Analysis

The graphical analysis of Twitter Data Analytics is done to conclude results made on particular topic as per the 3 categories kept under consideration.

Figure 35 and 36 shows the process and output of tweet analysis in RapidMiner using Aylien Extension. The output gives the polarity and subjectivity confidence values.



**Figure 35: Process in rapidminer for analysing tweets**

Row No.	Created-At	From-User	From-User-Id	To-User	To-User-Id	Language	Source	Text	Geo-Locatio...	Geo
1	Feb 17, 2016 ...	Aayush Jain	63126829	?	-1	en	<a href="http://...	Freedom 251...	?	?
2	Feb 17, 2016 ...	The Social N...	529417759	?	-1	en	<a href="http://...	The way thing...	?	?
3	Feb 17, 2016 ...	LOOREY™	1205442050	?	-1	en	<a href="http://...	Freedom 251...	?	?
4	Feb 18, 2016 ...	Dopamine Drip	39472247	?	-1	en	<a href="http://...	Freedom 251...	?	?
5	Feb 18, 2016 ...	sabarirajaa	277428028	?	-1	en	<a href="http://...	Dear Girls ..	?	?
6	Feb 18, 2016 ...	Irfan Khan	4925056616	?	-1	en	<a href="http://...	FREEDOM S...	?	?
7	Feb 18, 2016 ...	The Bear Jew	4402296612	friendlii_gh...	588797722	en	<a href="http://...	@friendlii_gh...	?	?
8	Feb 17, 2016 ...	Aditya Choud...	3071646384	?	-1	en	<a href="http://...	Freedom 251...	?	?
9	Feb 17, 2016 ...	Hiya Jain	3090457028	?	-1	en	<a href="http://...	Freedom 251...	?	?
10	Feb 17, 2016 ...	Ritika Jeswani	3061329248	?	-1	en	<a href="http://...	Freedom 251...	?	?
11	Feb 17, 2016 ...	Tina Kohli	3041662826	?	-1	en	<a href="http://...	Freedom 251...	?	?
12	Feb 17, 2016 ...	Veer Gupta	2998304725	?	-1	en	<a href="http://...	Freedom 251...	?	?
13	Feb 17, 2016 ...	Pooja Agarwal	2988522253	?	-1	en	<a href="http://...	Freedom 251...	?	?
14	Feb 17, 2016 ...	Dhruv Hegde	3093045462	?	-1	en	<a href="http://...	Freedom 251...	?	?
15	Feb 17, 2016 ...	Megha Dua	3197246558	?	-1	en	<a href="http://...	Freedom 251...	?	?
16	Feb 17, 2016 ...	Reshmaa Th...	3141079791	?	-1	en	<a href="http://...	Freedom 251...	?	?

**Figure 36: output on tweet analysis**

#### 4.1.1 News Analysis

The graphs for this shows many valuable results along with some proofs.

1. **Polarity analysis** to show whether a particular news is positive,negative or neutral by showing the count & from that we can infer the thoughts of people.This analysis is based on ontological criteria of Aylien extension.

As in figure 37, count of polarity shows that people have neutral feeling(113 neutral tweets out of 158) about Freedom 251,hence they don't possess any liking or disliking to freedom 251 launch.

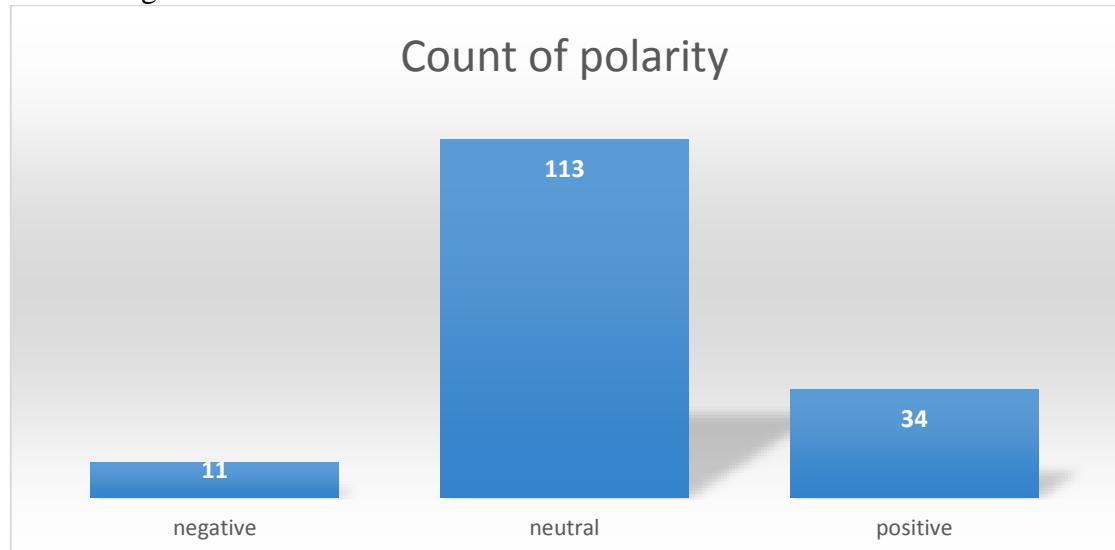


Figure 37: count of polarity

Proof:

- I. **method1** :applying ensemble technique-“vote”

Here the data mining techniques which are ensemble are decision tree, neural networks and LDA

Figure 38 shows the output of the decision tree interpreting the result as neutral. The output of LDA is shown in figure 39 inferring the support that freedom251 is a neutral news.The figure 40 shows the output of neural networks where the input is polarity\_confidence and subjectivity\_confidence and output nodes are positive,negative and neutral.The figure 41 shows the stacking model of the vote ensemble technique and the result shows the neutral support of the public regarding feedom251.

