

# Opinion Mining With Information Extraction

<sup>1</sup>Dr. Prerna Mahajan,<sup>2</sup>Geetika Gaba

<sup>1</sup>Head Of Department,,<sup>2</sup>Associate Software Engineer  
Information Technology,<sup>1</sup>IITM,<sup>2</sup>Orane Consulting Pvt Ltd,New Delhi,India

---

**Abstract:** The aim of this paper is to discuss and analyze a new and innovative research challenge on opinion mining. This research challenge has been developed in opinion analysis of social community and their opinions on different topic. In this paper we are meant to develop a system that can identify and analyze opinion as represented in the electronic text. People are more eager to express and share their views on any topic regarding day-to-day activities and social issues as well in this paper a precise method for predicting people sentiments is used that enable us to extract opinion from the interviews, survey and internet that predict reviews of each and every person, which could prove valuable for social growth and research. So, current scenario needs to be analyzed with the help of interviews, blogs/forums and natural language processing. The results further validate and support the main issue of paper.

---

## 1 Introduction

Generally we like to ask other's opinion when we make opinion on something .The area of opinion analysis in general. (Zhen et al,2014). A lot of people have actually studied the task of information extraction, but the area of opinion mining is actually much broader than that. Here are some fascinating tasks related to opinion analysis which I have discovered over the course of my survey over a known issue of India and Pakistan which is running since 1947.Opinion mining involves in making a system to gather and analyze opinions about issue through information extraction with the help of interviews, blogs, forums, tweets, comments and reviews. Opinion mining can be useful in several ways to judge person's sentiments and to predict future expectations. Basically opinion mining focuses to identify the attitude or sentiments of a speaker or a writer with respect to some topic or the overall relative orientation of a document. When we talk about opinion mining we have two aspects in it, one is opinion lexicon and other is opinion targets .An opinion lexicon is group of opinion words such as good, bad, positive, negative opinions about any topic. The aim of opinion mining is related to the topics on which opinions are expressed. It has also become useful socially and economically to create summaries of person's sentiments and opinions that consist of subjective sentiments extracted from reviews or from a reviewer's point that is positive or negative. (Jin and Ho,2009) Nowadays, current system needs to be analyzed using comments and natural language processing. In this paper the importance of social network and socialization is preferred for opinion mining and sentiment analysis In this paper the related work and definitions used in opinion mining and through information extraction is presented.

## 2 Related work

The document, text, blog and twitter data can be analyzed for opinion and sentiment analysis with the help of information extraction .The document level means sentence level data in which text data has been analyzed which is used to examine the opinions that have been generated with the help of data source. (Jakob and Gurevych,2010)Automated methods for information analysis have been widely used and have increasingly using from many years. The research theme is based on large and established computer science domains, like Natural Language Processing, Text Mining, and Automated Content Analysis etc. What is trending today is the rapid increase in the quality of data extracted from large amount of data, basically due to adoption of Information technology, that are available for machine learning algorithm and scripts that are trained.(Wu and Shen,2009)

Social media content by nature reflects opinions and sentiments, while older text analysis tended to focus on identifying topics As such; it deals with more complicated problems. Because of the increase in the width of data available and more complicated concepts to analyze, in recent years there has been a decrease in interest on traditional analytical based application, and a move towards greater use of statistics and visualization. Like any other discipline, also automated content analysis is becoming a data-intensive science.

## 3 Literature review

Wei Jin and Hung Hay Ho introduced machine learning approach for web opinion mining and information extraction. This approach gives output for server issues that have not been introduced in earlier approaches. (Jin and Ho,2009) This system can fetch new vocabularies based on the pattern, which is implemented through text and web mining. Here bootstrapping approach is used to conquer scenario in which gathering a large set could be expensive and difficult to fulfill. In this paper the efficiency of introduced approach in web opinion mining and extraction from product review are implemented in result.

Guang Qiu, Bing Liu, Jiajun Bu and Chun Chen emphasis on to main tasks in opinion mining that are opinion lexicon and information extraction. In this paper, an approach is used to extract opinion vocabulary and information iteratively. (Qiu et al,2015) The relations that are resulted between opinion words and targets are used for extraction a method for new opinion words polarity and domain priming is introduced. The outcome of this paper shows that this approach exceeds point of view of other approaches. Here author examine the relation between domain detection and words polarity classification. The detection of domain related words can compress reviews in shorter stanzas and still retains polarity at a level that can be compared to that of the full review.

Forman discussed about polarity classifier that extracts the subjects that are shown to be more effective than the original document. The research work shows that the reduction in framework results in the emergence of efficient algorithm for opinion analysis. (Forman,2004) Through this framework information can lead to statistical improvement in polarity with classified and accurate information.

(Jacob and Gurevych,2010) It shows how a CRF-based approach for opinion mining in a sole and multi domain. In this paper a relative assessment of this approach on data set from four different subjects are presented. The CRF-based approach emphasizes on a organized baseline on all dataset on a single domain .it also concludes promising results in the multi domain.

