

SENTIMENT ANALYSIS ON TWITTER USING STREAMING API

M.Trupthi¹,
Research Scholar,
Computer Science Department,
JNTUH, Telangana State, India.

Suresh Pabboju²,
Information Technology Department,
CBIT, Hyderabad, Telangana State,
India.

G.Narasimha³,
Computer Science Department, JNTUH,
Jagital, Telangana State, India.

Abstract—In general, opinion mining has been used to know about what people think and feel about their products and services in social media platforms. Millions of users share opinions on different aspects of life every day. Spurred by that growth, companies and media organizations are increasingly seeking way to mine information. It requires efficient techniques to collect a large amount of social media data and extract meaningful information from them.

This paper aims to provide an interactive automatic system which predicts the sentiment of the review/tweets of the people posted in social media using hadoop, which can process the huge amount of data. Till now, there are few different problems predominating in this research community, namely, sentiment classification, feature based classification and handling negations. A precise method is used for predicting sentiment polarity, which helps to improve marketing strategies. This paper deals with the challenges that appear in the process of Sentiment Analysis, real time tweets are considered as they are rich sources of data for opinion mining and sentiment analysis. This paper focus on Sentiment analysis, Feature based Sentiment classification and Opinion Summarization.

The main objective of this paper is to perform real time sentimental analysis on the tweets that are extracted from the twitter and provide time based analytics to the user.

Keywords—Sentiment Analysis, Streaming API, Twitter

I. INTRODUCTION

Microblogging sites are rich in sources for a varied kind of information. This is a common place where people exchange their opinions on various issues it could be on ongoing trends. Based on their experiences they share a comment or complaint on any product and express their thoughts in terms of positive or negative sentiment. Many upcoming organizations require feedback analysis on their products to improve further. Most of the time, Organizers analyze the user responses and answer them back on social media. So here is a challenge to analyze or detect and accomplish the global sentiment.

Huge unstructured data is available in many forms like tweets, reviews or news articles etc. which can be classified as positive, neutral or negative polarity according to the sentiment that is expressed by them. The main focus of the paper is to build a system which can build a classifier from a large data set at the start up and then can perform classification of tweets based on the classifier built.

Performing sentimental analysis on the tweets based on naïve Bayes' classifier trained from a large data set and provide time variant analytics based on the results obtained. Training should be performed on a large data set which is the major criteria to get efficient classifier.

There are many challenges in Sentiment Analysis. Firstly, an opinion word which is considered to be positive in

one state may be considered negative in a different situation. Second one, people may not always express opinions in the similar manner. Eg: "the picture was a great one" differs completely from "the picture was not so great". Opinion of people may be contradictory in their statements. It is more difficult for a machine to analyze. Most of the time people find it difficult to understand what others mean within a short sentence of text because it lacks context.

Sentiment analysis is done on three levels

- Document Level: Analysis is done on the whole document and then express whether the document positive or negative sentiment [3].
- Sentence level: It is related to find sentiment polarity from short sentences. Sentence level is merely close to subjectivity classification.
- Entity /Aspect Level sentiment analysis performs augmented analysis. The aim is to find sentiment on entities or aspects. eg: consider a statement "My Samsung S5 phone's picture quality is good but its phone storage capacity is low". Samsung camera and the quality of display has positive sentiment but phone's storage memory sentiment is negative [3].

This paper focuses on the short sentences and entity level sentiment analysis and classify the streamed tweets into positive, neutral and negative tweets using standard classifier.

A. Natural Language Processing – NLTK

NLTK is a suite of text processing libraries for tokenization, stemming, classification, tagging, parsing, and semantic reasoning. It also has lexical resources such as WordNet. It does the following:

- Stop words removal
- Unstructured to structured
Tweets are mostly unstructured i.e. 'rip' is written 'rest in peace', 'oooooooooooo' to actually 'good'. Conversion to structured is done by dynamic data
- Emoticons: The symbolic emoticons are converted in to words i.e. to sad

B. Naïve Bayes' Classification

In machine learning, naïve Bayes classifier uses Bayes' theorem with strong (naïve) independence assumptions between the features which were word frequencies.