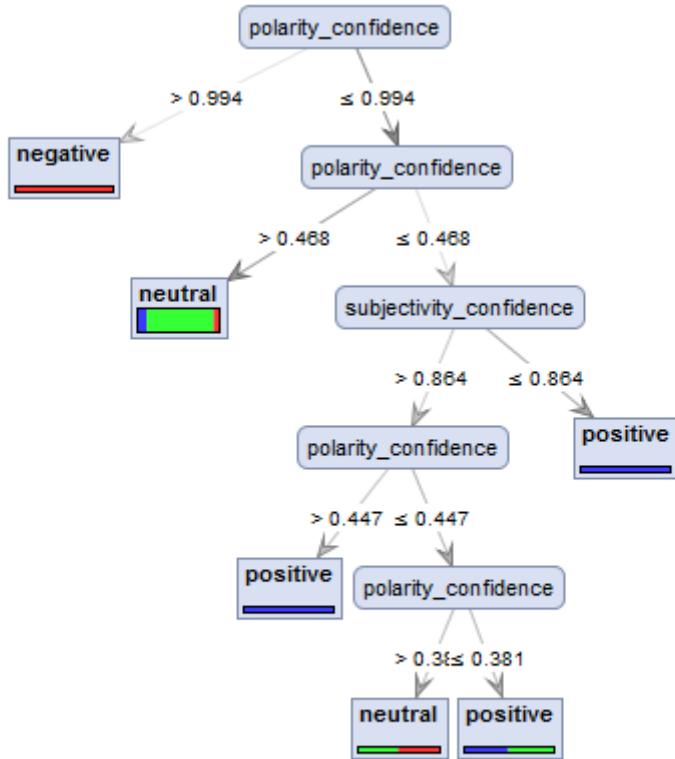
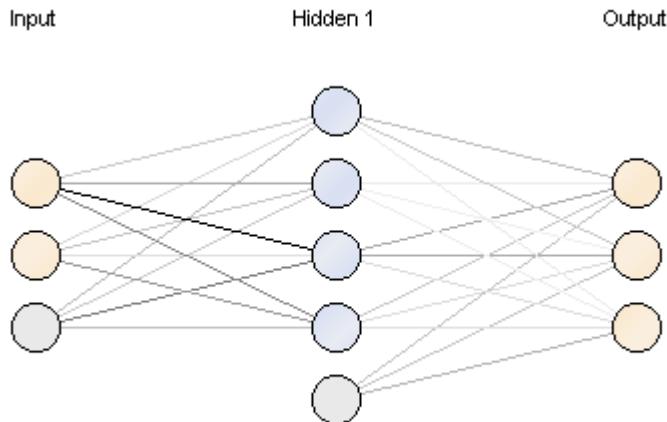


Figure 38: Output of decision tree

## Quadratic Discriminant Model

Apriori probabilities:  
 positive 0.2152  
 neutral 0.7152  
 negative 0.0696

Figure 39: Output of LDA



**Figure 40: Output of Neural Networks**

## AttributeBasedVoting

```
Using the majority of the following attributes for prediction:  

base_prediction0  

base_prediction1  

base_prediction2
```

The default value is neutral

**Figure 41: Overall output of vote ensemble technique**

- II. **method2:** retweet frequency count: through which it can be inferred that which particular tweet is tweeted highest and supports which kind of polarity.  
The figure 42 shows the highest counts for neutral polarity and hence again proves the neutrality of freedom251 news.

Text	Retweet-Cou	polarity
Freedom 251 launched as 'world's cheapest smartphone' at Rs 251https://t.co/W4tUNJr75	198.0	neutral
Here's an exclusive first look at Freedom 251, the Rs. 251 smartphone that's making the headlines today. #Freedom251 https://t.co/	186.0	neutral
Indian company Ringing Bells launches the world's cheapest smartphone Freedom 251 priced at ₹ 251 https://t.co/Llt1GN5oz8	104.0	neutral
Tommorow one new indian company Freedom is launching Lollipop Smartphone with quad core processor at Rs.251. This is Jalwa of 251₹ की कीमत के साथ FreeDom कम्पनी ने लॉन्च किया बेहतरीन फीचर्स वाला Android स्मार्टफोन.....https://t.co/hRPTCQGBFs https://t.	8.0	neutral
Aao Chutiyon Ye Dekho 251 wala Freedom Phone Ab Yahi Aukat Hai Tum Sabki https://t.co/gB98ayvQQm	5.0	neutral
Did any one booked freedom 251 ??	4.0	neutral
Freedom 251!!!!!!	2.0	neutral
The way things are going, I'm pretty sure tomorrow morning at 6 AM the freedom 251 website will crash ?? ;Rs 251 @makeinindia	1.0	positive
₹.251एके लाख में! भविष्युर्दृश्यमन्त्य शेषमंथलीhttps://t.co/JrLOFPwZ2 #Freedom251 https://t.co/wuwYOx5Pi5	1.0	neutral
के.पि. ओली भारत जान लागेकाछन, एक थान Freedom 251 पाइपलाईनबाट घर सम्म पूर्याउने ब्यवस्था मिलाइदिन भन्नु पन्यो :डी	1.0	neutral
I'm getting badly trolled on Twitter and Facebook because I called out the Freedom 251 phone and its Made in India claims. This is ho	1.0	negative
Freedom 251. Smartphone at Rs. 251. Pretty cool. World's cheapest phone ever.	.0	positive
Freedom 251 seems to be the Akash Tab of smartphones!	.0	positive
Dear Girls .. !!# Marry that Boy who was successfully Able to Place the Order of Freedom 251 .. !!;) ; # Patience Level _/\_	.0	positive
FREEDOM SMART DHONE AT 251 Than 2G Rorcharra १३०० +?????????????	0	positive

Figure 42: retweet count of freedom 251 posts

2. **Subjectivity analysis** to analyse the public opinion using object keywords like “RINGING BELLS”, “FREEDOM251”, “IPHONE”, “600000 REQUESTS “ETC. to discuss the generic v/s particular news support of the public. This analysis is based on ontological criteria of Aylien extension.  
 Figure 43 shows that out of 158 tweets, 85 are objective and 73 are subjective.

