# Automatic Detection and Classification of Brand Sentiments on Social Media using Machine Learning Algorithms

## Reshma Gulwani[1]* Rohit Singhal[2]

[1] Research Scholar, Sunrise University, Alwar Rajasthan, India

[2] Research Guide, Sunrise University, Alwar Rajasthan, India

*Abstract – Everyday a large amount of the data is generated by social media, blogs and other media on internet. This huge data contains opinions about the different topics or products or subjects. Opinion of people matters a lot to analyse how the spread of information influence the lives in a large-scale network like Twitter. Data generated by theses websites are unstructured and unorganized which requires processing to generate insights. Natural language processing is used to understand the structure and meaning of human language. Sentiment analysis is one of the major tasks of Natural language processing, where machine learning models are trained to classify text by polarity of opinion (positive, negative). Data used in this research are collected from twitter for reviews on railway services. Different machine learning techniques such as Naïve Bayes, Multinomial Naïve Bayes and support vector machines are used for training and testing the data. The performance of these models are evaluated and compared by using accuracy, precision, recall and F-measure.*

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - x - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

## I.    INTRODUCTION

Indian railways play very important role in transportation. People generally prefer travelling by train. One of the biggest advantages of using public transport is that it may help in reducing city traffic jams and pollution. Nowadays people are using social media sites like Twitter, Facebook to express their views about the public transport service. So it generates large amount of the data. To understand the public opinion towards railway services, sentimental analysis of that data can be very helpful for decision making.

Twitter is a microblogging and social networking service founded in 2006 on which users post messages or tweets. Tweets are limited to 140 characters, leaving 20 characters for the username. Twitter users may subscribe to the tweets posted by other users, an action referred to as "following". The service can be accessed through the Twitter website or through applications. Twitter users have adopted different conventions such as replies, retweets, and hashtags in their tweets. Twitter replies, denoted as @username, indicate that the tweet is a response to a tweet posted by another user. Retweets are used to re-publish the content of another tweet using the format RT @username. Hashtags are used to denote the context of the message by prefixing a word with a hash symbol. t.co is a URL shortening service created by Twitter. [1] It is only available for links posted to Twitter and not available for general use. [1] All links posted to twitter use a t.co wrapper. [2] Twitter hopes that the service will be able to protect users from malicious sites, [1] and will use it to track clicks on links within tweets. Words and phrases that are frequently used during a particular time period are known as "trending topics". These topics are listed by the platform for different regions of the world, and can also be personalize to the users [3].

## II.    OBJECTIVES

- To process and classify huge amount of the data generated by twitter.

- Automatically perform sentiment analysis on twitter posts

- To classify opinions in twitter text messages into categories like "positive" or "negative".

- Carry out the experiment to evaluate the performance of different classification models such as Naïve Bayes, Multinomial Naïve Bayes and support vector machines in terms of accuracy.

## III.    LITERATURE SURVEY

N. T. Renukadevi, S. Nandhinidevi, S. Karunakaran & R. Santhosh Kumar proposed a method for analysing nature of tweets on movie reviews. It also analyse polarity score in noisy twitter streams by means of two way classifier using Naive Baye's algorithm. Hash tag is keyword which is to be given as an input. Adeborna, E., Siau[5] performed Sentiment analysis on hotel reviews. Support vector machine is used to perform sentiment analysis on hotel reviews. They have used unigram feature along with TF-IDF (term frequency-inverted document frequency). They have described that use of TF-IDF is more effective than just frequency of words. Topic recognition model is proposed to perform the sentiment analysis on airline reviews to investigate and detect the polarity and the sentiments from the given text. Sachin Kumar, Marina,I. Nezhurina proposed sentiment Analysis on tweets for trains using machine learning techniques such as techniques such as support vector machines (SVM), Random forest (RF) and back propagation neural networks (BPNN). These techniques are used to analyze the hidden sentiments from tweet data. Rashmi Thakur, M.V. Deshpande [7] proposed a novel approach for sentiment classification on Train Reviews. Map reduce concept is used for sentiment classification of train reviews. [40]        Sahar A. El Rahman and Feddah Alhumaidi AlOtaibi and Wejdan Abdullah AlShehri[8] performed sentiment analysis of Twitter data on  two subjects McDonalds and KFC to show which restaurant has more popularity. Different machine learning algorithms were used .Various testing metrics like cross validation and f-score were used for testing the model.Adam Bermingham and

