

# Performance based Machine Learning Algorithm for Topic Oriented Text Categorization

Paruchuri Ramya, Geetha Guttikonda, Vemuri Sindhura, Vinod Kumar Gadde

**Abstract**—With the growth of societal news on the web, public opinions are given major importance in decision-making. Researchers of text-based mining have made number of evaluations and were diversified using different data mining methods so as to make the conclusions positive, negative and neutral. So, opinions of people are considered to mine the social information as people give superfluous interest to the reports. In this paper the newspaper data set is considered to find the opinion mining to evaluate the sentiment. Sentiment Analysis is used to compute the opinions of people before they judge on a particular issue. Machine Learning is one of the important approaches for analysis of sentiments. Different methods like Naïve Bayes, SVM, Maximum entropy and SLDA are used for classifying the sentiments. Predictions based on precision, f-measure, recall are done to determine which method best suits the classification.

**Keywords**— Opinion, Sentiment Analysis, Machine Learning, Naive Bayes, SVM, Maximum Entropy, SLDA, Prediction, Precision, Recall, F-Measure

## I. INTRODUCTION

Online journalism in India is a competing field as there are many opportunities created online for newspapers. The chance of strengthening the newspaper industry survival is increasing day by day and also with the help of relationship with advertisers. The movement away from the printing process can also help decrease costs.

The usage of social media is increased with the rise in population. As a responsibility it is necessary for people to know what is happening in and around the world. People pay more attention to opinions rather than to their preferences. The imperative understanding of opinion mining is constantly about the data gathering conduct of what other individuals think.

Sentiment analysis is tied in with social opinions and knowing how individuals respond to a specific situation. The past work of sentiment analysis was done on effectiveness of market and performance of stock price. Or maybe this work is for the most part focused on numerical highlights which don't compare specifically with gathering of information in regards to the articles that are irrelevant in classification of sentiment.

To track the temper of people in general about a particular issue or an item, the normal preparing system that is utilized is the Sentiment investigation. The sentiments about the item made in the discussion posts, remarks, visits, tweets and so on were gathered and analysed for examination. The supposition mining is likewise called as assumption investigation.

Assessment investigation can be useful from multiple points of view. In this paper, the attention is on the sentiment analysis of openly accessible news reports to give investigation to the general public.

## II. EXISTING SYSTEM

An Kim and Hovy[1] assess the estimation of a opinion holder (entity) utilizing WordNet to create arrangements of positive and negative words by growing seed records. They expect that equivalent words (antonyms) of a word have the same (inverse) extremity. The level of a word's equivalent words having a place with arrangements of either extremity was utilized as a proportion of its extremity quality, while those underneath an edge was esteemed nonpartisan or equivocal [9][10][11].

M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede [2] presented a lexicon-based sentiment classification. In this classification they utilized lexicons of positive or negative spellbound words. A semantic orientation calculator (SO-CAL) was built in light of these lexicons by fusing intensifiers and refutation words. This approach has been appeared to have 59.6% to 76.4% precision on 1900 reports of motion picture audit data set.

R. McDonald, K. Hannan, T. Neylon, M. Wells, and J. Reynar [3] developed a model for sentiment analysis at different levels of granularity at the same time. They use graphical models in which a document level assessment is associated with a couple of entry level inclinations, and each area level assumption is associated with a couple of sentences level slants. They apply the Viterbi count to construe the finish of each substance unit, constrained to ensure that the segment and report parts of the imprints a genuine way a comparative where they address a comparative entry/record. They report 62.6% precision at requesting sentences when the presentation of the record isn't given, and 82.8% exactness at ordering documents. The inconvenience of this computation is that it prompts genuine over fitting.