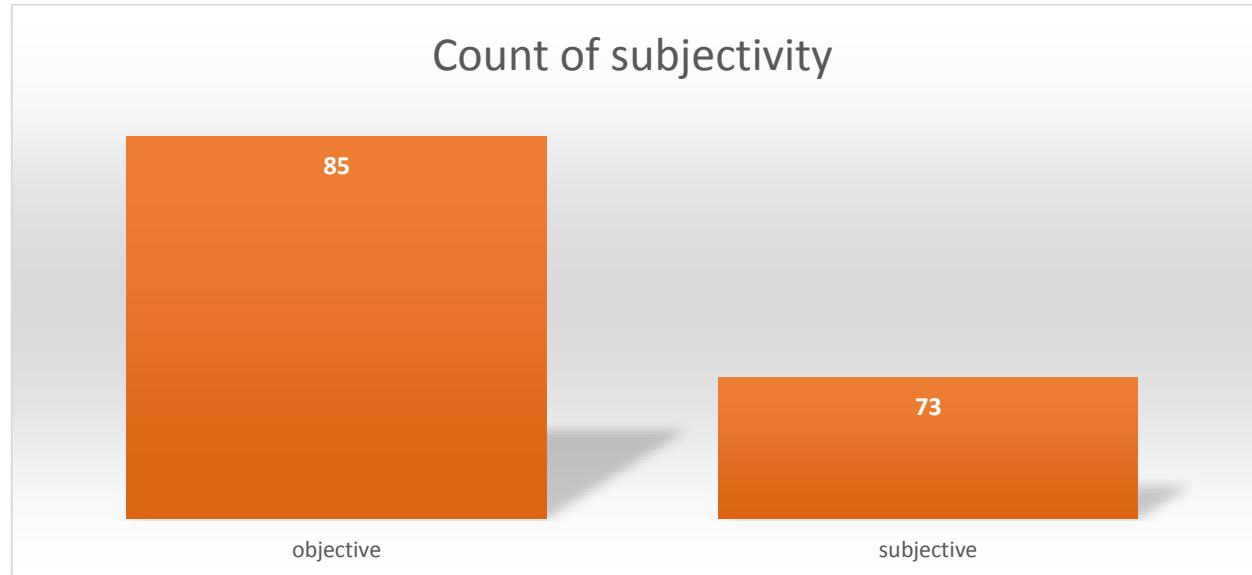
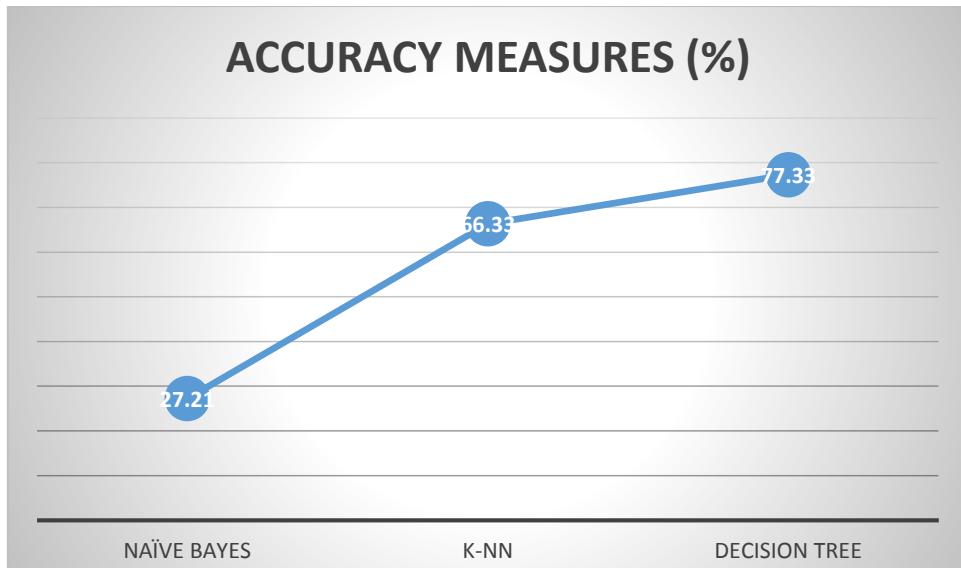


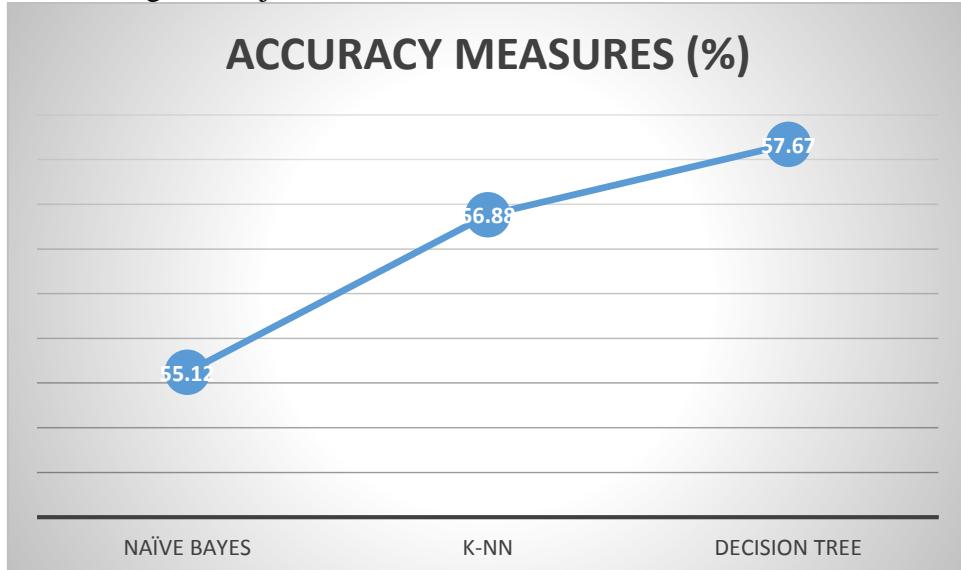
Figure 43: count of subjectivity

3. **ACCURACY ANALYSIS OF POLARITY BY COMPARING ACCURACY RATE OF DECISION TREE, NAÏVE BAYES AND K-NN**  
 The figure 44 shows that decision tree(77.33%) gives the highest accuracy in determining the neutral support of the tweets



**Figure 44:** % accuracy comparison of data mining algorithms for polarity confidence

4. ACCURACY ANALYSIS OF SUBJECTIVITY BY COMPARING ACCURACY RATE OF DECISION TREE, NAÏVE BAYES AND K-NN
- The figure 45 shows that decision tree(57.67%) gives the highest accuracy in determining the subjective nature of the tweets



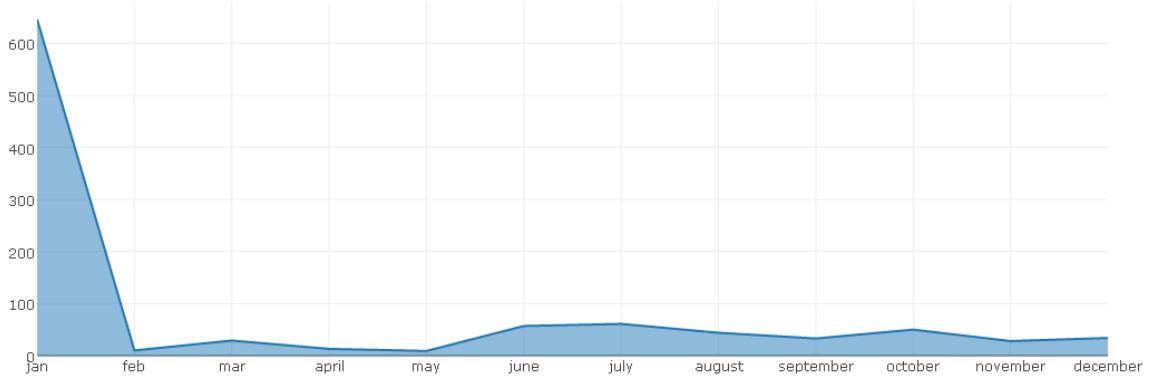
**Figure 45:** % accuracy comparison of data mining algorithms for subjectivity confidence