Alan F. Smeaton[9] provides a comparison of techniques of sentiment analysis in the analysis of political views by applying supervised machine-learning algorithms such as Naïve Bayes and support vector machines (SVM). Abdur Rasool et al[10] provides text mining and document based sentiment  on processed twitter data through machine learning techniques Sentiment analysis were performed on public opinion about international top two apparel international brands such as adidas and Nike. Opinions are categorized based on positive and negative attitude of common users about each brand. Positive reviews of Adidas are more than the Nike while there is the slight difference in negative reviews. Niu, Z., Yin, Z., & Kong, X [11] performed semantic orientation of documents, sentence, and words by using features like Part-of Speech, negation, Term frequency, Term Presence, and n-gram. Machine learning techniques such as Naïve Bayes (NB], Maximum Entropy (ME) and Support Vector Machines (SVM) are used to classify the text. Different machine learning techniques were compared and Bayesian algorithm provides better performance over the other models. Barbosa, L., & Feng, J. [12] used model for weight computation, classification and feature selection. The weights of

the classifier are adjusted by using the unique feature and representative feature. Information which distinguishes classes is known as 'Unique feature' and information representing a class is called 'Representing feature.' The performance of the Bayesian classification can be increased by using weights. 2-step automatic sentiment analysis is designed for classifying tweets .The first step is to classify the tweets as objective and subjective tweets and to next step is to classify the subjective tweet as positive or negative. Celikyilmaz, A., Hakkani-Tür,

D., & Feng, J [13] developed  a new method for noisy tweets which was pronunciation based clustering. The words with similar pronunciation are clustered in this approach. This can also be used to assign similar tokens to target organizations, user identifiers, HTML links and numbers used in text processing. Wu, Y., & Ren, F.[14] proposed a Very different probability for text classification of the tweets in which finding @username signifies a retweet and is considered as an influencing action which contributes to influencing probability. There is a correlation between these probabilities. N-gram and POS- tagger features were used for model. Learning based algorithm like the Multinomial Naïve Bayes classifier using Twitter corpus is used for sentiment classification.

So, there are many researches are performed on twitter sentiment analysis for different topics such as as product, movie, hotel reviews movie, sports etc. But in this research, automatic sentiment analysis is performed on twitter data for reviews on Indian railway seva. Machine learning techniques are used.to classifies the twitter data and finally we compared these classification techniques.

## IV.    PROPOSED METHODOLOGY

The primary intention of this research is to design and develop sentiment classification approach for Indian railway seva using machine learning Techniques In this research the work, first data is collected from twitter and split these into training set and testing set. Then pre-processing these data so that these can be fit for feature extraction, then classification techniques applied on the pre-processed dataset. Three classification techniques are applied and compared each other with accuracies. Following is the step by step procedure of using Natural Language Processing in sentiment analysis of railway text data on twitter.
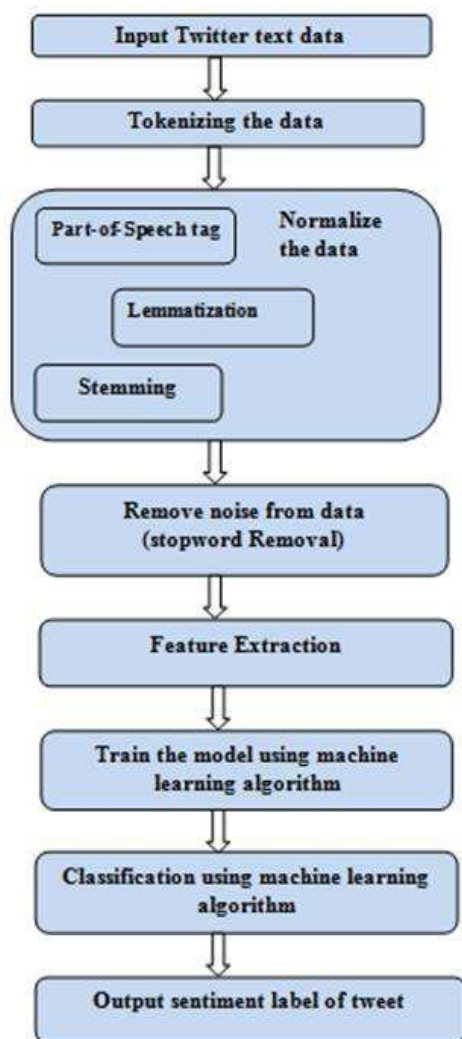
**Reshma Gulwani[1]* Rohit Singhal[2]**