Two mind-set following instruments, Opinion Finder and Google Profile of Mood States, were utilized to examine the content substance of day by day Twitter (Bollen et al., 2011) [4]. The previous estimates the positive and negative temperament. The last estimates the mind-set as far as six measurements (Calm, Alert, Sure, Vital, Kind, and Happy). They utilized the Self Organizing Fuzzy Neural Network model to anticipate DJIA esteems. The outcomes demonstrated 86.7% course precision (up or down) and 1.79% Mean Absolute Percentage Error. Despite the fact that they accomplished the high exactness, there were just 15 exchange dates (from December 1 to 19, 2008) in their test set. With such a brief period, it probably won't be adequate to finish up the viability of their technique.

Revised Manuscript Received on 10, September 2019.

\* Correspondence Author

Paruchuri Ramya \*, Department of Information Technology, VR Siddhartha Engineering College, Kanuru, Vijayawada, India, #4 Project Associate, CTS, Chennai, Tamilnadu, India, paruchuriramyaa2010@gmail.com

A keyword-based algorithm was proposed to distinguish the assessment of tweets as positive, neutral and negative for stock expectation (Bing Liu and Lei Zhang. 2012) [5]. Their model accomplished around 75% exactness. Nonetheless, their trial was short, from 8th to 26th in September 2012, containing just 14 exchange dates.

Constant Dirichlet Process Mixture (cDPM) show was utilized to take in the everyday subject arrangement of Twitter messages to anticipate the share trading system [6]. A sentiment time series was manufactured in light of these themes. Be that as it may, the day and age of their entire data set is fairly short, just three months.

In addition to the opinion words that are in general, the topic models that considers a set of opinion words that are aspect specific are proposed [7][8]. A hybrid model which uses Maximum entropy & LDA can discover the set of aspects and opinion words that are generally specific to an aspect for a restraint review data set (Zhao et al., 2010) [7]. The models CFACTS-R, FACTS-R, FACTS, CFACTS are proposed for analysis of sentiment for a product review data set (Lakkaraju et al., 2011) [8]. Major drawbacks of the above methods are that for a corresponding topic only a single opinion word distribution is considered. These methods make it difficult to know the sentiment (negative/ positive) articulated by the opinion words wrt particular topic.

Proposed System The current work has inspected the model's capacity to foresee stock market developments in light of assessment from a given day. The past work has additionally broken down the connection between the present stock return and positive or negative groupings of tomorrow's news articles, under the presumption that news articles posted the following day morning are ordinarily about occasions that happened today. The need of proposed work is to break down any sort of online news utilizing different kinds of scientific strategies and to pick which techniques gathering the best. From this the client can without much of a stretch comprehend the sentiment or opinions for a specific kind of news.

### III. PROPOSED SYSTEM

Opinion mining (now and then known as Sentiment Assessment or Emotion Artificial Intelligence) alludes to the utilization of natural language processing (NLP), analysis of text, computational semantics and biometrics to efficiently recognize, extricate, evaluate, ponder effective states and emotional data. This paper proposes an approach which can be extended to several sentiment analysis problems. The idea of this approach is that emotion examination calculations can perform better when the information, which are prepared, deal with a less wide category of topicPaper Submission Criteria The above figure represents the proposed system architecture. The entire system is implemented in three phases where phase I deals with the steps of pre-processing which produces a set of lexicons, phase II deals with topic extraction, and phase III handles the steps of sentiment analysis based on topic. The steps of proposed system are described below:

1. Split the news data into sentences and make a Bag of Sentences (BoS).
2. Process of segmentation/ tokenization is applied on Bag of Sentences which produces a set of "tokens" (Bag of Words).

3. After the segmentation and applying the steps of pre-processing like lemmatization, stop word removal, stemming, slang word removal a final set of lexicons are obtained on which topic extraction.
4. An LDA is model is applied on the pre-processed data which is used to acquire a probability-based word matrix and document- topic matrix which is utilized for topic summarization and obtaining the sentiment analysis for each individual topic.
5. A comprehensive dictionary (feature vector) of each important feature in the sentence is formed and classification is done using machine learning and lexicon-based approaches.
6. With use of this a polarity of each contextual feature in the news data and polarity of each news group is categorized as positive, negative.