#### 4.1.2 Celebrity Analysis

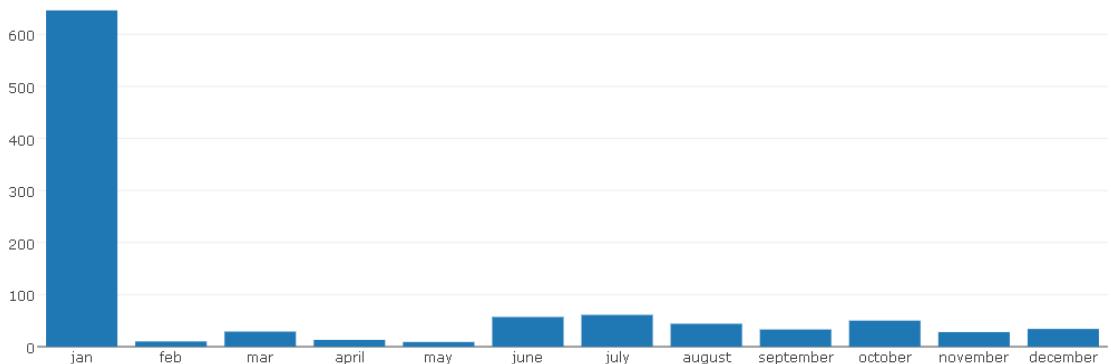
The tweet count of various followers in various months have been collected using java code which can be plotted using plotly.

Fig. 46 show that it was highest in January and steep fall in february which later on enhanced a little bit in march with sudden decrease in may and the almost a near about after july. Fig. 47 is an another representation of distributive which exactly displays the

relative amount by which the tweet count has been changing in various months i.e maximum in January and minimum in may.



**Figure 46:Distributive analysis of months vs follower tweet count**



**Figure 47:Bar Graph for month vs tweet count**

The graphical analysis further yield JSON code which can be used in future to build up a web service for online creating of graphs

```

• {
  o data:[
    ▪ {
      ▪ name:"MFEB",
      ▪ type:"bar",
      ▪ xsrc:"tanyabansal:3:DOBWYZIY0J3B4A0GDBEQTSLEF
        000EW86",
      ▪ ysrc:"tanyabansal:3:V43NE28ZRHDFOJJY8CR991WX41
        QXYYDU",
      ▪ uid:"dd6280"
    }
  ]
}
  
```

```

        }

    ○  ],
    ○  layout:{
        ■  yaxis:[
            ■  title:"",
            ■  type:"linear",
            ■  range:[
                ■  0,
                ■  680
            ],
            ■  autorange:true
        },
        ■  xaxis:{
            ■  title:"",
            ■  type:"category",
            ■  range:[
                ■  -0.5,
                ■  11.5
            ],
            ■  autorange:true
        },
        ■  height:497,
        ■  width:1121,
        ■  autosize:true,
        ■  dragmode:"zoom"
    }
}

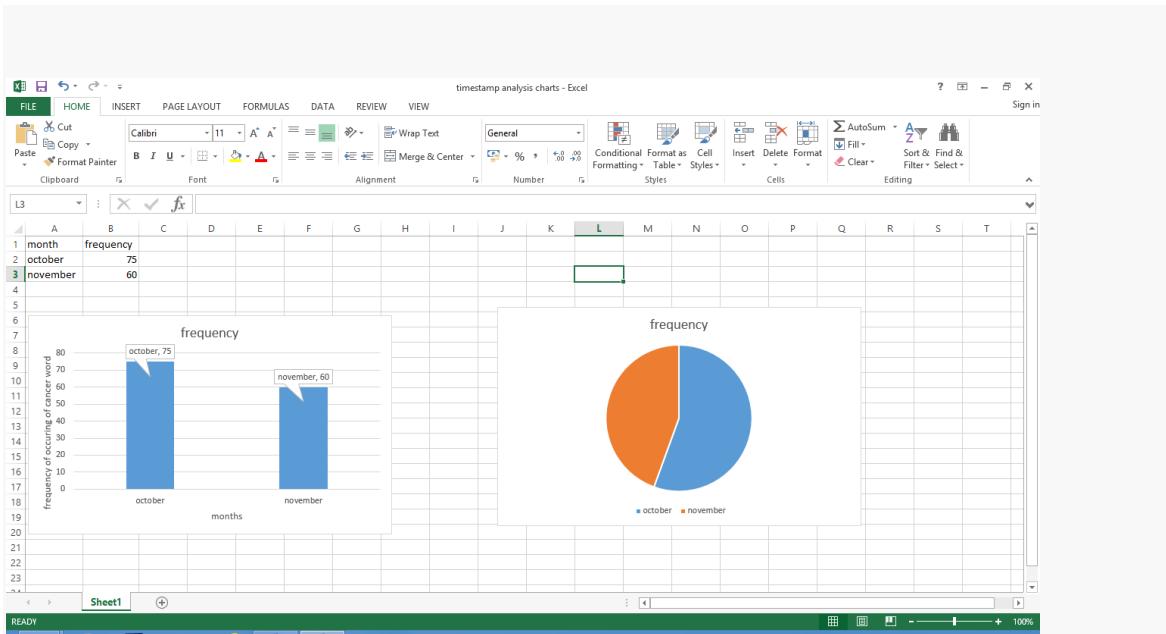
```

#### 4.1.3 Health Analysis

After storing the extracted attributes in the database, appropriate queries were used to calculate the frequency of tweets related to a particular disease for each month. This gave us an idea about the months in which the number of cases of a particular disease was prevalent. Also, seasons and months having chronic and acute disease cases can be later analysed to prevent such outbreak by taking appropriate measures beforehand.

The graphical analysis shows the frequency of tweets posted related to a particular keyword with respect to the months. Here in figure 48 frequency of cancer related posts was considered such that the cancer posts in the month of November (the blue part in pie

chart) were 60 and that in October (the orange part in pie chart) were 75, which means that the followers were more concerned about cancer in October than in November .also here in the fig48 only two months were considered since the tweets were collected in these two months only.



**Figure 48: Frequency of tweets**

## 4.2 Word Frequency Visualisation

Finally, we used cirrus tool to particularly analyse the data of tweet posts and extract the best known information from that. Cirrus is a visualization tool that displays a word cloud relating to the frequency of words appearing in one or more documents. One can click on any word appearing in the cloud to obtain detailed information about its relativity. Cirrus is a freely available visualization tool that generates a word cloud from a document or group of documents. The initial word cloud contains all of the most frequent words from the text, laid out graphically to fit together within an elliptical shape. Words may be oriented either vertically or horizontally, and they are rendered in different colours for ease of viewing. The user should be aware that word orientation and colour are strictly decorative elements, they do not indicate anything meaningful about the word. The size of the word, however, indicates the frequency with which it appears in the document. The larger the font size, the more frequent the word. The user can mouse over a word to see the precise number of times that it appears in the document.