Maniquing Hu and Bing Liu proposed multiple techniques have been proposed for feature extraction from product reviews based on data mining and natural language processing technique. To make feature based summary of a number of customers review on a product sold online is the main target of this paper. (Qui et al,2015) Opinion mining becomes growingly important as increased number of people are buying and expressing their opinion on the web as a form of reviews. Innovative result of this paper shows that the proposed techniques are effective in implement tasks. Opinion mining based on document, sentence, phrase or word level does not represent what exactly a person like or dislike.

Recently, there have been many researches to monitor public opinion and social trends(Akcora et al,2010; Diakopoulos and Shamma ,2010). They include election prediction using Twitter data ( Boutet and Yoneki ,2010), monitoring of customer sentiment on a certain brand (Liu et al,2007), movie performance prediction using Twitter (Baek et al,2014; Rui and Whinston,2013), disease and disaster tracking using Internet information(Sakaki, Okazaki and Matsuo,2010) and unemployment benefit prediction using Internet search information( D'Amuri and Marcucci,2009). Especially, public opinion monitoring is useful in sensing public opinion trends and reduction of potential social risks and conflicts(Lee et al,2008).

This segment portrays related work via web-based networking media use by scholarly libraries, and related work on Twitter, and content examination. Online networking ADOPTION BY ACADEMIC LIBRARIES Twitter and Facebook are the most ordinarily watched web-based social networking applications utilized by libraries (Palmer, 2014), which reflects the utilization of Twitter and Facebook as the most mainstream web-based social networking devices utilized by the overall population. One of the key components of these apparatuses is that they empower two-way correspondence and collaboration between a library and its benefactors. This sort of discourse regularly includes staff and benefactors examining books inside the accumulation, or staff reacting to general reference request. Scholarly libraries "utilizing them as effort apparatuses, a technique for advancing themselves and their administrations inside their groups and past" (Boateng and Quan Liu, 2014, p. 126). The capacity to extend client connection through online networking stages has to a great extent been viewed as a win for scholarly libraries, as they can possibly "take bookkeepers and an institutional nearness to the spots where the benefactors are, crossing over any barrier, and endeavoring to begin new discussions" (Gaha and Hall, 2015, p. 49). Jain (2014) likewise talked about the capability of online networking in showcasing library and data administrations. Twitter, specifically, "has been observed to be valuable in amplifying the scope of library and data administrations, too an exhibiting an a la mode picture of administrations arrangement" (Loudon and Hall, 2010, p.240). As indicated by Cuddy, Graham, and Morton-Owens (2010), Twitter stage additionally enables libraries to share news about offices, assets, downtime, occasions, and staff. The substance of a few tweets incorporates sees about loud development close to the library, and updates about administrations the library needs to advance, for example, text (IM) reference (Cuddy, Graham, and Morton-Owens, 2010). Loudon and Hall (2010) state that these sorts of tweets are gone for conveying news and improving current mindfulness identified with library occasions, assets and administrations, and changes to hours of operation. Del Bosque et al. (2012) found that scholastic libraries tweet about library assets, hours of operation, library occasions, and grounds and group occasions also. In maintaining and in addition enhancing their web-based social networking nearness, scholarly libraries frequently confront various difficulties that must be tended to (Ramsey and Vecchione, 2014; Zohoorian-Fooladi and Abrizah, 2014). Drawing in "adherents" for the library's Twitter sustains, for instance, was accounted for to be troublesome. In an examination by Chu and Du (2013), it was discovered that understudies once in a while add to long range informal communication devices utilized by libraries, and it was discovered that understudies disliked utilizing them. Also, Stuart (2010) found that libraries regularly disregard the social parts of Twitter and don't utilize it as a chance to associate with their supporters. They are "endeavoring to fabricate crowds on various online networking locales by making drawing in content and surveying their web-based social networking nearness" (Burgert, Nann, and Sterling , 2014, p. 1). Loudon and Hall (2010) guaranteed that there is by all accounts little direction or support from administration in selection of online networking in scholastic libraries. Or maybe, libraries "are for the most part following up on their own drive with regards to contriving and actualizing, a technique for the sending of Twitter in the work environment" (Loudon and Hall, 2010, p. 239). The constrained support of clients was likewise evident to be an obstruction for libraries (Chu and Du, 2013). To conquer these challenges and draw in library supporters to their online networking nearness, Burgert et al. (2014, p.3) showed that "it is essential to keep content new, use associations with grounds

accomplices, and generally announce the library's online networking nearness". Twitter account holders at the library ought to be empowering two-way correspondence (Cuddy et al., 2010). To do this, "custodians need to give predictable updates a congenial, easygoing tone, and tune in to their clients on the web-based social networking destinations" (Burgert et al., 2014, p. 6). Kim, Abels, and Yang (2012), p. 1, detailed that "some scholastic libraries have no unmistakable distributed destinations for utilizing Twitter or other online networking." Therefore, content mining could be a potential arrangement toward understanding the best substance sorts posted by libraries and their clients, with an end goal to fulfill the requirements of these clients and successfully draw in them in utilizing library administrations.