**Fig.1: Sentiment Analysis of Twitter data using Natural Language Processing and Machine Learning**

### A. Dataset Description

Twitter API is used to collect tweets from a specific time period, user, or hashtag. Tweepy is used for fetching the data-client for Twitter Application Programming Interface (API)[15]. It can be installed using pip command: pip install tweepy. To fetch tweets from the Twitter API an app needs to be registered through your Twitter account. After creating API access, we collect customer key, customer secret key, access token key and access secret key and we collect the data set with the help of these keys.

### B. Data Pre-processing

Before applying any machine learning algorithm for sentiment analysis, it is important to do text pre-processing. It is essential to reduce the noise in human-text to improve accuracy. Data is processed with the help of a natural language processing pipeline. To make the language understand by machine, first task performed on text is tokenization i.e. splitting strings into smaller parts called tokens. Normaization techniques such as stemming and

lemmatization are used to convert the words to base form. Stop words are removed while processing language, because they are generally irrelevant.

### C. Feature Extraction

To analyze a pre-processed data, it needs to be converted into features. Depending upon the usage, text features can be constructed using different techniques. In this research, word density and TF-IDF techniques are used. Term frequency-Inverse Document frequency is an efficient approach. TF-IDF is a numerical statistic that reflects the value of a word for the whole document (here, tweet).

### D. Machine Learning Techniques

The most conceptual view that performs sentiment analysis using machine learning can be shown as in Fig .2. Natural language processing transforms the text in something a machine can understand using text vectorization, then machine learning algorithms are fed training data and expected outputs labels to train the machines to make associations between a particular input and its corresponding output.



**Fig.2: Sentiment Analysis using Machine Learning**

### • Multinomial Naive Bayes

Naive Bayes is one of probabilistic algorithms for sentiment analysis classification. The Naive Bayes Classifier is a well-known supervised machine learning classifier with applications in Natural Language Processing (NLP). It assigns a probability that a given word or phrase should be considered positive or negative. Naive Bayes algorithm is based on the Bayes rule, which can be represented as follows:

Posterior = likelihood * proposition/evidence

$$P(B_j|A) = \frac{P(B_j)P(A|B_j)}{P(A)} \quad (1)$$

The Naïve Bayes Classifier algorithm can be divided into two types i.e. multivariate Bernoulli and multinomial Naïve Bayes [4]. In this research, model of multinomial Naïve Bayes is used since it assumed the mutual independence of each word for all classes and $P(x1,…,xn)=1$ or constant.

$$P(c_i|x_1,…,x_n) = P(x_1,…,x_n|c_i)P(c_i). \quad (2)$$

The posterior probabilities values can be obtained based on the probability of each class (prior probabilities) and the probability of each word

**Reshma Gulwani[1]\* Rohit Singhal[2]**

(conditional probabilities) in the training data by simplifying above equation as follows [6].

$$P(c_i|x_1, \ldots, x_n) = \prod_{j=1}^{j=n} P(x_j|c_i) P(c_i) \qquad (3)$$

In the Naïve Bayes Classifier, the testing data enter the class $ci$ that has a maximum posteriori (MAP) or $cMAP$. The calculation of $cMAP$ value is defined as follows:

$$c_{MAP} = \underset{c_i \in C}{\mathrm{argmax}}\, P(c_i) \prod_{j=1}^{j=n} P(x_j|c_i) \qquad (4)$$