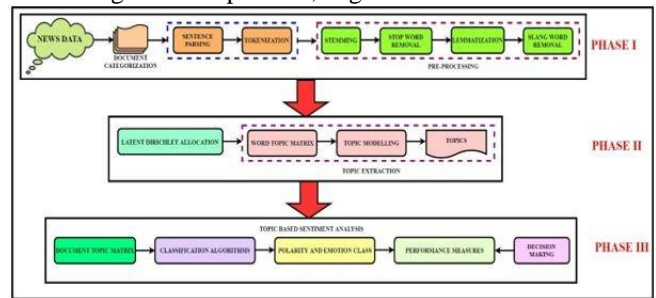


Fig 1: Architecture of Proposed System

#### A. Data Collection and Categorization

The data set utilized in this paper is a Usenet Newsgroup data set which is a gathering of roughly 20,000 newsgroup reports, divided (about) equally crosswise over 20 distinctive newsgroups. The Usenet bulletin board in this data set that incorporate newsgroups for themes like governmental issues, religion, automobiles, games, and cryptography, and offer a rich arrangement of content composed by numerous clients.

#### B. Pre-processing

The objective behind pre-processing is to represent to each record in the document as feature vector/tokens, to isolate the content into singular words which decides the quality of the next stage, the classification stage and it is done utilizing the following steps.

- 1) *Sentence Parsing*: Sentence parsing utilizes the whole textual input information and partition it into set of sentences based on a period at end of each sentence.
- 2) *Tokenization*: Tokenization partitions the whole set of obtained sentences into smaller units called "Tokens" isolated by spaces.
- 3) *Stemming*: Stemming is a procedure to expel the prefixes, additions from a word and change it to its base form.
- 4) *Stop Word Removal*: The idea is essentially expelling the words that happen generally over every one of the documents. When all is said in done, articles and pronouns are generally named stop words.
- 5) *Lemmatization*: Lemmatization is the algorithmic procedure of deciding the lemma for a given word that is by one means or another like stemming, as it maps a few words into one normal root.
- 6) *Slang Word Removal*: The majority of information includes larger part of

slang words which are to be changed into standard words to make free content.

### C. Topic Extraction

Term Frequency – Inverse Document Frequency (TF-IDF) is the most ordinarily utilized method in Natural Language Processing (NLP) which is utilized to sift through the high recurrence words that don't contain profitable data while stressing critical low recurrence words. Topic Modelling is an unsupervised approach utilized for finding and watching the group of words (called "topics") in expansive bunches of texts.

Blei et al. developed the Latent Dirichlet Allocation (LDA) model, a generative probabilistic model, to gather discrete content information. The LDA show is a three-level (archives subject's words) various levelled Bayesian model which alludes to the examination to draw inferences from unlabelled datasets and generally used to remove the idle point data from huge set of documents. LDA model's "D" is generated according to the following process:

- 1) Choose a multinomial distribution "φ" for a topic "t" from a Dirichlet distribution with parameter β.
- 2) Choose a multinomial distribution "θ" for a document "d" from a Dirichlet distribution with parameter α.
- 3) For a word "w<sub>n</sub>" in document "d",
  - "z<sub>n</sub>" is selected from "θ<sub>d</sub>".
  - "w<sub>n</sub>" is selected from "φ<sub>z<sub>n</sub></sub>".

$$p(D|\alpha, \beta) = \int \prod_{d=1}^M \int \prod_{n=1}^{N_d} p(\theta_d | \alpha) \left( \sum_{z=1}^{N_z} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \phi) P(\phi | \beta) \right) d\theta_d d\phi$$

In LDA, the latent topics are created with the use of vocabulary generated from the accumulated documents. Documents obviously contain a blend of topics, where probability distribution over a group of terms is treated as a subject. Each document is then observed as a probability distribution over the group of topics. Assume the data as originating from a generative process that is characterized by the probability distribution over what is watched and what is covered up

#### Topic based sentiment analysis

The document-topic matrix, another yield of LDA (Latent Dirichlet Allocation), speaks to the co-occurrence likelihood of a document and a latent subject. For the extremity and classification of emotion, dole out the documents that are centred around a couple of topics to the particular topics as per an edge esteem, which is reliant on the dissemination of document topic probabilities.