Naïve Bayes classifiers are highly accessible, requires number of parameters which are linear in the number of variables (features/predictors) in the learning problem. Training of Maximum-likelihood can be used for evaluation of a closed form expression which considers linear time, rather than expensive approximate iteration that is used for different types of classifiers.

Naive Bayes is a classifier technique used for building classifiers: Models assigns class labels to instances, represented feature values as vectors, in which class labels are extracted from some finite set. For example, the fruit which is round may be considered as an orange if its colour is orange, round, and it is about 4" in radius. Naive Bayes classifier independently considers each of these features to find the probability whether the fruit is an orange, regardless of any possible relationships between the features like roundness, colour and diameter.

C. Twitter Application Programming Interface

The interface TwitterAPI is used to collect streaming Tweets from Twitter which also stores tweet scores along with its timestamp.

Publicly posted Tweets published by users are extracted. In order to create a POST request to the twitter API and fetch the search results as a stream it uses CreateStreamingConnection() method. In one connection 5,000 Twitter user ids are allowed to submit for an application. Only publicly published Tweets can be captured using the API. The Streaming API searches for hashtags, keywords and geographic bounding boxes simultaneously. The filter API helps for searching and delivers the continuous stream of Tweets which matches the filter tag. POST method is preferred while creating the request, because long URLs are truncated and GET method is used to retrieve the results.

II. LITERATURE SURVEY

Turney et.al.[1] used bag-of-words method in which the relationships between words was not considered at all for sentiment analysis and a sentence is simply considered as a collection of words. To determine the sentiment for the whole sentence, sentiment of every individual word was determined separately and those values are aggregated using some aggregation functions.

Pak and Paroubek [2] proposed a model to classify the tweets as positive and negative. By using Twitter API they created a twitter corpus by collecting tweets and automatically annotating those tweets using emoticons. Using that corpus, the multinomial Naive Bayes sentiment classifier method was developed which uses features like POS-tags and N-gram. The training set used in the experiment was less efficient because they considered only tweets which have emoticons.

Po-Wei Liang et.al. [3] used Twitter API to collect data from twitter. Tweets which contain opinions were filtered out. Unigram Naive Bayes model was developed for polarity identification. They also worked for elimination of unwanted features by using the Mutual Information and Chi square feature extraction method. Finally, the approach for predicting the tweets as positive or negative did not give better accuracy by this method.

Thet [4],proposes a linguistic approach system for aspect based opinion mining, which is a clause/Sentence level sentiment analysis for opinionated texts. For every message

post sentence it generates a syntactic dependency tree, and splits the sentence into clauses. It then determines the contextual based sentiment score for each clause using grammar dependency of words and uses SentiWordNet which has prior sentiment scores for the words and also from domain specific lexicons.

Hussein[5], this paper explains the previous works, the goal is to identify the most significant. challenges in sentiment and explore how to improve the accuracy results that are relevant to the used techniques.

All the above mentioned work uses the corpus data in this paper the real streaming data based on the filteris used and it does not require any memory to store the tweets.

III. PROPOSED METHODOLOGY

The proposed system extracts the data from SNS services which is done using Streaming API of twitter. The extracted tweets are loaded into hadoop and it is been preprocessed using map reduce.This task is followed by classification which uses NLP and machine learning techniques. The classification used here is uni-word naïve bayes' classification.

