**International Journal for Research in Applied Science & Engineering Technology (IJRASET)**

# Enhanced Twitter Sentiment Analysis by using Open NLP with Positive Negative and Neutral Tweet

Neha Upadhyay[1], Assistant Prof. Angad Singh[2]

[1,2]Research Scholar Information and Technology,  Research Guide Information and Technology,Bhopal

*Abstract: Due to the vast opinion of rich web resources such as discussion forum, review sites , blogs and news corpora on the market in digital form, a lot of analysis is focusing on the area of sentiment analysis. People are trying to develop a system that can identify and classify opinion or sentiment as represented in an electronic text. A correct method to predict sentiments could modify us, to extract opinions from the internet and predict online customer's preferences,  might prove valuable for economic and market research. Till now, there are few different issues predominating in this analysis, i.e, sentiment classification, feature based classification and handling negations. In this thesis, we try to estimate the twitter posts about electronic products like mobiles, laptops etc.we collect the tweets of mobile from twitter and preprocess the tweets, in preprocessing remove irrelevent words such as name, symbols etc. after preprocessing we catogarized the tweets in two parts opinion or opinionless, we compare each tweets to database of positive, negative and average words those tweets carry positive, negative and average words, is in opinion cateogity other's discarted. now we check the polarity of opinion tweets, count the polarity of tweets and divide in positive, negative and average category. Estimating the probablity of positive tweets we use Open NLP for the better probability result.*

*Keywords: Tweets,Sentiment Analysis, Machine Learning Techniques, Open NLP Document Categorizer Tool, Maximum Entrop*
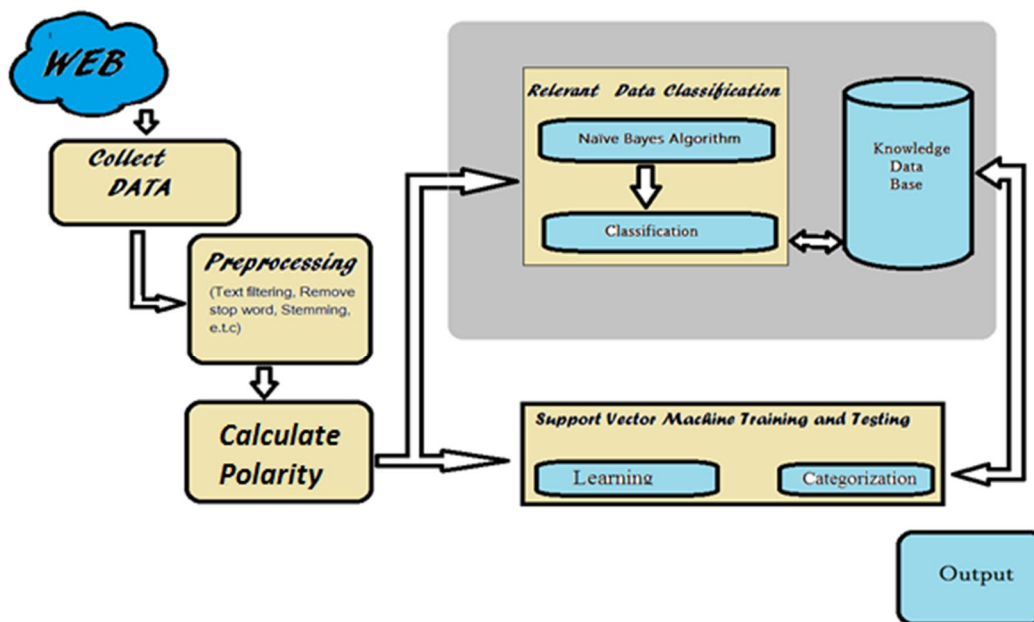
## I.    INTRODUCTION

The period of Internet has changed the way individuals express their perspectives. It is currently done through blog entries, online disscusion forums, item review sites etc. people rely on this client created data as it were. When people needs to purchase an item, they will lookup its surveys online before taking a choice. The measure of client created content is too large for an ordinary client to analyse. So to solve this, different sentiment analysis have used[1]. Symbolic methods or Knowledge base approach and Machine learning strategies are the two primary procedures are use as a part of opinion analysis. Knowledge base approach requires an expansive database of predefined emotions and a efficient learning representation for identifying sentiments[2]. Machine learning approach makes use of a training set to build up a sentiment classifier that classify opinions. Since a predefined database of whole feelings is not required for machine learning approach, it is  more simpler than Knowledge base methodology. In this paper, We used different machine learning procedures for classifyng tweets. Sentiment analysis is normally used at various levels fluctuating from coarse level to fine level.[3] Coarse level assumption examination manages deciding the estimation of a whole archive and Fine level manages quality level slant investigation. Sentence level estimation investigation comes in the middle of these two. There are numerous looks into on the territory of opinion investigation of client audits. Past examines demonstrate that the exhibitions of slant classifiers are subject to points. Due to that we can't say that one classifier is the best for all subjects since one classifier doesnt reliably outflanks the other. Conclusion Analysis in twitter is very troublesome because of its short length. Nearness of emoticons, slang words and incorrect spellings in tweets compelled to have a preprocessing venture before highlight extraction. There are diverse component extraction strategies for gathering pertinent elements from content which can be connected to tweets moreover. In any case, the element extraction is to be done in two stages to extricate applicable elements. In the primary stage, twitter particular components are removed. At that point these elements are expelled from the tweets to make ordinary content. After that, again include extraction is done to get more components. This is the ticket utilized as a part of this paper to create an effective component vector for investigating twitter assessment.[4][5] Since no standard dataset is accessible for twitter posts of electronic gadgets, we made a dataset by gathering tweets for a specific timeframe.