With the prior probability values as follows:

$$P(c_i) = \frac{N_{c_i}}{N} \qquad (5)$$

Where $Nci$ is the amount of training data that has class $ci$ and $N$ is the number of data used in the training data. The conditional probabilities values as follows:

$$P(x_j|c_i) = \frac{n_j}{n} \qquad (6)$$

Where $n_j$ is the number of occurrences of the word $x_j$ in class $ci$ while $n$ is number of words contained in class $c_i$.

Sometimes there are words that never appear in any of the classes during the classification process so that the resulting $(x_j|ci)$ value is zero. To prevent the occurrence of division by zero, then Laplace smoothing is used by adding the word frequency as much as 1(add-one) so that the calculation of $(x_j|ci)$ becomes

$$P(x_j|c_i) = \frac{n_j + 1}{n + n_k} \qquad (7)$$

where $n_k$ are the number of different (unique) words that appear in the training data.

•        **Support Vector Machine**

SVM was introduced by Vapnik in 1992. It was originally designed for linear classification and also be extended to multi-class by combining multiples SVMs for non-linear classification. It gives an excellent result for text categorization tasks such as sentiment analysis [50][51][52].This is an efficient algorithm for regression as well classification purpose. It draws a hyperplane to separate classes. The training samples are represented as points in the feature space.SVM performs classification by separating the points with a set of margin planes.
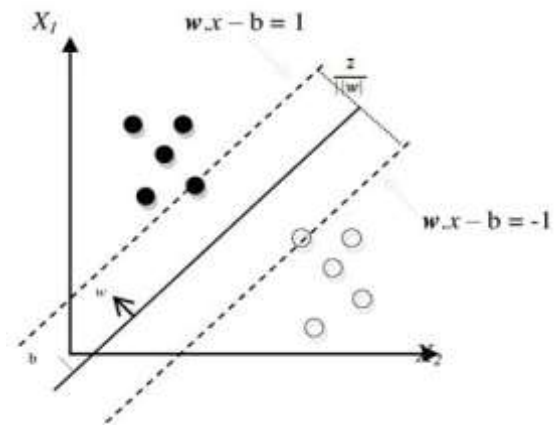


**Figure 3: Hyperplane in SVM**

The boundary hyperplane is chosen which maximizes the distance to the training samples. Support vectors are the points that determine the margin planes. Data which can be separated linearly is classified using Linear kernel and the data which is not linearly separable is classified using RBF kernel.

Decision function that separates two classes is defined as

$$f(x) = w^T x + b \qquad (8)$$

Where $w$ is the weight factor and $b$ is the bias of the function.

The optimal hyperplane is defined as

$$|w^T x + b| = 1 \qquad (9)$$

$x$ is a training example. The distance between $x$ and $wT$, $b$ais calculated as

$$Distance = \frac{|w^T x + b|}{\|w\|} \qquad (10)$$

The margin $MM$ is the distance with closeness,

The kernel function is a dot product of data inputs. The frequently used kernel functions are Linear, Polynomial, RBF and sigmoid. The linear kernel function is

$$M = \frac{2}{\|w\|} \qquad (11)$$

This algorithm works extremely well with regression, the effect of SVM increases as we increase dimensional space. SVM also perform well when the dimension number is larger than the sample number [49]. There exists a drawback too it does not perform well on huge datasets. SVM extensively uses cross-validation to increase its computational efficiency. It gives best results than Naive Bayes algorithm. The

**Reshma Gulwani[1]\* Rohit Singhal[2]**

basic idea is to find the hyperplane which is represented as the vector w which separates document vector in one class from the vectors in other class.

A support vector machine is another supervised machine learning model, similar to linear regression but more advanced. SVM uses algorithms to train and classify text within our sentiment polarity model, taking it a step beyond X/Y prediction.