One of the problems with word frequency visualization is the prevalence of conjunctions and articles in most texts. These less semantically interesting words often drown out the adjectives, adverbs, nouns, and verbs that are more likely to be of interest to the researcher. Cirrus offers an easy solution to this problem. By clicking the options icon in the upper right hand menu, a user can apply a list of stop words to the visualization. These stop words will be treated as noise, and stripped from the cloud, allowing words of greater interest to surface. Cirrus offers two pre-built lists of stop words, one for common English words, and one for common French words. The user can view the lists, so she know which words will be removed before applying them.

Cirrus offers users the ability to search within a word cloud for specific words, which is especially useful for seeing all of the word variations in a text. If I enter “comput” in the search box, Cirrus returns a word cloud indicating the word frequency of related terms like computers, computational, computing, computer-generated, computer-literate, microcomputers, human-computer, and supercomputing. It is also possible to enter comma separated terms in the search box in order to compare the frequency of multiple terms. By entering “comput, digit” for example, the user sees all of the variations of computer, as well as digital, digitization, etc. in the resulting word cloud. Any number of terms can be compared in this way. The Cirrus export menu exposes several options for sharing a word cloud. A permanent URL can be generated for a particular visualization, but the software also provides code snippets that allow the word cloud to be embedded in a website as a button or as a full visualization. Using this tool we basically analysed the certain important keywords used in a particular disease or for a person and its frequency of occurrence made our analysis criteria even better.

Once again, the example of Cancer is considered. Using this disease, further screenshots of the results obtained are given. These are the results obtained through Cirrus Tool. Using the word frequency visualization tool (here Cirrus), the keyword density of certain most frequently occurring words can help in analysing the words mostly tweeted and thereby helps know the trends through the size and frequency of its occurrence. Words are oriented either vertically or horizontally, and they are rendered in different colours just for ease of viewing and hence solves no further purpose regarding the position and placements as they do not indicate anything meaningful about the word. The size of the word, however, indicates the frequency with which it appears in the document. The larger the font size, the more frequent the word. The user can mouse over a word to see the precise number of times that it appears in the document. Here in figure 49 the word “cancer” is given the highest size which means that its occurrence in the tweets is the highest followed by other words which helps to know the keywords related to a particular disease, as an example see in figure 49 that words like awareness, breast, coffee etc. are also present which helps to relate these keywords to the disease specified in terms of its symptoms, type of disease and its treatments.

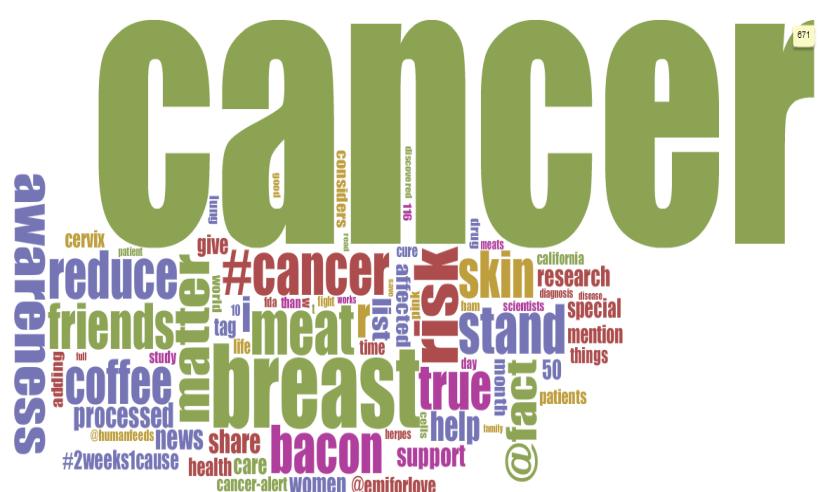
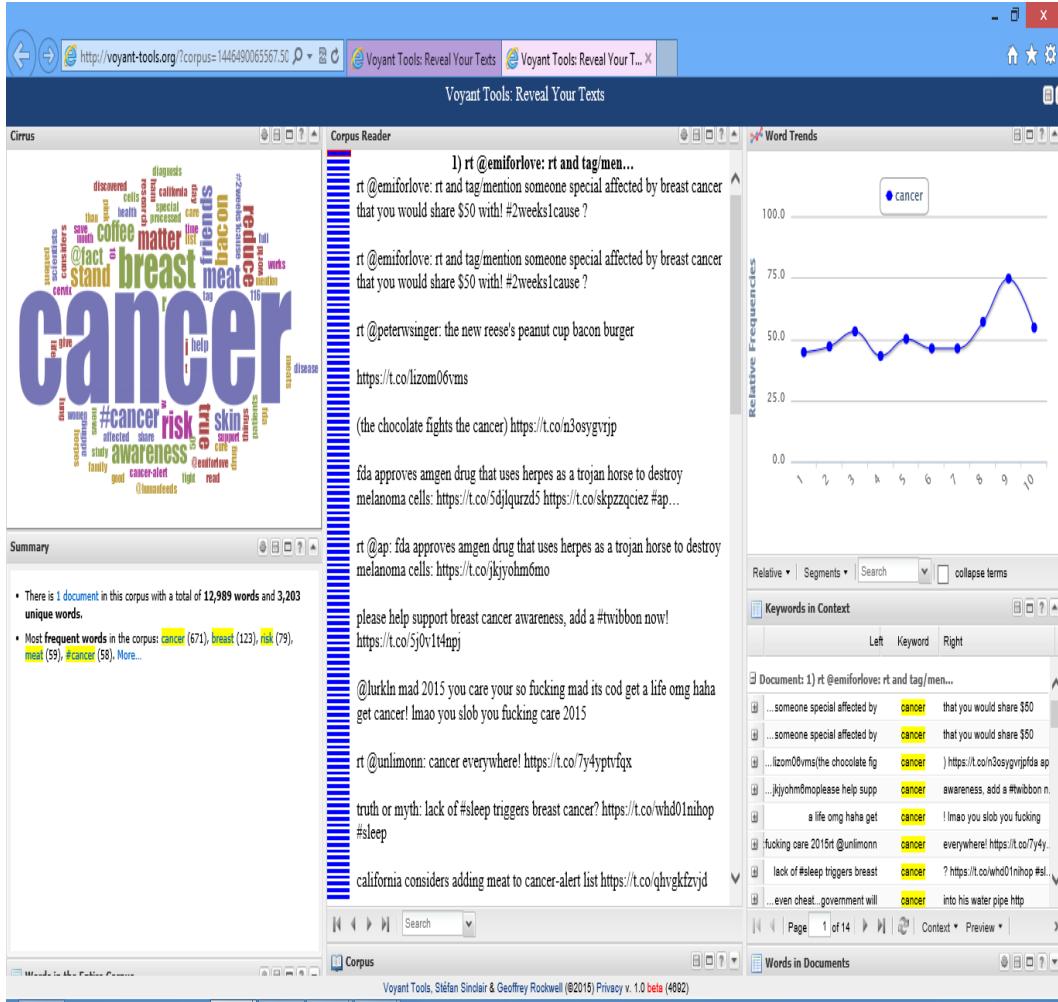


Figure 49: Word Cloud for extracted tweets

On clicking any particular word in the cloud (here cancer)we can get detailed analysis of that word in terms of “corpus reader”, “word trend”, “keywords in context” and “summary”. In the summary part (which is bottom left in figure 50),the complete summary of all the most frequently occurring words is given in key value pair where the word name is specified along with its frequency in brackets. The word trends (which is top right in figure 50)shows the graphs of occurrence of the word in total number of tweets .Corpus reader and keywords in context shows the tweets where the word occurs.