*A.   Related Work*

# International Journal for Research in Applied Science & Engineering Technology (IJRASET)

1) *Machine Learning Techniques*: Machine Learning strategies used a training set and a test set for a classification. Training set contains information feature vectors and their corresponding class labels. By using this training set, a classification model is created which tries to classify the information feature vectors into corresponding class labels.[7] At that point a test set is used to accept the model by predicting the class labels of unseen feature vectors. Various machine learning merthods like Naive Bayes (NB), Maximum Entropy (ME), and Support Vector Machines (SVM) are used to classify reviews. We are using the Open NLP Document Categorizer Tool for the better Result, introduction of words, expressions, sentences and that of documents. Semantic introduction is the polarity which may be either positive or negative or neutral.

*B.* *Proposed Architecture*

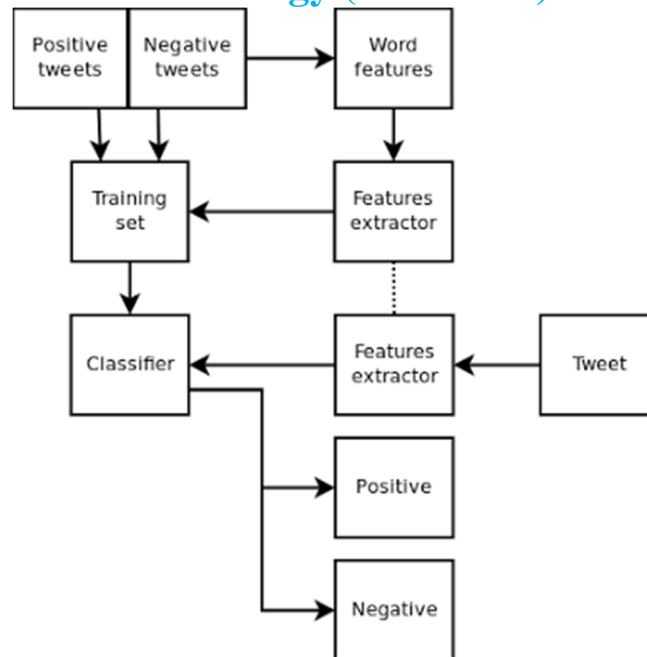# International Journal for Research in Applied Science & Engineering Technology (IJRASET)



Fig1 and 2: Proposed Architecture

1) *Text Pre-Processing:* Text Pre-preparing for extract significant data.
a) *Stop Words:* an, is, the, with and so on. The full list of stop words can be originate at Stop Word List. These words don't show any sentiments and can be removed.
b) *Repeating Letters:* In the event that you take a look at the tweets, in some cases peaple repeat letters to stretch the feeling. E.g. hunggrryyy, huuuuuuungry for 'hungry'. We can search for 2 or more repeative letters in words and replace them by 2 of the same.
c) *Punctuation:* We can remove punctuations like comma, single/double quote, question marks toward the begin and end of every word. E.g. excellent!!!!!! replaced with excellent

## II.    METHODOLOGY USED

A. *Open NLP Document Categorizer Tool*
OpenNLP is an arrangement of java-based Natural Language Processing tools using the Maximum Entropy , a statistical learning approach**.**

B. *OpenNLP and Machine Learning*
OpenNLPcan be used for:
1) Sentence Detection
2) Tokenization
3) Named-Entity Detection
4) Sentence Parsing
5) Coreference
6) Document Classification

C. *Sentence Detection*
The OpenNLP Sentence Detector can identify that a punctuation character denote the end of a sentence or not. In this sense a sentence is characterized as the longest white space trimmed character grouping between two punctuation marks. The first and last sentence make an exception case to this rule. The first non-whitespace character is thought to be the start of a sentence, and the last non whitespace character is thought to be a sentence end.

