# iJRASET

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

www.ijraset.com

Call: ○ 08813907089  |  E-mail ID: ijraset@gmail.com

# Sentiment Analysis Using Hybrid Approach

Ganesh K. Shinde[1], Vaibhav N. Lokhande[2], Rasika T. Kalyane[3], Vikas B. Gore[4], Umesh M. Raut[5]

[1, 2, 3, 4, 5]Department of CSIT, Dr B.A.M. University Aurangabad Maharashtra India

Abstract: Most important part of information gathering is to focus on how people think. There are so many opinion resources such as online review sites and personal blogs are available. In this paper we focused on the Twitter. Twitter allow user to express his opinion on variety of entities. We performed sentiment analysis on tweets using Text Mining methods such as Lexicon and Machine Learning Approach. We performed Sentiment Analysis in two steps, first by searching the polarity words from the pool of words that are already predefined in lexicon dictionary and in Second step training the machine learning algorithm using polarities given in the first step.

Keywords: Sentiment analysis, Social Media, Twitter, Lexicon Dictionary, Machine Learning Classifiers, SVM.

## I. INTRODUCTION

Social Media such as Twitter, Facebook, Blogs are become important tools where user willing to share their important sentiments on different topics. With this kind of platforms various opportunities and challenges arise to actively use various techniques to extract and understand the sentiments of others. Sentiment Analysis of twitter data has many usages like review of customer towards movie, product, services and application. Opinion mining of tweets includes the classification of tweets as Positive, Negative or Neutral.

Lexicon Based Sentiment Analysis based on the presence of certain word in document. Lexicon contains features including the part of speech tagging of word, their sentiment values, subjectivity of word etc. The Sentiment Analysis of tweets are annotate using this features provided by these lexicons. Values are tagged against each word separately. Using that we can obtain polarity of whole tweet by averaging the sentiment values of words. The Machine Learning based Sentiment Analysis technique is a model by training the classifier with labeled examples. First we require to gather a dataset with positive, negative and neutral classes, extract the features/words from that dataset & then train the algorithm based on the examples. This is most easy method for sentiment analysis. We are using both these approaches Lexicon Based Approach and Machine Learning Approach. We are showing the result of sentiment analysis by combining these two approaches. Usually Lexicon based approach perform entity level sentiment analysis and it gives high precision but low recall. To improve the performance measurements such as Recall, F-Score, Accuracy. Machine learning algorithm is train using the polarity given by lexicon based approach. Our hypothesis is that the accuracy given by such approach is get increase with increase in size of training data.

## II. LITERATURE REVIEW

The aim of system is to classify the twitter messages as positive, negative or neutral. For this we use two approaches: Lexicon based approach and Machine Learning approach.Lexicon based approach includes performing the sentiment analysis at document and sentence level by searching polarity of word from predefined word list.[1,2,3,4] determine polarity of sentence using predefined dictionary. Examples of such a Lexicon dictionaries are MPQA [5] and SentiWordNet 3.0 [6] Machine Learning approach includes the three algorithms 1) Maximum Entropy, 2) Support Vector Machine (SVM), 3) Naïve Bayes (NB). This includes the training the classification algorithm [7]. In this paper we are using combination of both approaches, it also called Hybrid approach. Normally combinations of both are used for subjectivity classification and then apply it to the learning algorithm [8]. Similar approach used in [9], which classify sentence in only two classes positive & negative, no neutral class it creates problem. [10] Uses same approach with different features. In [13] lexicon and machine learning approaches are combine. But they use different sentiment analysis methods. These are different than our approach; we first preprocess the data to remove unwanted data from it. Then we perform polarity detection using Lexicon dictionary and then apply this result to the Learning algorithm.

## III. PROPOSED SYSTEM

Following figure 1 shows the workflow of the proposed system. Data acquisition is carried through the Twitter API. Twitter API allow user to interact to with its data i.e. tweets. User can download these tweets by creating twitter API. User request to API for the data and it returns data according to the query enter by user. The twitter data contain noisy data such as RT for Retweets, '#' hashtags for filtering tweets according to the topic, @usernames, external web links, and emoticons. The task of preprocessing removes all noisy data, so that data will be clean and it is easy to perform operations on clean data.