**Figure 50: Cirrus Tool**

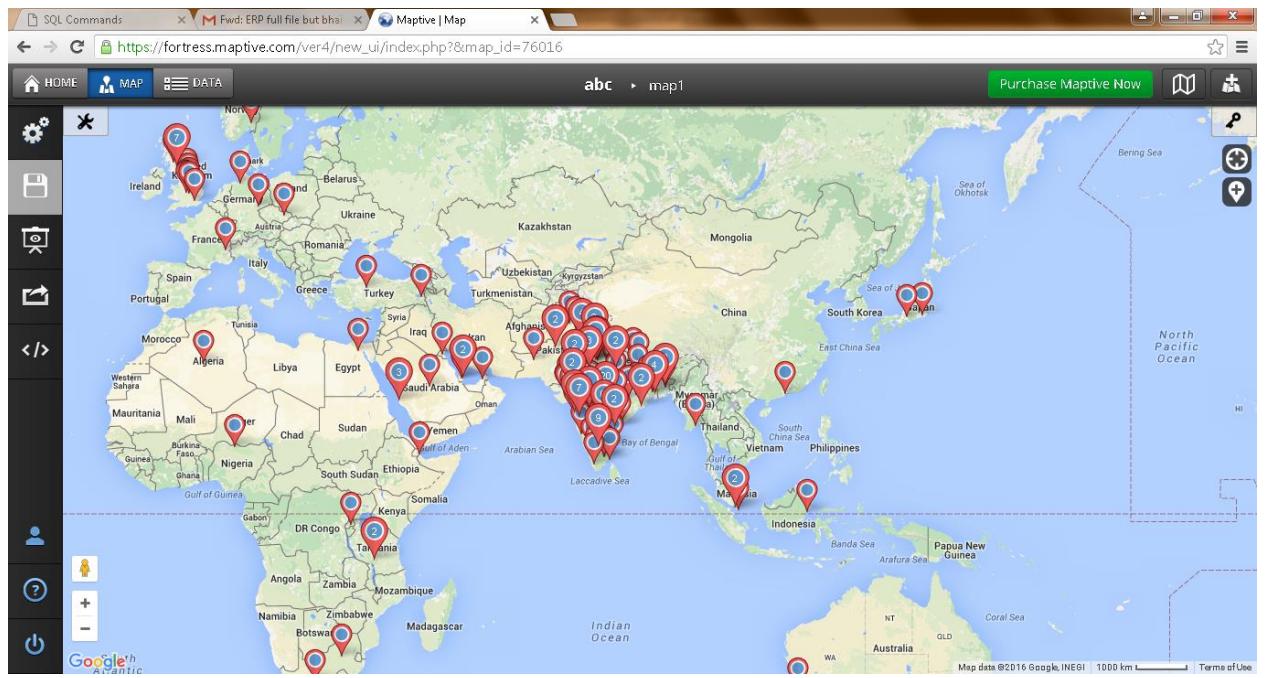
## 4.3 Geo Mapping

Geo mapping in this project is done both dynamically and statically. In celebrity analysis we locate the locations in dynamic way whereas in health analysis we used static way

### 4.3.1 Dynamic Mapping

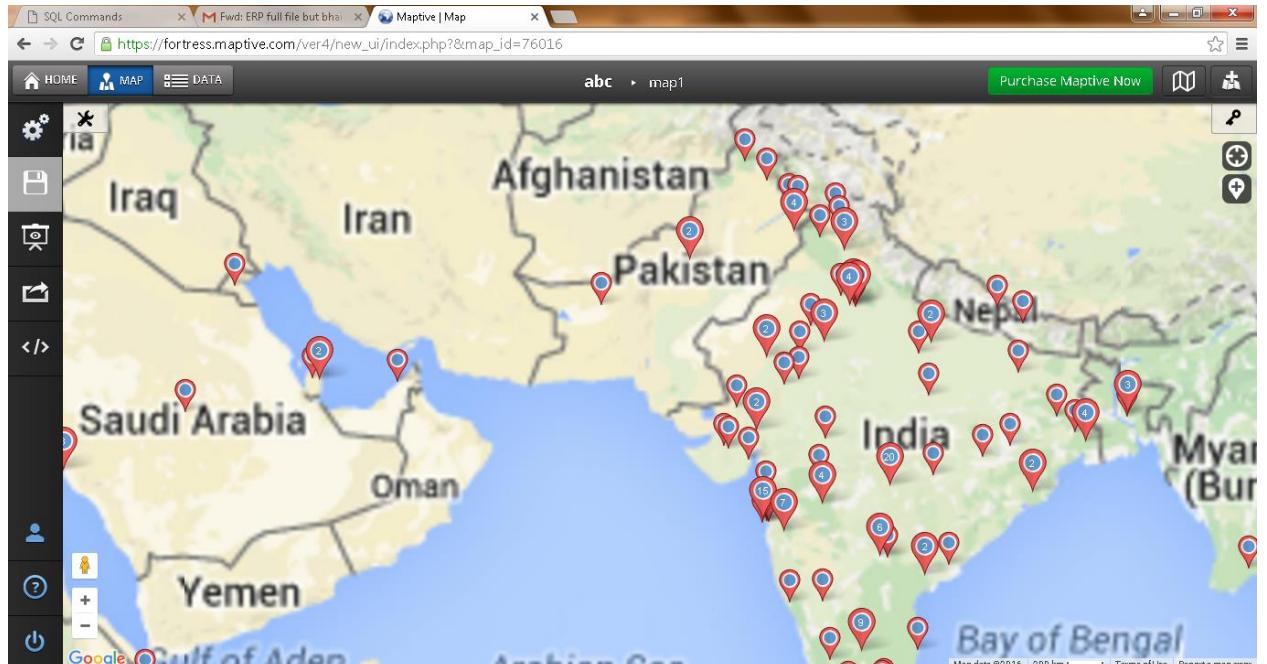
The follower location which was extracted can be geolocated using maptive a tool.