## V.    RESULTS AND DISCUSSIONS

The dataset consists of tweets about Indian railway seva. These data are pre-processed in filtering tokens, stop words and negative words. After that, the system performs feature selection and extraction process. The extracted features will be given as input to Multinomial Naïve Bayes and Support Vector Machine classifiers. In order to view the efficiency of sentiment classification of tweets using machine learning approach, accuracy is measured on twitter data set of different sizes. As the objective of this research, is to analyze the sentiments of tweets posted for Indian Railway Seva. Tweets are collected for a specific time period, user, or hashtag with the help of Twitter API. Following are the sample data sets.





**Figure 4: Data set Sample**

**Table I: Accuracy of Multinomial Naïve Bayes on Twitter data set**

| Training Size | Accuracy |
|---|---|
| 1000 | 60 |
| 2000 | 62.5 |
| 3000 | 63.4 |
| 4000 | 64.6 |
| 5000 | 65.09 |
| 6000 | 67.5 |
| 7000 | 69.6 |
| 8000 | 72 |
| 9000 | 74.5 |
| 10000 | 75.03 |

Table I and Figure 4 shows the accuracy of Multinomial Naive Bayes on twitter data set of different training size
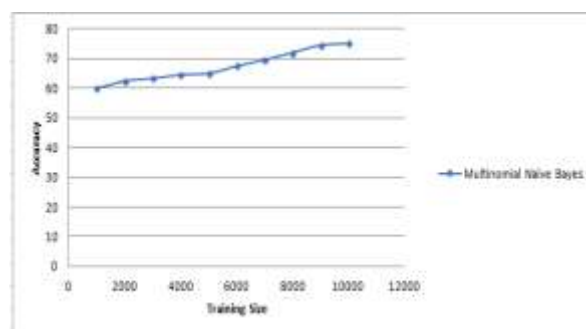


**Figure 4: Graph Representing Different results obtained for Multinomial Naïve Bayes Algorithm on Twitter Data Set**

Table II and Figure 5 shows the accuracy of Multinomial Naive Bayes on corpus data set of different training size

**Table II: Accuracy of Multinomial Naïve Bayes on Corpus data set**

| Training Size | Accuracy |
|---|---|
| 10000 | 58 |
| 20000 | 64.6 |
| 30000 | 66.4 |
| 40000 | 68.5 |
| 50000 | 69.7 |
| 60000 | 70.4 |
| 70000 | 74.6 |
| 80000 | 78.3 |
| 90000 | 80.5 |
| 100000 | 82.6 |



**Figure 5: Graph Representing Different results obtained for Multinomial Naïve Bayes Algorithm on Corpus Data set**

In the Classification using Support Vector Machine algorithm, linear kernel function is used. ($cost$) is the parameter of the penalty of the error in the classification and the value is determined by the researchers. In this study the user determines the value $C$, the values are 0, 1; 1; and 10 for modeling the data classification of training.

**Reshma Gulwani[1]* Rohit Singhal[2]**

**Table III: Accuracy of Support Vector Machine on Twitter data set**

| Training Size | Accuracy |
|---|---|
| 1000 | 60 |
| 2000 | 63.4 |
| 3000 | 65.09 |
| 4000 | 66 |
| 5000 | 67.5 |
| 6000 | 68.5 |
| 7000 | 69.6 |
| 8000 | 72 |
| 9000 | 74.4 |
| 10000 | 75.06 |

Table III and Figure 6 shows the accuracy of support vector Machine on twitter data set of different training size
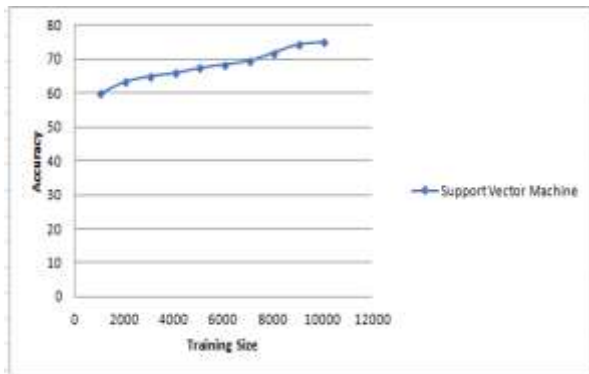