# International Journal for Research in Applied Science & Engineering Technology (IJRASET)

*D. Tokenization*

The OpenNLP Tokenizers segments an input character sequence into token. Tokens are generally words, punctuation, numbers and so on.

OpenNLP propose multiple tokenizer implementations:

1) *Whitespace Tokenizer:* A whitespace tokenizer, non-whitespace sequence are classified as tokens.
2) *Simple Tokenizer:* A character class tokenizer, sequences of the similar character class are token
3) *Learnable Tokenizer:* A maximum entropy tokenizer, detects token boundaries depends on probability model

Most part-of-speech taggers, parsers etc, work with content tokenized in this way. It is essential to ensure that your tokenizer produce tokens of the sort expected by your later text processing segments.

With OpenNLP (as with numerous frameworks), tokenization is a two-phase process: in the first place, sentence boundaries are classified, then tokens inside every sentence are classified.

First step: Create Open NLP Categorizer object.
Second step: Give input datasets of Tweets.
Third Step: Train open NLP Model
Fourth Step: Give Tweets to be Analyse
Fifth Step: Categorizer will start finding best category. Category will obtain in the form of
Sixth Step: Using 100 Iteration
Seventh Step: Show best category of Tweets, it will be positive or negative or Average .

*E. Maximum Entropy Classifier*

In Maximum Entropy Classifier, no assumptions are taken regarding to the relationship between features. This classifier always try to maximize the entropy of the system by assessing the conditional circulation of the class label. The conditional distribution is characterized as

$$P_\lambda(y|X) = 1/Z(X)exp\left\{\sum_i \lambda_i f_i(X,y)\right\}$$

'X' is the feature vector and 'y' is the class label. $Z(X)$ is the normalization factor and $\_i$ is the weight coefficient. $f_i(X,y)$ is the feature function which is defined as

$$f_i(X,y) = \begin{cases} 1, & X=x_i \text{ and } y = y_i \\ 0, & \text{otherwise} \end{cases}$$

In our feature vector, the relationships between part of speech tag, emotional keyword and negation are utilized effectively for classification.

## V. RESULT AND ANALYSIS

This part examines and evaluates the observational results of tests to validate the exhibited framework proposed. The reasonable motive of the work is accomplished in this section

Firstly, choose the Product for which we require to analyse the sentiment and after that calculate the Positive Probability of Product survey through any of them classifiers.

In this part we figure out the importance of the proposed approach. In proposed approach we are using two Open NLP for the Classification.

Almost 2000 tweets were classified into three different classes: Positive Sentiment, Negative Sentiment and Average Sentiment. Figure 4 shows classification of Product tweets into three classes (positive, negative and Average) the help of segment charts. We have taken Electronic Product tweets for our research work.

# International Journal for Research in Applied Science & Engineering Technology (IJRASET)
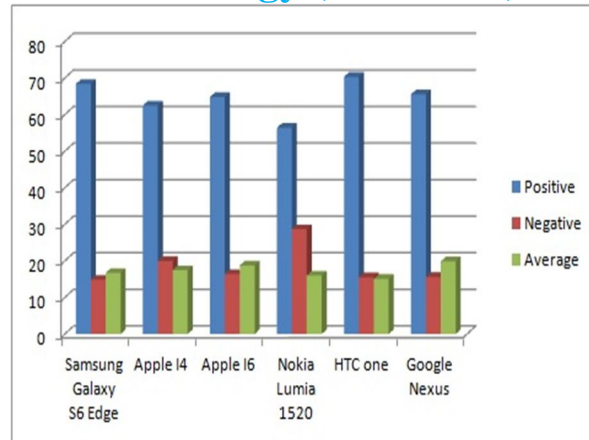


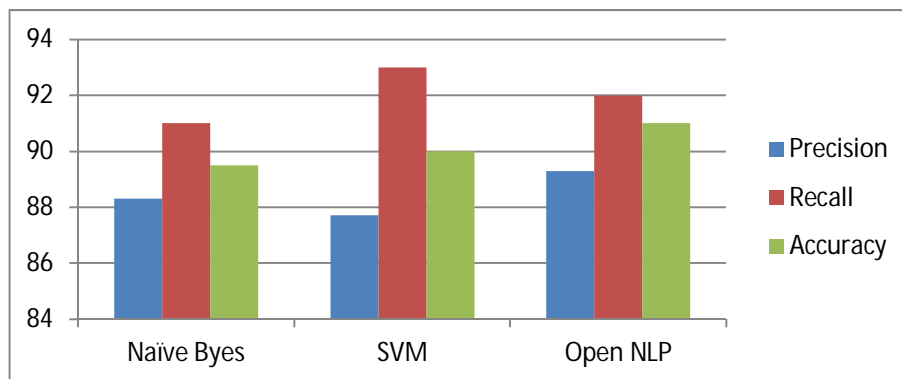Fig3: Classification of individual product tweets

We can figure our outcome through Naïve Bayes and S.V.M, and open NLP and compare the outcome with the help of positive sentiments, negative sentiments and neutral sentiments. We can compute our outcome through Naïve Bayes and S.V.M and Open NLP Document Categorizer Tool classifiers and compare the result with the help of positive sentiments, negative sentiments and neutral sentiments.