Fig.51 shows the geolocation of followers located in different parts of the world. The label represents their position and the number on them represents the frequency of that particular region like 2 for ahmedabad



**Figure 51:geolocation of followers**

Fig. 52 represents the zoomed view of maps indicating that many of the followers are in india as compared to others.pakistan has less number.moreover terror countries like iran has no followers reflecting salman khan a better image. Fig. 53 represents that india is the country which is most popular among salman khan followers



**Figure 52:close view of geolocation**

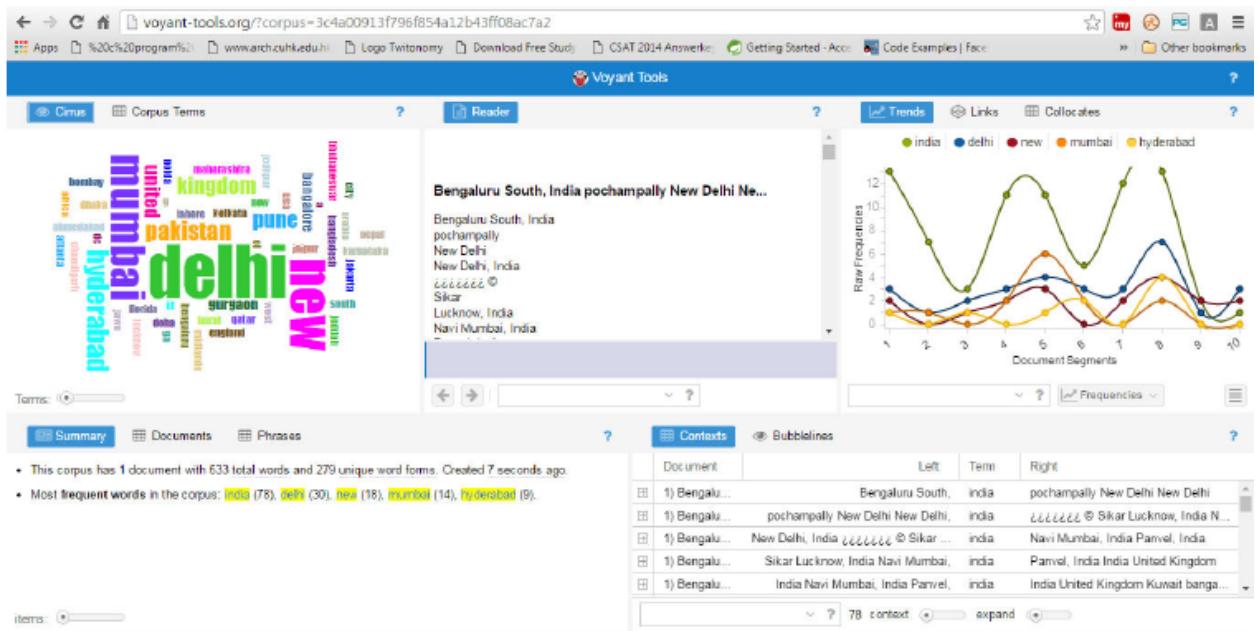


Figure 53:country wise analysis of followers

#### 4.3.2 Static Mapping

By using the database, a tweet frequency is calculated for each country for a particular disease. This calculated frequency is then used by JxBrowser and Google API service of maps marker for Geo-Mapping.

Figure 54 shows one such geo mapping. Balloon marks (marker position) represents the frequency of tweets for that particular area for example France: 1, marker has been placed at France and 1 represents its frequency, these marker has been placed on basis of their longitude and latitude. A tooltip is associated with one of the balloon representing the frequency of extracted tweets regarding a particular disease like 1 in case of France. Figure 55 shows the zoomed image for the same.

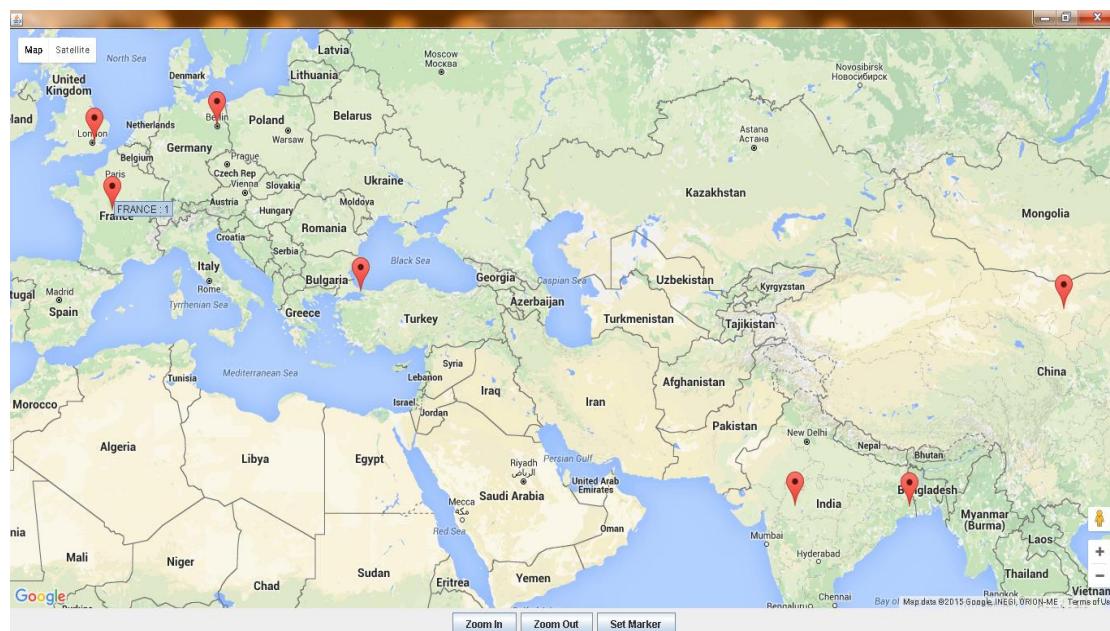
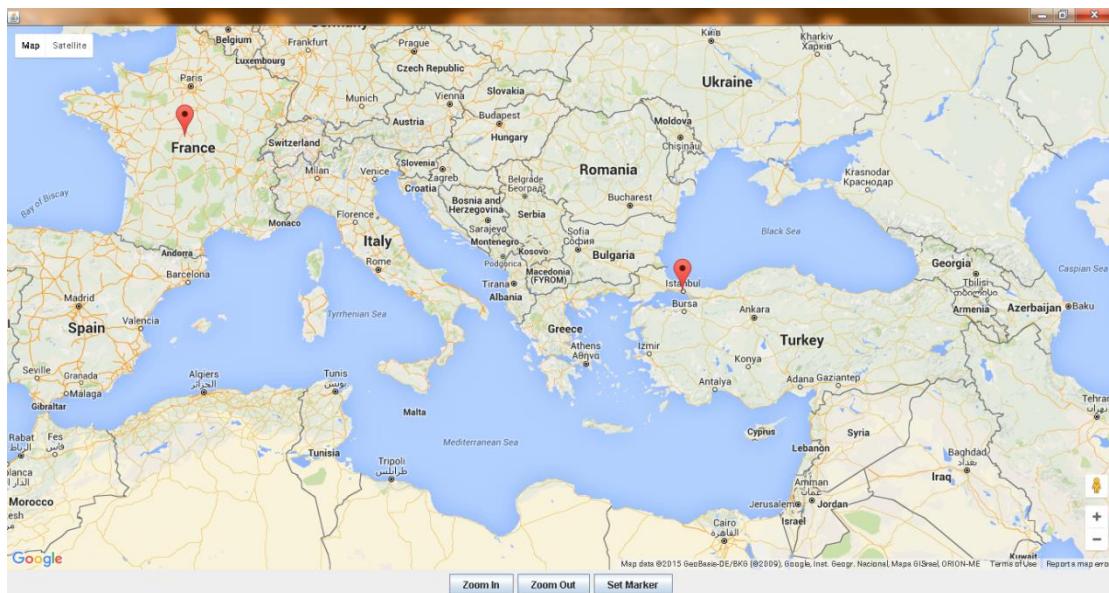