**Figure 6: Graph Representing Different results obtained for Support Vector Machine Algorithm on Twitter Data Set**

**Table IV: Accuracy of Support Vector Machine on Corpus Data Set**

| Training Size | Accuracy |
|---|---|
| 10000 | 58.02 |
| 20000 | 60.03 |
| 30000 | 64.05 |
| 40000 | 68.5 |
| 50000 | 69.3 |
| 60000 | 72.4 |
| 70000 | 74.6 |
| 80000 | 79.03 |
| 90000 | 82.4 |
| 100000 | 84.3 |

Table IV and Figure 7 shows the accuracy of support vector Machine on Corpus data set of different training size
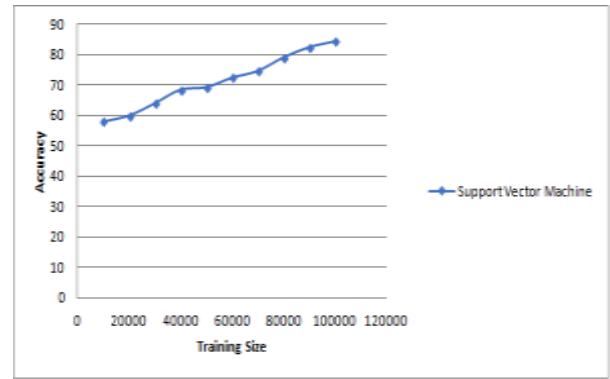


**Figure 7: Graph Representing Different results obtained for Support Vector Machine Algorithm on Corpus Data Set**

**Table V: Accuracy, precision and recall value with Multinomial NBC and SVM on Twitter Data set**

| Evaluation Parameters | | Algorithms | |
|---|---|---|---|
| | | Multinomial Naïve Bayes | Support Vector Machine |
| Accuracy | | 75.033% | 75.066% |
| Precision | Positive | 74% | 74% |
| | Negative | 76% | 76% |
| Recall | Positive | 77% | 77% |
| | Negative | 73% | 73% |
| F1 Score | Positive | 75% | 76% |
| | Negative | 74% | 75% |

**Table VI: Accuracy, precision and recall value with Multinomial NBC and SVM on Corpus Data set**

| Evaluation Parameters | | Algorithms | |
|---|---|---|---|
| | | Multinomial Naïve Bayes | Support Vector Machine |
| Accuracy | | 82.6% | 84.3% |
| Precision | Positive | 81. 57% | 83.24% |
| | Negative | 83.57% | 85.24% |
| Recall | Positive | 84.6% | 86.24% |
| | Negative | 80.6% | 82.24% |
| F1 Score | Positive | 83.05% | 84.7% |
| | Negative | 82.05% | 83.7% |

## VI. CONCLUSION

In this paper, many existing techniques for sentiment analysis using machine learning algorithms are presented. As day by day more people prefer to share the post on twitter. So it would be beneficial to analyze the twitter data set. Method for processing linguistic data set using machine learning algorithms has been given. Data cleaning, preprocessing and removing noise data such as hyperlinks or URL's, twitter handles in replies are preceded by a @ symbol and punctuation and special characters for e.g.: (\ | [ ];