#### 4 Opinion mining

Opinion mining is a kind of natural language processing judging the views of a public on any kind of issue.(Sannella et al,2006) Opinion mining is also called sentiment analysis involves making a system to gather and analyze opinion about any topic. Automated opinion mining frequently uses machine learning, a type of artificial intelligence to mine opinions from text. Opinion mining can be useful in many ways. It can help people in many ways, In this paper as India and Pakistan issue is taken so it will help socially and economically to know the judgment of public over this issue which is taken in consideration from many years since 1947. Being able to identify this kind of information in a systematic ways gives a particular conclusion and it also gives a clear picture of public opinion. There are several challenges that have been considered in opinion mining. (Forman,2004)The word that is considered to be positive in one scenario may be considered negative in other scenario. In this paper the people who are saying that the war between India and Pakistan should take place are also saying that the war should take place for sake of our future generation so that this all issue should come to an end and many people are saying the war should take place for the sake of taking revenge from Pakistan. The intentions in both the scenario is quite much different but the statement is same.

The second challenge is that people don't always express the same way. (Esuli and Sebastiani,2005) As some times the person's intentions are positive but the statement he/she is saying depicts negative image. It is basically finding the opinion of the text with the use of natural language processing and text analysis to identify and extract the opinions from text that have been spoken by public. Opinion mining is widely applied to reviews for a range of applications, ranging among social, economical, marketing etc.

Generally speaking, opinion mining focuses on determining the attitude of speaker or a writer with respect to come topic or the whole contextual polarity of document. The attitude with the way he is speaking may be his/her judgment or evaluation after analyzing the situation is determined

In opinion mining the method for determining opinions is the use of scaling system by words commonly linked with having a negative, neutral or positive opinion with them are given a number on a -10 to +10 scale. When a piece of text is analyzed with the help of natural language processing, the environment is given a score with each aspect of opinions. (Bai and Padman,2006) This allows the movement to a more specified understanding of intentions. It mainly concentrates on attitudes, but previously traditional text mining focused on the analysis of facts and it deals with classifying entire text on the basis opinions towards certain domain.

#### 5 Information Extraction

Method of filtering information from large volume of text is known as information extraction, refining of information from large text and tagging of particular terms. In this paper we have mined the opinions from large text of reviews which have been collected through various interviews and blogs text with the help of information extraction. We have extracted information from the collection of large text and mined opinions.(Xu et al,2010) The idea of reducing the information from a document to a tabular structure is not new but the identification of instances of a particular class of events is carried in this paper. There has been growing interest in developing information extraction as an enormous amount of information exists only in natural language processing. A mature information extraction technology would allow us to rapidly create extraction systems which will allow us to mine opinions from large text of information. (chaovalit et al,2005) The task of information extraction is used to identify a predefined set of texts of a specific domain and ignoring other irrelevant information where that domain consists of large amount of text.

The process of extraction such information involves the certain phrases or noun like words denoting a person's opinions and intensions and finding semantic relations between each others. However, in this case specific domain knowledge is required in order to accurately understand the opinions wile data gathering in a structured form.

#### 6 TF-IDF

In data recovery, tf-idf, short for term frequency-inverse record recurrence, is a numerical measurement that is proposed to reflect how imperative a word is to an archive in a gathering or corpus. It is frequently utilized as a weighting element in data recovery, content mining, and client demonstrating. The tf-idf esteem builds relatively to the quantity of times a word shows up in the record,

yet is regularly balanced by the recurrence of the word in the corpus, which alters for the way that a few words seem all the more every now and again when all is said in done. These days, tf-idf is a standout amongst the most well known term-weighting plans.

Tf-idf remains for term recurrence reverse record recurrence, and the tf-idf weight is a weight regularly utilized as a part of data recovery and content mining. This weight is a factual measure used to assess how critical a word is to a record in an accumulation or corpus. The significance builds relatively to the quantity of times a word shows up in the report however is balanced by the recurrence of the word in the corpus. Varieties of the tf-idf weighting plan are frequently utilized via web crawlers as a focal apparatus in scoring and positioning a record's pertinence given a client inquiry.

**TF:** Term Frequency, which measures how much of the time a term happens in a report. Since each report is distinctive long, it is conceivable that a term would seem significantly more circumstances in long records than shorter ones. Consequently, the term recurrence is regularly partitioned by the archive length (otherwise known as. the aggregate number of terms in the archive) as a method for standardization:

$TF(t) = (\text{Number of times term } t \text{ shows up in a report}) / (\text{Total number of terms in the record}).$

**IDF:** Inverse Document Frequency, which measures how vital a term is. While registering TF, all terms are considered similarly imperative. Nonetheless it is realized that specific terms, for example, "is", "of", and "that", may show up a considerable measure of times yet have little significance. Therefore we have to overload the successive terms while scale up the uncommon ones, by figuring the accompanying:

$IDF(t) = \log_e(\text{Total number of archives} / \text{Number of records with term } t \text{ in it}).$

**Consolidating these two we think of a tf score and an idf score. We figure these scores in the log-scale.**

The log term freq. of a term t in d is characterized as

$$1 + \log(1 + \text{tft}, d)$$

The log reverse archive recurrence which measures the usefulness of a term is characterized as:

$$\text{idf}t = \log_{10} N_d / \text{tft} = \log_{10} N_d / \text{tft}$$

where N is the aggregate number of archives in the gathering

Joining the two, the tf-idf score is given by

$$\text{wt}, d = (1 + \log(1 + \text{tft}, d)) \cdot \log_{10} N_d / \text{tft}$$

The tf.idf score increments with number of events inside an archive

The tf.idf score increments with uncommonness of terms in the accumulation

## 7 Machine Learning Algorithms for Opinion Mining

### 1. Decision Tree

Decision trees can be used to categorize text. They are designed so that the underlying data space is divided in a pattern in order to create class partitions which are more skewed in terms of their class distribution. For a given text example, a most suitable label is assigned to it and use it for the purposes of classification.