- 1) *Sentiments Lexicon Dataset*: In the sentiment's dataset, a wide variety of lexicon data sets are available out of which the below mentioned three are broadly used.

- AFINN from Finn Arup Nielsen,
- bing from Bing Liu and associates, and
- nrc from Saif Mohammad and Peter Turney.

All the above-mentioned lexicons are basically unigrams (single words) and are a huge collection of variety of words in English. Each of these words can be assigned with a score to obtain the negative/positive notion which expresses the feelings like bitterness, happiness, outrage and so forth. For topic i, the negative ratio is obtained as a ratio of negative sentiments in the total sentiment.

$$\text{Neg. ratio}_i = \frac{\text{Num\_neg\_sentiment}_i}{\text{Num\_total\_sentiment}_i} \quad (1)$$

The ratio of emotion is obtained as the proportion of sentiments to the total proportion for a particular topic "i" and emotion "j".

$$\text{Emotion ratio}_{ij} = \frac{\text{Num\_emotion}_{ij}}{\text{Num\_total\_emotion}_i} \quad (2)$$

- 2) *N-Gram Analysis*: The Usenet dataset is a much larger corpus of more modern text, so a bigram approach is used in order to determine this feature. Bigram or digram is an arrangement of two adjoining components from a series of tokens, which are ordinarily words, syllables or letters. A bigram is a n gram for n=2. The frequency distribution of each bigram in a string is normally utilized for straightforward statistical analysis of content in numerous applications, incorporating into speech recognition, cryptography, and computational semantics. Bigrams help provide the conditional probability of a token given the preceding token, when the relation of the conditional probability is applied.

$$P(W_n | W_{n-1}) = \frac{P(W_{n-1}, W_n)}{P(W_{n-1})} \quad (3)$$

The probability for a particular token W<sub>n</sub> wrt to token preceding it W<sub>n-1</sub> is equivalent to its bigram probability / tokens co-occurrence P(W<sub>n</sub> | W<sub>n-1</sub>) which is divided by the probability of token preceding it P(W<sub>n-1</sub>).

#### A. Machine Learning Algorithms

Symbolic approaches and machine learning techniques are two main methodologies in sentiment analysis. The approaches in which manually generated lexicons and rules are referred to as "Symbolic approaches". In this paper, techniques of machine learning are being used.

1) *Naïve Bayes Algorithm*: In classification of data a machine learning technique, Naive Bayes technique [11] is one of the best methods. For classification of data a "Multinomial Naive Bayes Technique" is used where "d" denotes a news data, "c\*" denotes a class which is assigned to data "d", where

$$P_{NB}(c|d) := \frac{P(c) \prod_{i=1}^m P(f_i|c)^{n_i(d)}}{P(d)} \quad (4)$$

In the above used equation, the terms "f" represents a feature, "n<sub>i</sub>(d)" represents feature count, "d" represents the data, features count is denoted by "m". "P(c)" and "P(f<sub>i</sub>|c)" represents the maximum likelihood estimates.

- 2) *Maximum Entropy*: Maximum Entropy for a text classification is implemented by word counts as features [9][11].

$$f_{w,z}(d,c) = \begin{cases} 0 & \text{if } c! = c^* \\ \frac{N(d,w)}{N(d)} & \text{otherwise} \end{cases} \quad (5)$$

the where "N(d,w)" represents the count of total number of times the word "w" occurs in a document "d" and "N(d)" represents the count of total number of words in "d". Maximum entropy uses the concept of bigrams and the features are added easily by maximum entropy without the concept of feature overlapping



1) **Support Vector Machine:** SVM considered as one of the most significant classification techniques where a set of vectors is taken as an input data with size “m”. Availability of a feature is represented as each entry in the vector. Each feature is formed by single word in the data for a unigram feature extractor. The value is assigned as “1” with the presence of feature, and in its absence the value is assigned as “0”. Overall processing speeds up as the presence of feature is used and there is no necessity to scale up the input data.