Figure 54: Geo Mapping



**Figure 15: Zoomed in view**

## **CHAPTER 5**

### **CONCLUSION, SUMMARY AND FUTURE SCOPE**

We designed and implemented a behavior analyzing system that uses twitter data to analyse behavior in real time. Twitter data is continuously downloaded using Twitter streaming API. Then the tweet texts are extracted along with the time stamps, and user locations and stored in a database, which is used for further analysis in the geographical and text terms. All the output data is visualized as an interactive map, bar charts, and word clouds. Our aim was to analyze the data present on social media. People often post messages about interesting news, their daily activities, thoughts and feelings. This may help to analyze their behavior and behavior of other peoples towards him/her. It makes easy to analyze positive and negative attitude from their tweets Twitter message, when used in combination, helps increase the accuracy of our prediction.

This project includes the extraction of tweets using Twitter API and then analyzing them in many possible ways. Apart from this, tools named RapidMiner, Cirrus, and maptive, plotly are used to take the analyzing process to another step.

We designed and implemented a novel disease surveillance system that uses twitter data to automatically track a particular disease activity in real time. We are also interested in using the social network information (e.g., friends and followers network) to predict the disease outbreak. We plan to gather additional online data (e.g., Facebook, blogs data, news feeds) for real-time disease surveillance. We set criteria to group certain symptoms and treatments in order to find the proximity of prevailing diseases and calamities

We will explore this project by studying more data mining tools such as twitonomy. And try to implement data mining algorithms such as Expectation Maximization (EM), Page Rank, AdaBoost and compare their accuracy with the other one.

Another possible research direction is to further integrate the modules into one Complete using Hadoop or R-programming. We may try to use Facebook as a data source to implement our analysis.

## REFERENCES

1. Bo Pang and Lillian Lee. "A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts." 2004. Proceedings of ACL, pp. 271–278.
2. Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. "Thumbs up? Sentiment Classification using Machine Learning Techniques." 2002. Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 79–86.
3. Barbosa, L., and Feng, J. 2010. Robust sentiment detection on twitter from biased and noisy data. In COLING.
4. T. Sakaki, M. Okazaki, and Y. Matsuo. Earthquake shakes twitter users: real-time event detection by social sensors. In Proceedings of the 19th international conference on World wide web, WWW '10, pages 851
5. Blei, D.; Ng, A.; and Jordan, M. 2003. Latent dirichlet allocation. Journal of Machine Learning Research (JMLR) 3.
6. Carneiro, H., and Mylonakis, E. 2009. Google trends: a web-based tool for real-time surveillance of disease outbreaks. Clin Infect Dis 49(10):1557–64.
7. Chew, C., and Eysenbach, G. 2010. Pandemics in the age of twitter: Content analysis of tweets during the 2009 h1n1 outbreak. PLoS ONE 5(11):e14118.
8. Culotta, A. 2010a. Detecting influenza epidemics by analyzing twitter messages. arXiv:1007.4748v1 [cs.IR].
9. Culotta, A. 2010b. Towards detecting influenza epidemics by analyzing twitter messages. In KDD Workshop on Social Media Analytics.
10. Eisenstein, J.; O'Connor, B.; Smith, N. A.; and Xing, E. P. 2010. A latent variable model for geographic lexical variation. In Empirical Natural Language Processing Conference (EMNLP).
11. "Tweet Tracker: An Analysis Tool for Humanitarian and Disaster Relief". *The 5th International AAAI Conference on Weblogs and Social Media*.
12. Eysenbach, G. 2009. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the internet. J Med Internet Res 11(1):e11.

13. T. Wilson, J. Wiebe, and P. Hoffmann. Recognizing contextual polarity in phrase-level sentiment analysis. In Proceedings of Human Language Technologies Conference/Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP 2005), Vancouver, CA, 2005.
14. Stopwords, Filtering and Data Sparsity for Sentiment Analysis of Twitter by Hassan Saif, Miriam Fernandez, Yulan He, Harith Alani Knowledge Media Institute, The Open University, UK
15. The Porter stemming algorithm: then and now. by Willett, P. (2006)a Program: electronic library and information systems, 40 (3). pp. 219-223.
16. J. Ginsberg, M.H. Mohebbi, R.S. Patel, L. Brammer, M.S. Smolinski, and L. Brilliant. 2008. Detecting influenza epidemics using search engine query data. *Nature*, 457(7232):1012–1014.
17. A Survey of Algorithms for Keyword Search on Graph Data -By Haixun Wang and Charu C. Aggarwal
18. Vasileios Lampos and Nello Cristianini. 2010. Tracking the flu pandemic by monitoring the social web. In IAPR 2nd Workshop on Cognitive Information Processing (CIP 2010).
19. Google flu trends. <http://www.google.org/> flutrends, 2012.
20. H. Achrekar, A. Gandhe, R. Lazarus, S.-H. Yu, and B. Liu. Predicting flu trends using twitter data. In Computer Communications Workshops (INFOCOM WKSHPS), 2011 IEEE Conference on, 2011.
21. M. J. Paul and M. Dredze. You are what you tweet: Analyzing twitter for public health. In L. A. Adamic, R. A. Baeza-Yates, and S. Counts, editors, ICWSM. The AAAI Press, 2011
22. Google chart tools. <https://developers.google.com/chart/interactive/docs/gallery>.
23. Signorini, A. M. Segre, and P. M. Polgreen. The use of twitter to track levels of disease activity and public concern in the u.s. during the influenza a h1n1 pandemic. *PLoS ONE*, 6(5):e19467, 05 2011.