### 2. SVM (Support Vector Machine)

It is a grouping technique. In this calculation, we plot every information thing as a point in n-dimensional space (where n is number of components you have) with the estimation of each element being the estimation of a specific arrange). For instance, on the off chance that we just had two components like Height and Hair length of an individual, initially plot these two factors in two dimensional space where each point has two co-ordinates (these co-ordinates are known as Support Vectors)

### 3. Naive Bayes

It is an order method in view of Bayes' hypothesis with a suspicion of autonomy between indicators. In straightforward terms, a Naive Bayes classifier expect that the nearness of a specific element in a class is irrelevant to the nearness of whatever other element. For instance, a natural product might be thought to be an apple in the event that it is red, round, and around 3 creeps in distance across. Regardless of the possibility that these components rely on upon each other or upon the presence of alternate elements, a gullible Bayes classifier would consider these properties to freely add to the likelihood that this natural product is an apple. Naive Bayesian model is anything but difficult to fabricate and especially valuable for extensive informational collections. Alongside straightforwardness, Naive Bayes is known to beat even exceedingly refined arrangement strategies.

### 4. KNN (K-Nearest Neighbors)

It can be utilized for both characterization and relapse issues. Nonetheless, it is all the more broadly utilized as a part of characterization issues in the business. K closest neighbors is a straightforward calculation that stores all accessible cases and arranges new cases by a larger part vote of its k neighbors. The case being doled out to the class is most regular among its K closest neighbors measured by a separation work. These separation capacities can be Euclidean, Manhattan, Minkowski and Hamming separation. Initial three capacities are utilized for persistent capacity and fourth one (Hamming) for straight out factors. On the off chance that  $K = 1$ , at that point the case is essentially allocated to the class of its closest neighbor. Now and again, picking K ends up being a test while performing KNN demonstrating.

### 5. K-Means

It is a sort of unsupervised calculation which tackles the bunching issue. Its strategy takes after a straightforward and simple approach to arrange a given informational collection through a specific number of bunches (expect k groups). Information focuses inside a bunch are homogeneous and heterogeneous to associate gatherings.

In K-means, we have bunches and each group has its own centroid. Aggregate of square of contrast amongst centroid and the information focuses inside a group constitutes inside entirety of square an incentive for that bunch. Likewise, when the entirety of square esteems for every one of the groups are included, it winds up plainly add up to inside total of square an incentive for the bunch arrangement. We realize that as the quantity of bunch expands, this esteem continues diminishing yet in the event that you plot the outcome you may see that the aggregate of squared separation diminishes pointedly up to some estimation of k, and afterward considerably more gradually after that. Here, we can locate the ideal number of bunch.

### 6. Random Forest

Random Forest is a trademark term for a gathering of choice trees. In Random Forest, we've gathering of choice trees. To arrange another question in light of qualities, each tree gives an order and we say the tree "votes" for that class. The timberland picks the arrangement having the most votes (over every one of the trees in the woodland). Each tree is planted and developed as takes after: In the event that the quantity of cases in the preparation set is N, at that point test of N cases is brought aimlessly yet with substitution. This example will be the preparation set for developing the tree. On the off chance that there are M input factors, a number  $m \ll M$  is determined with the end goal that at every hub, m factors are chosen indiscriminately out of the M and the best part on these m is utilized to part the hub. The estimation of m is held consistent amid the timberland developing. Each tree is developed to the biggest degree conceivable. There is no pruning.

## 8 Extraction of opinions from large text

The major feature of Opinion mining is opinion extraction from large volume of text with the vast variety of existing work. In this paper an unsupervised approach is used to extract features by mining syntactic patterns of features. In this paper syntactic relations among opinion words in reviews and data in sentences by using information extraction. The results that have been executed by using R Studio and the various packages like tm, ggplot, stringi, wordcloud and nnet are used to depict distributional characteristics on opinion features. The mentioned packages also exhibit graphically and depict the density and strength of public opinions. The approach is used to mine the relationship between opinion words and topic related features. The approaches for feature extraction only use the patterns that have been mined from multiple reviews. After collecting large amount of various data from different sources huge number of unnecessary features are produced during feature collection phase which needs to be avoided. Feature filtering process removes unwanted surplus features by applying various algorithms like pruning algorithms in natural language processing.

## 9 Experimental Result and Finding

This research demonstrates a text-analytic approach in the analysis of data collected from various sources via different means reflecting people sentiments and reflect their opinions about the government policies and relationship with neighboring country. The methodology for data collection is described next.