**SLDA [Sentiment Latent Dirichlet Allocation]:** SLDA for topic inference [8] neglects the positions of individual words. It specifies that all the sentences formed from the group of words are stipulated from a particular topic. The generative process for SLDA is as follows.

1. For an aspect “z”, generate a word distribution as “ $\varphi_z \sim \text{Dirichlet}(\beta)$ ”
2. For each aspect, for a review “d”
  - a. Generate the review’s aspect distribution “ $\theta_d \sim \text{Dirichlet}(\alpha)$ ”
  - b. For a particular sentence
    - i. Choose an aspect  $z \sim \text{Multinomial}(\theta_d)$
    - ii. Generate words  $w \sim \text{Multinomial}(\varphi_z)$

#### D. Performance Analysis

Sentiment Analysis techniques can be analysed by using techniques for accuracy calculation i.e., by calculating F-measure, recall, precision along with analysing the time taken for the execution of opinion mining techniques.

- 1) **Precision:** It is the count of significant documents recovered divided by the total count of recovered documents.
- 2) **Recall:** Recall is defined as the count of significant documents recovered divided by the total count of significant documents.

Suppose P= Count of retrieved relevant records.

Q= Count of relevant records not retrieved

R= Count of retrieved irrelevant.

$$\text{Precision} = \left( \frac{P}{(P+R)} \right) * 100\% \tag{6}$$

$$\text{Recall} = \left( \frac{P}{(P+Q)} \right) * 100\% \tag{7}$$

3) **F-measure:** The F-measure (F-score / F1-score) is a degree of accuracy and is characterized as the harmonic weighted mean of the test accuracy and recall.

$$F = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \tag{8}$$

#### IV. RESULTS AND OBSERVATIONS

All The 20news-bydate.tar.gz file is the input file that is collected from a civic website “<http://qwone.com/~jason/20Newsgroups/>” which contains both training and test data sets as shown below

```

# A tibble: 511,655 x 3
  newsgroup id
  <chr> <chr>
1 alt.atheism 49960
2 alt.atheism 49960
3 alt.atheism 49960
4 alt.atheism 49960
5 alt.atheism 49960
6 alt.atheism 49960
7 alt.atheism 49960
8 alt.atheism 49960
9 alt.atheism 49960
10 alt.atheism 49960
# ... with 511,645 more rows, and 1 more variables: text <chr>

```

Fig 2: Reading the data from different file

1) **Words in Newsgroups:** Most regular words are obtained from the entire dataset with the removal of the headers, signatures, and formatting. A word cloud of most regularly used words is generated which represents as a cloud of words as shown in below fig.



Fig 3: WORD CLOUD OF COMMON WORDS IN THE ENTIRE DATASET

To represent newsgroup the below output differs in the occurrence of words so as to quantify using tf-idf metric.

```

> tf_idf
# A tibble: 173,913 x 6
  newsgroup      word     n     tf     idf  tf_idf
  <chr> <chr> <dbl> <dbl> <dbl> <dbl>
1 comp.sys.ibm.pc.hardware  scsi   483  0.017616807  1.203973  0.02121016
2 talk.politics.mideast   armenian 582  0.008048902  2.302585  0.01853328
3 rec.motorcycles        bike     324  0.013898421  1.203973  0.01673332
4 talk.politics.mideast   armenians 509  0.007038932  2.302585  0.01620866
5                               encryption 410  0.008139651  1.897120  0.01548238
6                               nhl      157  0.004396651  2.995732  0.01317119
7 talk.politics.misc     stephanopolis 158  0.004162276  2.995732  0.01246906
8                               bikes     97  0.004160947  2.995732  0.01246508
9                               hockey   270  0.007561119  1.609438  0.01216915
10                               comp.windows.x  oname 136  0.003535498  2.995732  0.01059141
# ... with 173,903 more rows

```

Fig 4: Quantifying the data using tf-idf values

To remove words that are related to the topic, observe the top tf-idf values for certain groups by using ggplot as shown in fig below.