This line diagram explain the comparision of Naïve Bayes And Open NLP Document Categorizer for positive tweets of three Product
 (Samsung Glaxy S6, Apple I4, Apple I6).

| parameter | Naïve Byes | SVM | OpenNLP |
|-----------|------------|------|---------|
| Precision | 88.3 | 87.7 | 89.3 |
| Recall | 91 | 93 | 92 |
| Accuracy | 89.5 | 90 | 91 |

Fig.4,5  Comparison between Naive Bayes and Open NLP on Positive  review



## VI.    CONCLUSION

There are different  procedures to identify the sentiments from content. In this paper,our analysis represents  that Machine Learning strategies are less difficult and more efficient. This technique can be applied for twitter sentiment analysis. Theresome issues while dealing with identifying  setiment keyword from tweets having  multiple keywords. It is also challanging to solve missspelt and slang words.To manage this issue,to solve the difficulties in data, an efficient feature vector is made by doing feature extraction in two stages after appropriate preprocessing. Classification accuracy of the feature vector is tested by using different classifiers like Nave Bayes, Open NLP .This feature vector performs useful for electronic products.The message communicated in Twitter can be identified with the human behaviour, nature, personality and attitude. Classification of tweets into positive sentiments, negative sentiments or neutral sentiments indicates the views of people on Product. That's help people to choose best item and they easily

# International Journal for Research in Applied Science & Engineering Technology (IJRASET)

decide that which item is famous in the overall business market.

## REFERENCES

[1] H. Isah, P. Trundle and D. Neagu, "Social networking media identifying for product safety using data mining and sentiment analysis," *Computational Intelligence (UKCI), 2014 14th UK Workshop on*, Bradford, 2014, pp. 1-7.

[2] M. Anjaria and R. M. R. Guddeti"Influence factor based opinion mining of Twitter using supervised learning process," *(COMSNETS) Communication Systems and Networks, 2014 Sixth International Conference on*, Bangalore, 2014, pp. 1

[3] Tiara, M. K. Sabariah and V. Effendy, "Sentiment analysis on Twitter using lexicon-based and support vector machine methodology for assessing the performance of a television program," *( ICoICT ) Information and Communication Technology, 3rd International Conference on 2015*, Nusa Dua, 2015, pp. 386-390.

[4] R.Rajshree and M S. Neethu," machine learning techniques used in sentiment analysis in twitter," *(ICCCNT)Computing, Communications and Networking Technologies, Fourth International Conference 2013 on*, Tiruchengode, 2013, pp. 1-5.

[5] G. Gautam and D. Yadav, ."Sentiment analysis in twitter using semantic analysis and machine learningapproaches," *Contemporary Computing (IC3), 2014 Seventh International Conference on*, Noida, 2014, pp. 437-442.

[6] J. Akaichi, Z. Dhouioui and M. J. Lopez-Huertas Perez, "For sentiment classification text mining face book updates on," *(ICSTCC) System Theory, Control and Computing, 2013 17th International Conference*, Sinaia, 2013, pp. 640-645.

[7] M. S. Schlichtkrull, "Learning affective projections for emoticons on Twitter,*6th IEEE International Conference*," *(CogInfoCom) Cognitive Infocommunicationson*, Gyor, 2015, pp. 539-543.

[8] D. Hakkani-Tür and A. Celikyilmaz, J. Feng, "Probabilistic model-based sentiment classification of twitter texts," *(SLTTechnology of Spoken Language Workshop, IEEE 2010*, Berkeley, CA, 2010, pp. 79-84.

[9] Q. H. Vuong and A. Takasu," Transfer Learning for Emotional PolarityClassification," *(IAT) Intelligent Agent Technologies and (WI) Web Intelligence , 2014 IEEE/WIC/ACM International Joint Conferences on*, Warsaw, 2014, pp. 94-101.

[10] N. Azam, Jahiruddin, M. Abulaish and N. A. H. Haldar ,"Twitter Data Mining of Events Analysis and Classification," *2015 Second International Conference on Soft Computing and Machine Intelligence (ISCMI)*, Hong Kong, 2015, pp. 79-83.

[11] M. Kanakaraj and R. M. R. Guddeti," *Semantic Computing (ICSC), 2015 IEEE International Conference on*, Anaheim, CA, 2015, pp. 169-170.

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 ⊙ (24*7 Support on Whatsapp)