### Data source

The data is extracted from twitter platform. The tweets collected related to text string India and Pakistan issues. Later the tweets were preprocessed using a set of rules to remove retweets, similar tweets, short length and symbol based tweets. Approximately a collection of 8000 tweets was analysed. Additional text data was also collected by interviewing a set of people over several blogs related to issues.. After cleaning approx 500-600 tweets were removed. training and test data is 70% and 30 % respectively chosen by classifier.

### Screen shot of CSV file

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
1	As Indian govt told Pakistani artist to leave india. According to me there is no fault of Pakistani artist as they are here in india just to make there career. this is wrong to create hurdles in anyone's carer																			
2	According to me whatever Indian govt is doing is absolutely correct because whenever Indian people are against Pakistani . pakistani artist's never support india. After being earning Indian currency they never support indian govt. they are earning from India and they are not even supporting ev																			
3	we should not blame actors and actresses and they are not a part of terrorist community they are like separate organization .they have different lifestyle as compare to Pakistan as they are living in India																			
4	I think as artist are not involved in all this terrorism so Indian govt. should not Indulge them in all this issue ,think because of all this there career should not get affect																			
5	There is a war between india and Pakistan and because of this somones career should not get affect.India and Pakistan should make some kind of agreement so that these all thing should not affect others life																			
6																				
7	Like we know we import dryfruits from Pakistan and US have banned Chinese good so this is wrong that Indian govt have told Pakistani artist to leave india if every county will do like this then relations between countries itself will affect. govt should find other ways to control this issue																			
8																				
9	Whatever terrorist community is doing they are getting successful in their motive they are capturing and indulging all these things with their issues which is wrong																			
10																				
11	Counter attack should be done to make Pakistan realize																			
12																				
13	The friendly conversation what till today was happening will not going to take this india Pakistan scenario to anywhere else .what I seriously feel is a strict action should be taken																			
14																				
15	After this uri attack which is really disastrous .and even after 2 days to uri attack some suspicious attack were also seen in Chandigarh.and not about Chandigarh and uri where terrorist want to attack they target at a small and crowded place																			
16																				
17	As recently said by chetan bhagat that Pakistan should be renamed as terroristan though I seriously support chetan bhagat and believe that strict action should be taken																			
18																				
19	Whatever actions have been taken by Indian army after uri attack should be taken quite often should that we could stand at a part																			
20																				
21	Pakistan and Indian issues are very serious and we should take an action but we should keep it separate to politics itself.that cultural knot should not be broken																			
22																				
23	According to me this action is not correct ,we should bridge the gap between Indian and Pakistan but after this it will make a huge gap																			
24																				
25	War is not a solution of anything and there is not religion of any artist. their religion is just to entertain people. they are talented because of their talent this is not like they are in india that's why they are talented if they could be at some other county even then they would be talented																			
26																				
27	There is not link between terrorism and Pakistani artist Indian govt is doing wrong																			
28																				
29	The Government should prepare a communication package on J&K issue justifying its stance as well as why the issue is closed from its perspective except for hand-over of PoK. Kashmiris call Indians as aliens. Did we have alien as first PM of India. How is it that Hindu Mattoo is Indian and Kashmiri																			
30																				
31	Congress too played the same Pakistan card for every debate and now it is the BJP repeating history. In fact we ignore history and pretend to have forgotten it but if we dig the history we will find the solution to Kashmir problem because answer lies within the problem by itself.																			
32																				
33	The constant vigilance of the Dastgiris Indians have at least reminded the leftist circles to warm towards ramabharina the night of Dandite of Kashmir. It is astonishing that a naren sitting outside Dahi does not see the difficulties of ordinary Indians who are less fortunate than Kashmiris an																			

**Primary source:** A focused group interview has been structured regarding the people sentiments and feedback on the concerned issues. Face to face interviews have been scheduled and taken as a primary source by asking reviews against India and Pakistan on different issues, public opinions and youth opinions. These interviews were video recorded to analyze and record the statements. The research paper endeavors to capture the major issues which have been prevailing in society from previous generations and will be affecting our future generation. After that converting audio and video data to text data is converted in csv format and then R Studio is used to process the data for data analytics.

**Secondary source:** The realtime data available from social media offers a wealth of information that could be used toward enhancing and developing our understanding the issues and sentiments further.

Blogs and twitter reviews have been collected to know about reviews of public over India Pakistan issues which have been running since many years. Blogs pages have been become a major source to express one's opinions and emotions. Bloggers update the daily events on the blog and express their emotions, opinions and sentiments through blog.

The other source that has been used as a secondary data is reviews; the user generated reviews are compiled from varied sources and websites on internet. The reviewer's data is widely used to judge the emotions, intentions, expectations and opinion on a particular domain. The twitter data is being extracted using ROAuth package and data is processed to extract texts and approximately 320 tweets were extracted from different handles.

## Results and interpretation:

### Word Cloud:

A rich variety of text mining methods are available in R which allow us to highlight the most frequently used keywords in texts. We can create a word cloud, also referred as *text cloud* or *tag cloud*, which reflects the visual representation of text data. The procedure of creating word clouds involves the use of the text mining package (*tm*) and the word cloud generator package (*wordcloud*) available in R for helping us to analyze texts and to quickly visualize the keywords as a word cloud.