**Reshma Gulwani[1]* Rohit Singhal[2]**

{} - + ( ) < >?! @ # % *) are demonstrated using Natural Language Processing (NLP) techniques which provide accurate data as well as reduce the size of dataset. NLTK toolkit is used to preprocess the twitter data. NLTK toolkit has different functions for tokenize each word in the tweet, which allowed performing unigram analysis of the word. Stemming and Lemmatization is used to obtain the root form of the word. WordNet is a lexical database for the English language that helps to determine the base word. Part-of-Speech tag is used to identify the context and relative position of the word in tweet and obtained the synsets term by analyzing the WordNet. TF_IDF and count features are extracted from preprocessed data. Since, the analyses assigned sentiment label to each tweet. That analysis helps to understand the sentiment of the user when reacting on the Twitter platform, which derives not only the opinion from the user but allow the business to know the feedback about the product or events. That feedback helps them to make intelligent decision in future. Indian Railways can be benefited by incorporating such text analytics techniques to handle the huge amount of reviews received for satisfying the customers. Further, analysis of data using machine learning techniques like Naïve Bayes, Multinomial Naïve Bayes and Support vector machines are performed for measuring consistency, accuracy and reliability of classified sentiment analysis data. According to the accuracy, the precision and the recall value, the performance of Support Vector Machine technique is better than Multinomial Naïve Bayes technique. The Visualization of the sentiment analysis result using python NLTK and Sklearn platform and compared the result using machine learning algorithm.

## REFERENCES

[1]     About Twitter's Link Service<http://t.co>".Twitter Help Center (module of Twitter).Archived from the original on Februray 25, 2011. Retrieved February 23, 2011..

[2]     Garrett, Sean (June 8, 2010). "Links and Twitter: Length Shouldn't Matter". Twitter Blog (blog of Twitter). Retrieved February 23, 2011.

[3]     https://support.twitter.com/articles/101125.

[4]     Pang B. & Lee L. (2008) Opinion mining and sentiment analysis. Found Trends Inf Retr 2(1-2): pp. 1–135

[5]     Adeborna, E., Siau, K. (2014). "An approach to sentiment analysis – the case of airline quality rating". In: PACIS 2014 Proceedings, Paper 363, Chengdu, 24–28 June (2014).

[6]     Wilson, T., J. Wiebe, and R. Hwa (2006). Recognizing strong and weak opinion clauses. Computational Intelligence, 22(2): pp. 73-99.

[7]     Rashmi Thakur, M.V. Deshpande (2017). "A Novel Approach for Sentiment Classification on Train Reviews" International Journal of Advanced Research in Basic Engineering Sciences and Technology (IJARBEST). ISSN (Online) : 2456-5717 Vol. 3, Special Issue 37

[8]     Sahar A. El Rahman and Feddah Alhumaidi AlOtaibi and Wejdan Abdullah AlShehri (2019). "Sentiment Analysis of Twitter Data" International Conference on Computer and Information Sciences (ICCIS), pages: 1-4

[9]     Adam Bermingham and Alan F. Smeaton (2014). "On Using Twitter to Monitor Political Sentiment and Predict Election Results" https://www.researchgate.net/publication/ 267250109

[10]    Abdur Rasool et. al. (2019). "Twitter Sentiment Analysis: A Case Study for Apparel Brands "2019 J. IOP Conf. Series: Journal of Physics: Conf. Series 1176, 022015 IOP Publishing doi:10.1088/1742-6596/1176/2/022015

[11]    Niu, Z., Yin, Z., & Kong, X. (2012). "Sentiment classification for microblog by machine learning." Proceedings - 4th International Conference on Computational and Information Sciences, ICCIS 2012, pp. 286–289.

[12]    Barbosa, L., & Feng, J. (2010). Robust Sentiment Detection on Twitter from Biased and Noisy Data. Coling, pp. 36–44

[13]    Celikyilmaz, A., Hakkani-Tür, D., & Feng, J. (2010). Probabilistic model-based sentiment analysis of twitter messages. 2010 IEEE Workshop on Spoken Language Technology, SLT 2010 - Proceedings, pp. 79–84.

[14]    Wu, Y., & Ren, F. (2011). Learning sentimental influence in twitter. Proceedings - 2011 International Conference on Future Computer Sciences and Application, ICFCSA 2011, pp. 119–122.

[15]    M. B. Myneni., L. V. N. Prasad and J. S. Devi (2017). "A Framework for Sementic Level Social Sentiment Analysis Model", JATIT. vol. 96, no. 16, pp. 1992-8645

**Reshma Gulwani[1]* Rohit Singhal[2]**

**Corresponding Author**

**Reshma Gulwani***

Research Scholar, Sunrise University, Alwar Rajasthan, India