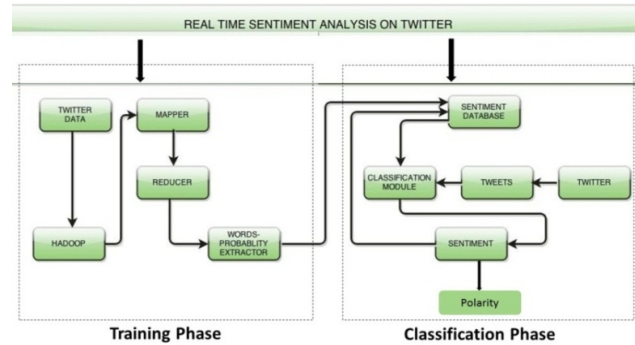


Figure1. Proposed Block Diagram

Consider the number of all positive tweets, positive words and negative words from our training phase. Then calculate the probability of a tweet being positive.

$$P(C) = \frac{\text{No of Positive Tweets}}{\text{Total Number of Tweets}} \quad (1)$$

For each word in each tweet that is being streamed is checked for the probability of it being used given that it is positive.

$$P(D|C) = \frac{\text{Positive Score of the Word}}{\text{Total Number of Positive Words}} \quad (2)$$

Then we checked the word itself being used irrespective of whether or not it is positive.

$$P(D) = \frac{\text{Positive Score of the word} + \text{Negative score of the word}}{\text{Total Number of Words}} \quad (3)$$

In-order to check the probability of word being positive given that is used in a tweet which is given as follows:

$$P(C|D) = \frac{P(C) * P(D|C)}{P(D)} \quad (4)$$

The probability of a word is then passed to the Sentiment function which then classifies, if the probability of the word is greater than 0.6 then it is positive, as neutral if the probability is between 0.4 and 0.6 and negative if it is lesser than 0.

IV. IMPLEMENTATION

The implementation of the work is done in three folds. First the pre-processing, and training of the data set. a) Training Phase, second classification and scoring for the tweets based on the filters set. b) Classification Phase and then followed by representation of the data on the web application. c) Application phase and User Interface.

A. Training & Streaming phase

In the first phase Training the classifier is the atmost important task. AS an input to this module 20,00,000 tweets were collected from several sources which are already classified and the job of this module is to build a classifier by training on the large data set. Nltk is used to remove the words with POS tags which are not usefulto build the classifier. Hadoop is used to extract information from it and MapReduce is used to easily extract several words with their positive and negative probabilities. The output of reducer is several numbers of words with their positive and negative scores. These scores are stored in database using MongoDB, which inturn is used by the classification module. The classification module is used to classify the tweets from twitter.

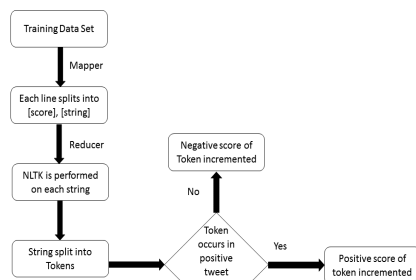


Figure2. Training Phase

The dataset that is already classified is given a sentiment score of 2 or 5 for each tweet, indicating that is negative for a score 2 and positive for a score 5. The dataset considered for training is offered by Stanford University and the classification is first done by human.

16269	5	@GreenFalcon805 thank u for your help..
16270	2	@greenfee I have Sky, Setanta the lot! no womens golf since golf channel pulled pit of the UK. Enjoying UK!
16271	2	@greenfelttip Social networking and the fundamental limit of 24 hours in a day makes everyone a bad pers
16272	5	@greenflash008 I guess you can tell I've got some pretty eclectic tastes. Jackie Chan and Horace Walpole.
16273	2	@greengalz I'm sorry I hope it gets better. http://myloc.me/4y5l
16274	2	@greengalz WE need a new schedule
16275	2	@greengeco29 I am sorry for your loss
16276	5	@greengiri74 @MrJamesMillar I'm so working on that in a month!!
16277	5	@greengvingco and breathe..... Hope time slows down for you !
16278	5	@greengooddy Happy early three-oh. So far, it's my favourite decade by a long shot!!!
16279	2	@greenhouseffekt Thanks jon for following up
16280	2	@greeniebach Fed won 6-4 6-4. I didn't see it, but am guessing Rafa was tired from yesterday? Sorry love.