A list of data pre-processing steps are conducting such as statement extraction: firstly, from each post, sentences are individuated, that are parts of text ending with a full stop, comma, question mark, exclamation mark or semicolon. Next step is tokenization where each statement is divided into tokens, which are parts of text bounded by a separator (space, tab or end of line). Stemming and Lemmatization is also done in order to reduce the number of different terms, each token is transformed reducing its inflectional forms to a common base form. Stop-words are also eliminated so that common categories are eliminated to properly distinguish among significant statements. Hence, in this step articles, prepositions, pronouns and conjunctions are first recognized and then removed the reviews that have been collected through interviews and blogs are processed and a word cloud has been generated through it. In this paper the words like Pakistan India, Uri are the words that contain highest frequency that is why in the word cloud they are highlighted more.

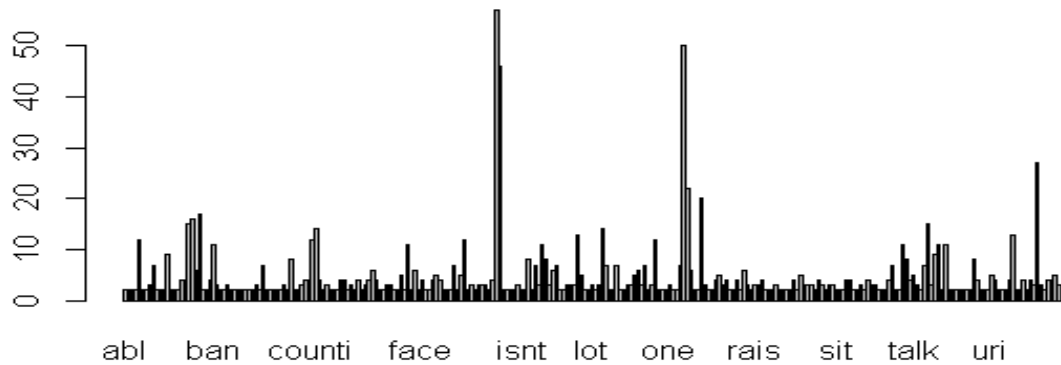
Figure 1: Word Cloud



### Bar plot:

A bar plot or bar graph is a graph that represent grouped data with rectangular bars with x-axis and y-axis. The bars can be plotted vertically or horizontally which shows the grouped data representation on graph. Basically it is a chart that is used to show comparisons among different categories with respect to topics on which data has been collected. On axis the chart shows specific categories being compared and the other axis it represents the discrete value. It also represents the bars clustered in groups of more than one category. In this grouped bar plot which is derived using R packages, scripts have been coded to derive the bar plot using collected data.

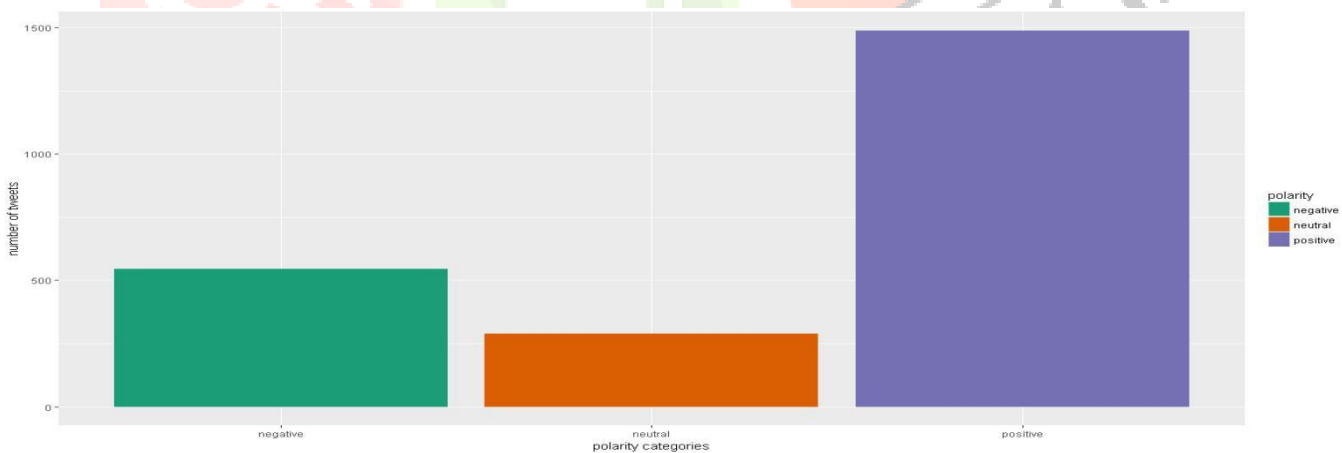
Figure 2:ggplot of term frequencies



**Polarity Graph:**

Twitter, as a web-based social networking is an extremely mainstream method for communicating feelings and cooperating with other individuals in the online world. At the point when taken in accumulation tweets can give an impression of open conclusion towards occasions. In this paper, we give a positive or negative conclusion on Twitter posts utilizing a notable machine learning strategy for content order. What's more, we utilize physically marked (positive/negative) tweets to assemble a prepared strategy to finish an undertaking. The errand is searching for a connection between twitter estimation and occasions that have happened. The prepared model depends on the Bayesian Logistic Regression (BLR) order technique. We utilized outer dictionaries to identify subjective or target tweets, included Unigram and Bigram highlights and utilized TF-IDF (Term Frequency-Inverse Document Frequency) to sift through the components.