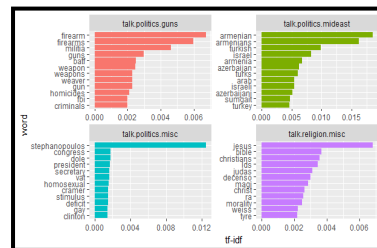


Fig 5: Top 12 terms with highest tf-idf for a newsgroup

In fig below portrays a pairwise correlation for four main clusters of newsgroups: computers/electronics, politics/religion, motor vehicles and sports.

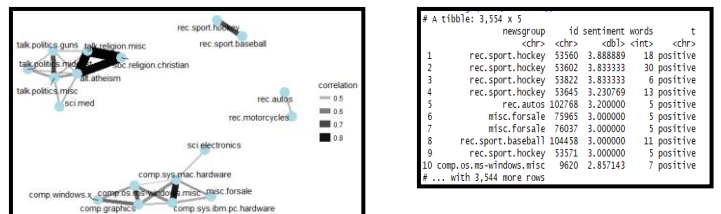


Fig 6: Pair wise Correlation and Visualized Stronger Correlations

**Topic Modelling:** First task in Topic Modelling is to divide the messages into four science-related groups and then generate a document-term matrix and model is fitted with LDA() function as shown in fig below

```
<<DocumentTermMatrix (documents: 2306, terms: 555)>>
Non-/sparse entries: 34946/1244884
Sparsity : 97%
Maximal term length: 14
weighting : term frequency (tf)
```

Fig 7: Document term frequency matrix generation

The mathematical matrix that depict the occurrence of terms that appear in a collection of documents is called a document-term matrix or term-document matrix is a. In this, rows match up to documents in the collection and columns match to terms as shown in fig below.

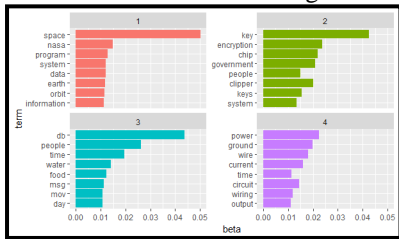


Fig 8: The top 8 words from each topic fit by LDA associated with science-related newsgroups

1) *Sentiment Analysis*: A score is assigned by AFINN lexicon for the words which includes the values 5 and -5, where the negative sentiment is replicated by negative scores and positive sentiment is replicated by positive scores. Sentiment analysis techniques are used to observe how often positive and negative words obtain in these Usenet posts. A positivity score for each word is obtained using AFINN sentiment lexicon and picture it with a bar plot as shown in fig below.



Fig 9: Positive score using AFINN sentiment lexicon

From this analysis, the “misc.forsale” newsgroup was the most positive as it included many positive terms regarding the products that a user wanted to put up for sale! Recognize why some newsgroups ruined up more positive or negative rather than others. Observe the total positive and negative assistance of each word as shown in fig below.

```
# A tibble: 1,909 x 4
  word occurrences contribution t
  <chr> <int> <int> <chr>
1 abandon 13 -26 negative
2 abandoned 19 -38 negative
3 abandons 3 -6 negative
4 abduction 2 -4 negative
5 abhor 4 -12 negative
6 abhorred 1 -3 negative
7 abhorrent 2 -6 negative
8 abilities 16 32 positive
9 ability 177 354 positive
10 aboard 8 8 positive
# ... with 1,899 more rows
```