We perform 1) Remove Duplicate tweets, 2) Remove Retweets, 3) Remove URL's, 4) Remove Unnecessary Space, 5) Remove twitter hashtags, 6) Remove Punctuation Marks,7) Remove Numbers and, 8) Remove twitter username starts with @ symbol
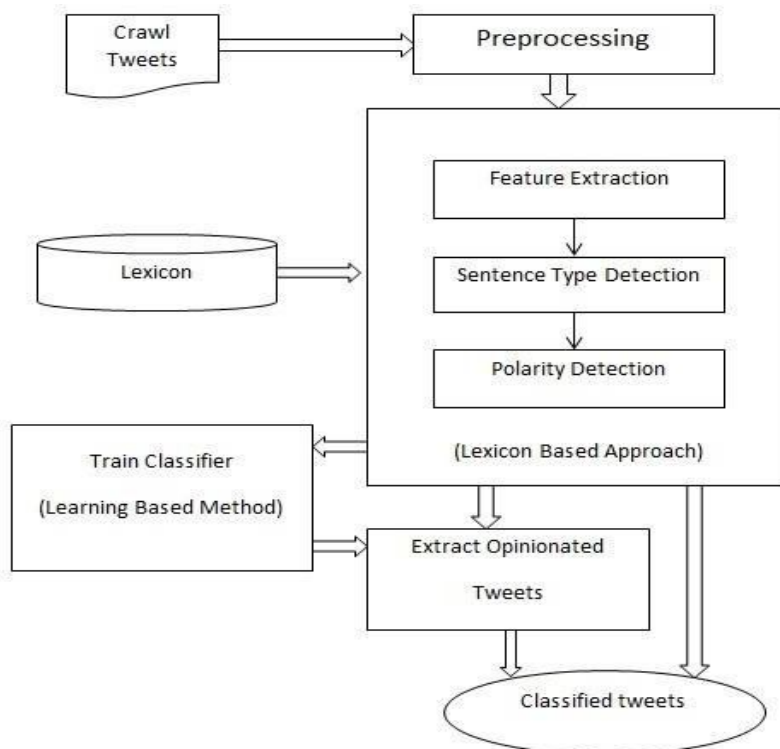


Fig -1: Workflow for sentiment Analysis

As tweeter contain much more unnecessary data, so we need to find data that contain opinion, which we use for sentiment analysis. So feature selection is way to find out this. Normally we find out the tweets that contain the Adjective, because presence of adjective in tweets indicate that the tweet contain the opinion about something in the world.

Next step is feature selection for find out the subjective tweets. Subjective tweets are the tweets that contain the user emotion, view about something in the world. So it is necessary to find the subjective tweets, for that we need to classify tweets as Subjective and Objective tweets we just find the polarity of tweets by searching the occurrences of the word in the lexicon dictionary, and simple replace the word position with the polarity value shows by the lexicon dictionary. The Polarity of whole tweet is calculated by the aggregation of the word polarity present in that tweet. Before that polarity is tagged against each word.

Negation Handling is one of the major issue in sentiment analysis. Because many sentence contain the negation word that shifts the polarity of the sentence. Many classifiers remove the negation words by considering it as stop words. We had overcome this problem, when we find any negation term in sentence then we simply replace that negation term with punctuation symbol '!'. For that we had just made some changes in the lexicon simply add symbol '!' before each word in lexicon and just shifts the polarity of that words. The polarity given by the lexicon dictionary for the each sentence can be considered as training data. These training data is given to Machine Learning Classifier to train the classifier. By training using this training data we calculate polarity of other data which can be passed as a testing data to the classifier. This implements the performance of the Sentiment Analysis. Training and testing data is used for experiments.

## IV.EXPERIMENTS

The initially we collect dataset using twitter API. The query we fired while collecting data from API is 'car' i.e. the API gives all data (tweets) that contain the word car. For our experiments we collect 28000 tweets.

Next step is to remove the noise form the collected data. Noise such as Duplicate tweets, Retweets, punctuations, numbers, HTML links etc. Data we extracted contain unnecessary that is extra data. Hence to extract only those tweets that contain the some opinion, we perform feature selection. This includes extracting only those tweets that contain adjective. This can be done by using Tree Tagger Part of Speech Tagging (POS) technique [10]. And then we classify those tweets as Subjective tend Objective tweets and consider only the Subjective as main feature for Sentiment Analysis.

This can be performed using the MPQA Lexicon which contain the words with its subjectivity information i.e. whether word id Strong or Weak Positive [11]. After performing this feature selection we have 25000 tweets for further processing.