Figure 3:Polarity Graph

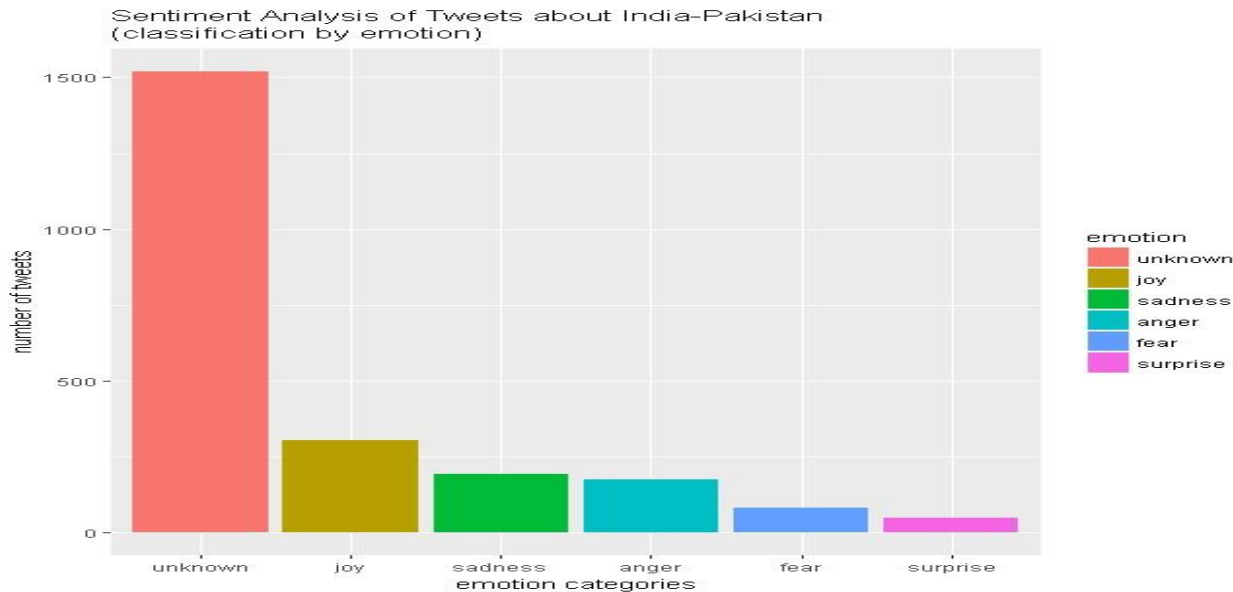




### Sentiments graph:

Sentiment graph shows with six emotions and three polarity dimension of user on tweets about demonetization in Fig. Result shows that people trust on the decision taken by the government. Public has taken Demonetization as a positive step by Indian Government and they support the step. Polarity of tweet can be positive, negative or neutral so after allotting the polarity result shown in the graph fig 2. Each tweet has been classified in one of eight emotions and the result shown in fig 3 and Fig 4. Representation of tweets according to allotted emotion to each tweet.

**Figure 4: Sentiments graph**



## 8 Conclusion and Summary

By considering this paper we are able to judge that opinion mining has a wide variety of applications in information systems, including classifying reviews, blogs, comments, it is concluded that opinion analysis has been done on the basis of a specific domain and different types of features has distinct distributions, and it is also found that different kind of scripts and algorithms has been written to analyzed a large amount of data and then finally opinions.

In future, more work is needed to improve the performance measures so that opinion mining can also be applied to new applications although the techniques that have been used in opinion mining are emerging very fast to at an extent that now we can also judge one's feelings however many limitations are still there.

## 9 References

1. Akcora, C.G., Bayir, M.A., Demirbas, M., Ferhatosmanoglu, H.: Identifying Breakpoints in Public Opinion. In: 1st Workshop on Social Media Analysis, pp. 62--66. Washington, DC (2010)
2. Baek, H.M., Ahn, J.H., Oh, S.W.: Impact of Tweets on Movie Sales: Focusing on the Time when Tweets are Written. J. ETRI. (2014)
3. Bai and Pad.R(2006), "Markov blankets and meta-heuristic search: Sentiment extraction from unstructured text," Lecture Notes in Computer Science, vol. 3932, pp. 167–187, 2006.
4. Boateng, F., & Quan Liu, Y. (2014). Web 2.0 applications' usage and trends in top US academic libraries. Library Hi Tech, 32(1), 120–138.
5. Boutet, A., Kim, H., Yoneki, E.: What's in Your Tweets? I Know Who You Supported in the UK 2010 General Election. In: The International AAAI Conference on Weblogs and Social Media (2012) Advanced Science and Technology Letters Vol.51 (CES-CUBE 2014) Copyright © 2014 SERSC 227
6. Burgert, L., Nann, A., & Sterling, L. (2014). Ventures in social media. Codex: The Journal of the Louisiana Chapter of the ACRL, 3(1), 20–44.