Figure3. Stanford Annotated Score

B. Classification Phase

The tweets extracted by Streaming API are then classified into positive, negative or neutral tweet. If the words turn out to be positive, then the sentence is classified as positive. Mapper code when runs on this dataset will split the file into two parts namely the score, an integer and the tweet, which is string.

The Reducer will check for each word in the string and increment its positive score if the overall tweet is scored as 5. Similarly, the negative score for the word is incremented if the word happens to occur in another tweet which is scored as 2.

On the other hand, if the words turn out to be negative, then the sentence is classified as negative. If the words are mixture of both positive and negative, then we check the sum of positive and negative scores of words, appropriately the sentence is classified. The final output of the Reducer is stored in the format {[word], [negative_score], [positive_score]} as shown in fig.4. and The final scores uploaded onto MongoDB.

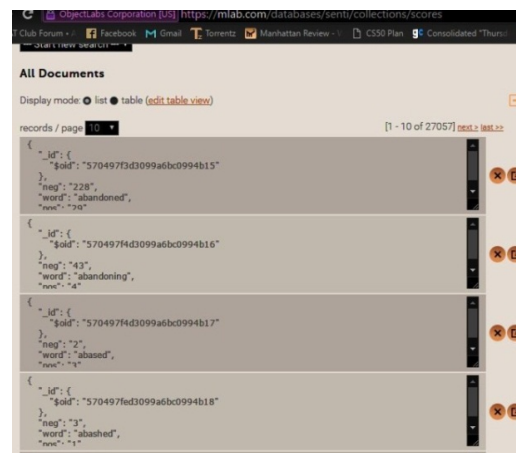
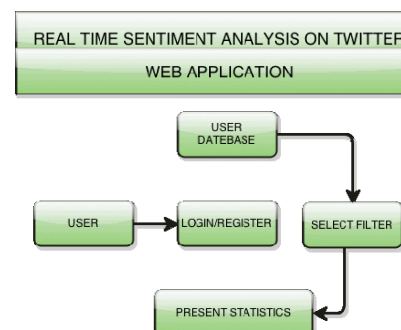


Figure4. Scores Records Stored in MongoDB

C. Application Phase & User Interface

The web Application allows all users to register by providing basic information themselves to the application. The details are stored in MongoDB under the users collection. Whenever the user tries to login to their account, their details



will be checked against the details stored in MongoDB for a match.

Figure 5. Web Application Phase

Initially when the user logs in, there will be no filters so the dashboard will be empty. The user has been provided options to add/delete filters. When the user starts adding filters then all the data can be analyzed in the UI Module. Hence, results for each filter can be viewed by the user.

In the UI module provides the interface to analyse the classified data. Here the user can set their own filter for which they can visualize the data through a Donut Chart for the depiction of time-based Analytics imported from the Morris API. The users can choose to display the data over hourly, daily or weekly durations.

V. RESULTS

A) Results for the filter "Obama"