Fig 10: Scores total positive and negative assistance of each word

Figs below represent the most positive and negative individual

messages, by grouping and summarizing by id rather than newsgroup

```
# A tibble: 3,554 x 5
  newsgroup id sentiment words t
  <chr> <chr> <dbl> <int> <chr>
1 rec.sport.hockey 53907 -3.000000 6 negative
2 sci.electronics 53899 -3.000000 5 negative
3 talk.politics.mideast 75918 -3.000000 7 negative
4 rec.autos 101627 -2.833333 6 negative
5 comp.graphics 37948 -2.600000 5 negative
6 comp.windows.x 67204 -2.700000 10 negative
7 talk.politics.guns 53362 -2.666667 6 negative
8 alt.atheism 51309 -2.600000 5 negative
9 comp.sys.mac.hardware 51513 -2.600000 5 negative
10 rec.autos 102883 -2.600000 5 negative
# ... with 3,544 more rows
```

Fig 11: List of top most Positive and Negative Messages

1) *N-Gram Analysis*: For example, a phrase like “don’t like” show the way to passages erroneously being labelled as positive. The Usenet dataset is a voluminous corpus of recent text; therefore, it could be attracted in how sentiment analysis may be inverted in this text as shown in figs below.

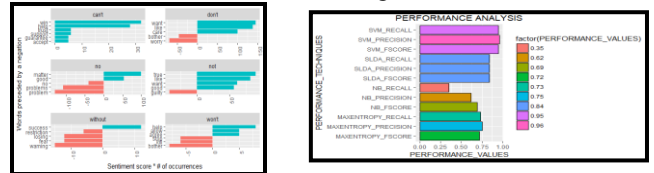
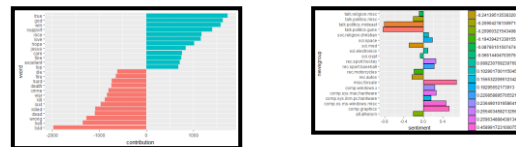


Fig 12: The most ‘negating’ sentiment word



Fig 13: Most common Words, Positive and Negative Words  
The most common word used in the entire dataset can be visualized using word cloud. Visualization can be done for the most frequently used negative and positive words, rather the size of words isn’t comparable for sentiments. The negative and positive word sets from the entire group of news data is generated and a plot is generated representing the positive and negative sets of data as shown in fig below.



An overall sentiment for each news group is then generated based on calculating each sentiment score for a news group as shown in fig below.

Topic	Total	Sentiment	Sentiment
1	talk	-0.48317649	negative
2	soc	0.10290170	positive
3	sci	-0.02249939	negative
4	rec	0.01011825	positive
5	misc	0.67283494	positive
6	comp	0.67283494	positive
7	alt	0.67283494	positive

Fig 15: Overall Sentiment of each news group

PERFORMANCE_TECHNIQUES	PERFORMANCE_VALUES
1	NB_PRECISION 0.62
2	NB_RECALL 0.35
3	NB_FSCORE 0.69
4	SVM_PRECISION 0.96
5	SVM_RECALL 0.95
6	SVM_FSCORE 0.95
7	SLDA_PRECISION 0.84
8	SLDA_RECALL 0.84
9	SLDA_FSCORE 0.84
10	MAXENTROPY_PRECISION 0.75
11	MAXENTROPY_RECALL 0.73
12	MAXENTROPY_FSCORE 0.72

Fig 16: Analysis measures for various machine learning techniques



Here positive and negative sentiments were categorized on the newsgroup dataset which is collected from different sources. The sentiments are measured using different machine learning techniques which are applied on news dataset. Predictions are done to determine which method best suits the classification. From identification of above measures, SVM method gives the best performance with 96%, 95%, 95% and 96% in terms of accuracy, F-score, recall and precision.

### V.CONCLUSIONS

With the increase of sentiment rich resources like online-news, personal blogs and individual web journals, new chances and difficulties emerge as people as of now will, and do, effectively utilize information innovations to chase out and see the opinions of others. Within the space of opinion mining, the major activity deals with opinion wrt machine and text judgement. In this paper different opinion mining procedures are proposed which centres around gathering information, a news data. Subsequently after gathering information it is changed into required arrangement. This information is pre-processed what's more, subjected to figure the opinion mining score utilizing different strategies. A higher level of accuracy is obtained for just a couple of the strategies can reach to a high-level accuracy. Hence, the results for opinion mining still have far to go before coming to the certainty level requested by down to earth applications.