7. Choi, H., Varian, H.: Predicting the Present with Google Trends. Technical Report, Google (2009)
8. Chao, Zhou, Lina (2005), Movie Review Mining: a Comparison between Supervised and Unsupervised Classification Approaches, Proceedings of the 38th Hawaii International Conference on System Sciences – 2005.
9. Cuddy, C., Graham, J., & Morton-Owens, E. (2010). Implementing Twitter in a health sciences library. *Medical Reference Services Quarterly*, 29(4), 320–330.
10. Chu, S., & Du, H. (2013). Social networking tools for academic libraries. *Journal of Librarianship and Information Science*, 45(1), 64–75.
11. D’Amuri, F., Marcucci, J.: “Google it!” Forecasting the US Unemployment Rate with a Google Job Search Index. In: Conference on Urban and Regional Economics (2009)
12. Diakopoulos, N., Shamma, D.A.: Characterizing Debate Performance via Aggregated Twitter Sentiment. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pp. 1195–1198. ACM (2010)
13. Esuli, A. and Fab. S. (2005), “Determining the semantic orientation of terms through gloss classification”, Proceedings of 14th ACM International Conference on Information and Knowledge Management, pp. 617-624, Bremen, Germany, 2005.
14. Forman, B. Pang and L. Lee (2004), “A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts,” Proc. 42nd Ann. Meeting on Assoc. for Computational Linguistics.
15. Gaha, U., & Hall, S. (2015). Sustainable use of social media in libraries. *Codex: The Journal of the Louisiana Chapter of the ACRL*, 3(2), 47–67.
16. H.Zhen, Kui, Jung-Jae K., and Yang C. (2014). “Identify Features in opinion mining via intrinsic and extrinsic Domain Relevance” Vol 26, no.3, March 2014.
17. Jakob, N. and G.I. (2010), “Extracting Opinion Targets in a Single and Cross-Domain Setting with Conditional Random Fields,” Proc. Conf. Empirical Methods in Natural Language Processing, pp. 1035-1045, 2010.
18. Jin W and Ho.H.H (2009), “A Novel Lexicalized HMM-Based Learning Framework for Web Opinion Mining,” Proc. 26th Ann. Int’l Conf. Machine Learning, pp. 465-472, 2009.
19. Jin.W and Ho.H.H (2009), “A Novel Lexicalized HMM-Based Learning Framework for Web Opinion Mining,” Proc. 26th Ann. Int’l Conf. Machine Learning, pp. 465-472, 2009.
20. Lee, C.H., Hur, J., Oh, H.J., Kim, H.J., Ryu, P.M., Kim, H.K.: Technology Trends of Issue Detection and Predictive Analysis on Social Big Data. *J. Electronics and Telecommunications Trends*. 28, 62--71 (2013)
21. Loudon, L., & Hall, H. (2010). From triviality to business tool: The case of Twitter in library and information services delivery. *Business Information Review*, 27(4), 236–241
22. Liu, Y., Huang, X., An, A., Yu, X.: ARSA: a Sentiment-aware Model for Predicting Sales Performance Using Blogs. In: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 607–614. ACM (2007)
23. Palmer, S. (2014). Characterizing university library use of social media: A case study of Twitter and Facebook from Australia. *The Journal of Academic Librarianship*, 40(6), 611–619.
24. Qiu, G., Liu, B., Bu, J., and Chen, C. (2015), “Opinion Word Expansion and Target Extraction through Double Propagation,” *Computational Linguistics*, vol. 37, pp. 9–27, 2011. *International Journal of Computer Applications* (0975 – 8887) Volume 109 – No. 4, January 2015 32
25. Ramsey, E., & Vecchione, A. (2014). Channeling passions: Developing a successful social media strategy. *Journal of Library Innovation*, 5(2), 61–82.
26. Rui, H., Liu, Y., Whinston, A.: Whose and What Chatter Matters? The Effect of Tweets on Movie Sales. *Decision Support Systems*. 55, 863–870 (2013)
27. Sakaki, T., Okazaki, M., Matsuo, Y.: Earthquake Shakes Twitter Users: Real-time Event Detection by Social Sensors. In: 19th International Conference on World Wide Web, pp. 851–860. ACM (2010)
28. Tumasjan, A., Sprenger, T.O., Sandner, P.G., Welpe, I. M.: Election Forecasts with Twitter How 140 Characters Reflect the Political Landscape. *J. Social Science Computer Review*. 29, 402–418 (2011)
29. Sannella, J.M., Kim, M. and Ho, E. (2006), “Extracting Opinions Holders, and Topics Expressed in Online News Media Text,” Proc. ACL/COLING Workshop Sentiment and Subjectivity in Text, 2006.
30. Wu, Chu, Shen, L. (2009), “A New Method of Using Contextual Information to Infer the Semantic Orientations of Context Dependent Opinions”, 2009 International Conference on Artificial Intelligence and Computational Intelligence
31. Xu, Bing, Zhao, Tie, Zheng, De. Quan, Wang, Shan (2010), “Product features mining based on conditional random fields model”, Proceedings of the Ninth International Conference on Machine Learning and Cybernetics, Qingdao, 11-14 July 2010