All the tweets that are tweeted from the time of execution that contain the word Obama are scored for sentiment analysis. Based on the scores obtained, the tweet is classified as positive, negative or neutral. These results are displayed as the summary opinion for each filter through a donut chart.

```

C:\Users\kruthi>python tweet.py
['obama', 'OBAMA', 'Obama']
obama The Latest: Obama says US and Europe must help refugees https://t.co/GqWqWqs1z positive
obama @camefromempire @jewelsinfo @realDonaldTrump Obama sealed records and now we know trump was right negati
ve
obama Herrenhausen #Quint #Roni #Mollande #Cameron #Obama #Merkel #Hannover #cosedilavore https://t.co/IYr
wo2JlCU https://t.co/vn88qiu2026 neutral
obama 'Obama did nothing for this country' https://t.co/FD36Kv64HF negative
obama #Obama now decides we're at war with #ISIS something every person in this country has said for a yr. His
incompetence is off the charts negative
obama Obama pleads for unity in Europe: President Barack Obama made a lengthy, personal plea for European unity
Monday asu2026 https://t.co/N6Q5d5Kmc neutral
obama Barack Obama says world needs a united Europe https://t.co/atU6uWry9 positive
obama glucioQuincioC yo ya perd'u00ed el capacidad de asombro ante tanta estupidez deber'u00edamos hacer una pe
tici'u00f3n en las calle pidiendo visitas d Obama neutral
obama BBC News - Syria conflict: Obama to deploy 250 more special forces troops https://t.co/0duK8FwK1 neut
ral
obama The Latest: Obama says US and Europe must help refugees https://t.co/WpIkPQwQV positive
obama Sterling hits 10-week high amid Obama bounce https://t.co/R8BEzx18s neutral
obama 'Obama did nothing for this country' https://t.co/FD36Kv64HF negative
obama ReNews: Odds move sharply towards Britain staying in EU after Obama warning -Reuters- https://t.co/ZY
N7Lw0iQc negative
obama Obama ve Kerry'den #u00d3zg'u00fcrBasu0131n etiketi ile a'u00e7u0131klama:Basu0131n, liderlerin hesa
p vermesini sa'u0131flar https://t.co/rdoKshW95 https://u2026 neutral
obama 'Obama did nothing for this country' https://t.co/FD36Kv64HF negative

```

Figure6. Streaming Twitter for the filter Obama



Figure7. Analytics for the Filter Obama

B. Results for the Filter ISIS

The analytics for ISIS, it is evident that Twitter verse feels mostly negative about ISIS. Most of the tweets contained links to articles involving ISIS tweeted by news agencies, so they were scored to be neutral. Very few tweets were classified as positive which were mostly tweeted by people who support the Islamic Front.

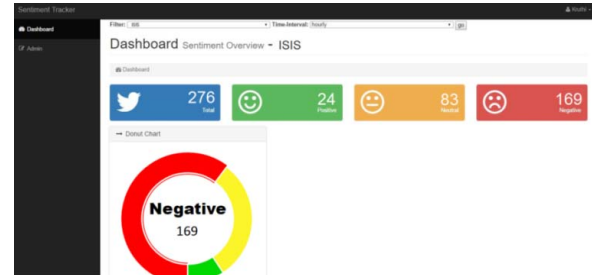


Figure8. Analytics for the Filter ISIS

C. Results for the filter Education

The results as expected were mostly positive. A few were neutral because of articles reported and very few tweets were classified as negative.

The results have been more accurate for filters which are contained in tweets tweeted in English language. As this classification is language dependent and only English Language is considered, when we take tweets specific to a country, say India, where we know people also tweet in languages like Hindi, the results are not entirely accurate. Most of them will be classified as Neutral or Negative. Majority of negative words carry more weight which gives rise to this classification.

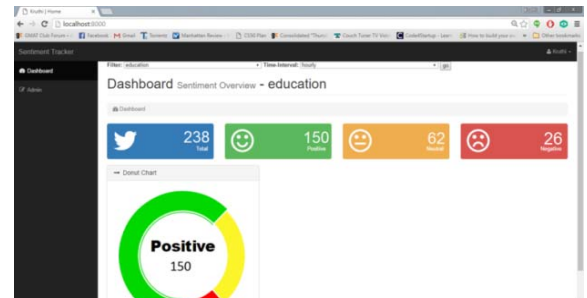


Figure 9. Analytics for the Filter Education

VI. CONCLUSION AND FUTURE SCOPE

This work is of tremendous use to the people and industries which are based on sentiment analysis. For example, Sales Marketing, Product Marketing etc. The key features of this system are the training module which is done with the help Hadoop and MapReduce, Classification based on

Naïve Bayes, Time Variant Analytics and the Continuous-learning System. The fact that the analysis is done real time is the major highlight of this paper. Several existing systems store old tweets and perform sentiment analysis on them which gives results on old data and uses up a lot of space. But in this system, the tweets are not stored which is cost effective as no storage space is needed. Also all the analysis is done on tweets real-time. So the user is assured that, getting new and relevant results.