### REFERENCES

1. M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-based methods for sentiment analysis," *Comput. Linguist.*, vol. 37, pp. 267-307, 2011.
2. R. McDonald, K. Hannan, T. Neylon, M. Wells, and J. Reynar, "Structured Models for Fine-to-Coarse Sentiment Analysis," in *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, 2007.
3. Johan Bollen, Huina Mao, and Xiaojun Zeng. 2011. Twitter mood predicts the stock market. *Journal of Computational Science*, 2(1):1-8
4. Bing Liu and Lei Zhang. 2012. A survey of opinion mining and sentiment analysis. In *Mining Text Data*, pages 415-463. Springer.
5. Thien Hai Nguyen and Kiyooki Shirai. 2013. Text classification of technical papers based on text segmentation. In Elisabeth Mtais, Farid Meziane, Mohamad Saraee, Vijayan Sugumaran, and Sunil Vadera, editors, *Natural Language Processing and Information Systems*, volume 7934 of *Lecture Notes in Computer Science*, pages 278-284. Springer Berlin Heidelberg
6. Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In *Proceedings of the Seventh Conference on International Language Resources and Evaluation (LREC'10)*, volume 10, pages 2200-2204.
7. Himabindu Lakkaraju, Chiranjib Bhattacharyya, Indrajit Bhattacharya, and Srujana Merugu. 2011. Exploiting coherence for the simultaneous discovery of latent facets and associated sentiments. In *Proceedings of the Eleventh SIAM International Conference on Data Mining, SDM 2011, April 28-30, 2011, Mesa, Arizona, USA*, pages 498-509. SIAM / Omnipress.
8. Chinthala S., Mande R., Manne S., Vemuri S. (2015) "Sentiment Analysis on Twitter Streaming Data" In: Satapathy S., Govardhan A., Raju K., Mandal J. (eds) *Emerging ICT for Bridging the Future - Proceedings of the 49th Annual Convention of the Computer Society of India (CSI) Volume 1. Advances in Intelligent Systems and Computing*, vol 337. Springer, Cham
9. Y. Sandeep, V. Sindhura (2015) "Methodological study of opinion retrieval techniques for Twitter social network" 2015 International Conference on Circuits, Power and Computing Technologies [ICCPCT-2015., ISBN: 978-1-4799-7075-9

10. Prasanna Komara, Geetha Guttikonda, Madhavi Katamaneni "Providing better medical healthcare services using data mining techniques" *International Journal of Innovations in Engineering and Technology (IJJET)* Volume 7 Issue 2 August 2016, ISSN: 2319 - 1058.
11. Vemuri Sindhura, Sandeep Yeliseti "Medical data Opinion retrieval on Twitter streaming data" 2015 IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), ISBN: 978-1-4799-6085-9

### AUTHORS PROFILE

**Paruchuri Ramya** Department of Information Technology, VR Siddhartha Engineering College, Kanuru, Vijayawada, India, #4 Project Associate, CTS, Chennai, Tamilnadu, India, paruchuriramya2010@gmail.com

**Geetha Guttikonda** Department of Information Technology, VR Siddhartha Engineering College, Kanuru, Vijayawada, India, Project Associate, CTS, Chennai, Tamilnadu, India, geetaguttikonda@gmail.com

**Vemuri Sindhura** Department of Information Technology, VR Siddhartha Engineering College, Kanuru, Vijayawada, India, Project Associate, CTS, Chennai, Tamilnadu, India, [vemurisindhura2233@gmail.com](mailto:vemurisindhura2233@gmail.com)

**Vinod Kumar Gadde** Department of Information Technology, VR Siddhartha Engineering College, Kanuru, Vijayawada, India, Project Associate, CTS, Chennai, Tamilnadu, India, g3vk0988@yahoo.co.in