We then perform the polarity detection using the MPQA Lexicon [11], where we search for the occurrence of the each word of tweet in a lexicon dictionary, when find then replace that word with the polarity value given by the lexicon. When we find that word is not occur in the lexicon then we replace it with polarity value zero that indicate the Neutral polarity. Finally we aggregate all polarity values of words in tweets that the aggregate value indicates the polarity of the tweet. This tweets are consider as the training data, which are used to train classifier

We train the classifier, so that it will assign the sentiment polarity to the newly sentiment tweets i.e. testing data. We use Support Vector Machine as our leaning algorithm. Training Data is data which can be labelled as positive, negative and neutral by the lexicon based method. Our basic features are Unigram, Bigram, and Trigram. We calculate the accuracy of polarity classification for the newly opinionated tweets using these classification features i.e. Unigram, Bigram and Trigram. Testing data is newly opinionated tweet that are to be classified based on training given to the classifier using the data classified by the lexicon based method

## V. EVALUATION

In Evaluations we performed the opinion analysis using the Support Vector Machine learning method. The reason for using this algorithm is that it gives better result than learning algorithms like Naïve Bayes, Maximum Entropy [12]. For the process of sentiment analysis, we divide training data into different parts, this is done to check the accuracy of sentiment classifier when the training data size increases. The aim is to just check the variations in the accuracy of sentiment classifier for same test data. We use measure Accuracy to evaluate the sentiment classification performance. This measure can be check against classification features such as Unigram, Bigram and Trigram with different training size. Table 1 shows the accuracy for all three classification features, with variation in training data size

Table I: Accuracy Results

| Training Data Size | Testing Data Size | Unigram | Bigram | Trigram |
|---|---|---|---|---|
| 5000 | 1000 | 60.53 | 59.38 | 57.13 |
| 10000 | 1000 | 61.62 | 60.53 | 57.26 |
| 15000 | 1000 | 62.28 | 61.20 | 58.95 |
| 20000 | 1000 | 62.42 | 61.27 | 59.02 |
| 25000 | 1000 | 63.23 | 62.23 | 59.98 |

## VI. CONCLUSIONS

Sentiment Analysis and opining mining for twitter data, when perform using Lexicon based approach it shows high precision but low recall so there is problem of performance. To increase that performance we combine both the approaches i.e. Lexicon Based Approach and Machine Learning Approach, this give better performance. We perform the evaluation for the different sized training datasets. And these evaluation results ensure that the proposed algorithm is highly effective and better for sentiment analysis of twitter messages. In this way we used hybrid approach for sentiment analysis             .

## REFERENCES

[1]  Taboada, M., Brooke, J., Tofiloski, M., Voll, K., & Stede, M. 2011. "*Lexicon- based methods for sentiment analysis*." Comput. Linguist. 37, (2): 267—307.

[2]  Hu, M., & Liu, B. 2004. "*Mining and summarizing customer reviews*. " In: Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '04). ACM, New York, NY, USA. pp. 168--177.

[3]  Kim, S., & Hovy, E. 2004. *Determining the sentiment of opinions*. In: Proceedings of the 20th International Conference on Computational Linguistics (COLING '04). Association for Computational Linguistics, Stroudsburg, PA, USA

[4]  Ding, X., Liu, B., & Yu, P.S. 2008. "*A holistic lexicon-based approach to opinion mining*." In: Proceedings of the International Conference on Web Search and Web Data Mining (WSDM '08). ACM, New York, NY, USA. pp. 231-240.

[5]  Pang, B., Lee, L., and Vaithyanathan, S. (2002)." *Thumbs up?: Sentiment classification using machine learning techniques*". In Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - Volume 10, EMNLP '02, pages 79–86, Stroudsburg, PA, USA. Association for Computational Linguistics.

[6] Multi Perspective Question Answering (MPQA). OnlineLexicon "http://www.cs.pitt.edu/mpqa/subj_lexicon.html.

[7] Stefano Baccianella, Andrea Esuli, Fabrizio Sebastiani. *"SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment Analysis s and Opinion Mining"*. In Proceedings of international conference on Language Resources and Evaluation (LREC), 2010.

[8] Wiebe, J. and Rilo_, E. 2005. "*Creating Subjective and Objective Sentence Classifiers from Unannotated Texts*. " CICLing 2005

[9] Tan, S., Wang, Y. and Cheng, X. 2008." *combing Learn- based and Lexicon-based Techniques for Sentiment Detection without Using Labeled Examples."* SIGIR 2008

[10] www.cis.uni-muenschen.de/~schmid-tools/TreeTagger/

[11] Multi Perspective Question Answering (MPQA) Online Lexicon <http://www.cs.pitt.edu/mpqa/subj_lexicon.html

[12] Go, A., Bhayani, & R., Huang, L. 2009. "*Twitter sentiment classification using distant supervision."* Technical report, Stanford.

[13] Lei Zhang , Riddhiman Ghosh, Mohamed Dekhil,, Meichun Hsu, Bing Liu 2011. "*Combining Lexicon-based and Learning-based Methods for Twitter Sentiment Analysis*". HPL Laboratories

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 ⟨24*7 Support on Whatsapp⟩