However, the proposed system has some limitations. First limitation is being Uni-gram Naïve i.e. training of the data was done based on word probabilities and used the same for classification. Future enhancement to this work might be to use n-gram classification rather than limiting to uni-gram which will require pattern filtering on Hadoop. When classifying the sentence, words are taken individually rather than the sentence in total. The semantic meaning is neglected that might be present between the words. Second Limitation this is only used for English Language. It might be possible to build a system which can perform sentiment analysis in all languages. The third limitation that the system may not provide actual intended meaning of the user. There might be some sort of sarcasm present in the sentences which the system ignores.

References

- [1] P. D. Turney, "Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews," in Proceedings of the 40th annual meeting on association for computational linguistics, pp. 417–424, Association for Computational Linguistics, 2002.
- [2] A. Pak and P. Paroubek. „Twitter as a Corpus for Sentiment Analysis and Opinion Mining". In Proceedings of the Seventh Conference on International Language Resources and Evaluation, 2010, pp.1320-1326.
- [3] J Po-Wei Liang, Bi-Ru Dai, "Opinion Mining on Social MediaData", IEEE 14th International Conference on Mobile Data Management,Milan, Italy, June 3- 6, 2013, pp 91-96, ISBN: 978-1-494673-6068-5, <http://doi.ieeecomputersociety.org/10.1109/MDM.2013>.
- [4] T. T. Thet, I.-C. Na, C. S. Khoo, and S. Shakthikumar, "Sentiment analysis of movie reviews on discussion boards using a linguistic approach," in Proceedings of the 1st international CIKM workshop on Topic-sentiment analysis for mass opinion. ACM, 2009, pp. 81-84.
- [5] Hussein, D.-M.E.D.M. A survey on sentiment analysis challenges. Journal of King Saud University – Engineering Sciences (2016), <http://dx.doi.org/10.1016/j.jksues.2016.04.002>
- [6] A Kowcika and Aditi Guptha "Sentiment Analysis for Social Media", International Journal of Advanced Research in Computer Science and Software Engineering, 216-221,Volume 3,Issue 7, July 2013.
- [7] G.Vinodini and RM.Chandrashekarana, "Sentiment Analysis and Opinion Mining: A Survey", International Journal of Advanced Research in Computer Science and Software Engineering, 283-294, Volume 2, Issue 6, June 2012.
- [8] Apoorv Agarwal Boyi Xie Ilia Vovsha Owen Rambow Rebecca Passonneau, "Sentiment Analysis of Twitter Data",Columbia University,New York.
- [9] Pablo Gamallo and Marcos Garcia "A Naive-Bayes Strategy for Sentiment Analysis on English Tweets" Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014), pages 171–175, Dublin, Ireland, Aug 23-24 2014.
- [10] Harry Zhang "The Optimality of Naive Bayes", FLAIRS2004 conference. (available online: PDF (<http://www.cs.unb.ca/profs/hzhang/publications/FLAIRS04ZhangH.pdf>)).
- [11] Ms.K.Mouthami, Ms.K.Nirmala Devi, Dr.V.Murali Bhaskaran, "Sentiment Analysis and Classification Based on Textual Review".
- [12] Sentiment Analysis Data Set Corpus<http://thinknook.com/twitter-sentiment-analysis-training-corpus-dataset-2012-09-22/>
- [13] Sang-Hyun Cho and Hang-Bong Kang, "Text Sentiment Classification for SNS-based Marketing Using Domain Sentiment Dictionary", IEEE International Conference on Consumer Electronics (ICCE), p.717-718, 2012.
- [14] Aurangzeb Khan and BaharumBaharudin, "Sentiment Classification Using Sentence-level Semantic Orientation of Opinion Terms from Blogs", 2011.
- [15] Popescu, A. M., Etzioni, O, *Extracting Product Features and Opinions from Reviews*, In Proc. Conf. Human Language Technology and Empirical Methods in Natural Language Processing, Vancouver, British Columbia, 2005, 339–346.