# REVIEW ON SENTIMENT ANALYSIS OF TWEETS CLASSIFICATION BY MACHINE LEARNING

Hemant Kumar Menaria [1] , Pritesh Nagar [2] , Mayank Patel[3]
*[1, 2,3] Geetanjali Institute of Technical Studies College, Udaipur, INDIA*

***Abstract-*** Sentiment analysis has turned out one of the most significant tools in natural language processing because it opens up numerous possibilities to understand people's sentiments on different topics. The purpose of an aspect-based sentiment analysis is to understand this further and find out what someone is talking about, and whether he likes it or does not like it. A real-world example of the perfect realm of this topic is the millions of available Indian welfare plans like Swachh Bharat Abhiyan and Jan Dhan Yojna. These welfare schemes were launched by the government at all levels in schools, states and center level. The schemes work with the collaboration of center and state government. These welfare schemes are introduced for the different levels according to the peoples and their lifestyle. The welfare scheme mostly introduced to develop the weaker and minority section of the society. Some schemes are introduced only for women and girls like Beti Bachao, Beti Padaho and Ujjwala Yojna. These schemes empower the women by helping them financially and providing the basic need facilities. There have been various ways to deal with handle this issue, utilizing machine learning. In this thesis, labeled data is used on the basis of polarity and Tweets preprocess and extract unigram features after preprocessing of the tweets.

***Keywords-*** *Sentiment analysis, Support Vector Machine, Random Forest*

## I. INTRODUCTION

The process of sentiment analysis refers to emotions, feelings, attitude or opinion. With development of technology based on World Wide Web, an individual or a group often explicit their opinions or sentiments over the technology of internet with the help of reviews, blogs, and ratings. With the help of such a rise in textual form of data, the organization feel the desire to evaluate and figure out this type of text and further evaluate the business-based insights. Advertising companies and business owners often engage the process of sentiment analysis in order to explore new strategies of business and campaign of advertising. The algorithms of machine leaning are mainly used for classifying and predicting the negative or positive kind of sentiments [1] [2]. The most famous categories of machine learning algorithm are the supervised and un-supervised categories. The supervised form of algorithm mainly uses a dataset of labelled form where each of the trained set document is marked or labelled with an appropriate form of sentiment. Whereas, the case of unsupervised learning includes an unlabelled form of dataset [3]. The algorithms of unsupervised learning are more of more complex form and it needs an additional algorithm of clustering in its initial implementation phase. This kind of study usually concerns with the techniques of supervised learning on a labelled form of dataset. The technique of SA can be operated on three distinct levels such as, aspect level, sentence level, and document level [10]. In case ofclassification of sentiments based on sentence level, there occurs an individual sentenced based polarity of a document whereas in case of classification base on aspect level, firstly it helps in identifying distinct forms of aspects. Document Level mainly aims in classification of the overall document or the topic as negative or positive.

### 1.1 Sentiment Analysis

The concept of sentiment analysis is understood by combining the terms "Senitiment" and "Analysis". The word sentiment represents feeling that can be joyful, confusing, irritating, distracting. The sentiments are the feelings based on certain attitudes and opinions rather than facts due to which sentiments are of subjective nature. The sentiment implies an emotion usually motivated by opinion or perception of a person. The psychologists attempts to present multitude of emotions classified into six distinct classes: joy, love, fear, sadness, surprise and anger. The emotions based on sadness and joy are experienced on daily basis at different levels. We are mainly concerned about sentiment analysis detecting a positive or a negative response or opinion [3]. The major significance of sentiment analysis is that every emotion is linked to human perception forming an ingrained part of all humans which means that every human has the potential to generate different opinions acting as a tool for sentiment analalysis. Sentiment analysis refers to the analysis automation of a known text determining the distinct types of feelings conveyed. The term sentiment analysis and opinion mining can be used interchangeably. Sentiment analysis as defined as an information extraction and natural language processing task with an aim to gain the feelings of writer expressed positively or negatively based on requests, comments or questions analysing large data-sets or documents. It basically intends to define writer's feeling regarding a specific topic

based on writer's own opinion. It models a branch that can help in providing a judgement over distinct fields. The measurement of sentiments is a biased technique with it is really complex to achieve high accuracy of automated systems [2] [5].

## 1.2 Online Social Media

In this modern era, social media is everywhere and for everyone, it is now such a hefty part of all the generations of life. Social media are web-based communication tools that empower people to discuss, share and consume their everyday activities on different social media which are encounters formally or informally.



Figure 1: Social Media[16]

There are diverse online networking platforms available, for example, twitter, face book, YouTube, Flicker and so forth. It includes gatherings and blogging, writing news articles which allow people to take part in healthy discussions over social media. An evolution of human social interaction by social media mass adoption .This provides a platform of great importance for users to express their opinions, feelings, issues, joy, views, struggle and emotions. Most of the persons share information and make connections all over the world through social media. On the personal level, online networking provides facilities to connect with friends to do innovative things, pick up information in the educational field, to create interests in new things, and as a source of entertainment. On the other side at the professional level, user utilize online networking to create a network of experts with professional from all, over the world to widen knowledge insight in a specific field.  At the organization level, online networking permits them to have a live communication with their audience even from distant areas, get feedbacks of customer and can promote their brand. Social media[5] having lots of unique features, some of them are as follows:

- Immediacy – This feature provides users can control the discussions and able to quickly react to wrong data.

- Interaction – Online networking likewise encourages association and engagement between the group and organization, enhancing correspondence and connections.
- Audience – Through cell phones and PCs different group of peoples can create numerous online networking locales.
- Scalable – Online networking can be custom-made to meet the one of a kind needs of any association and can be incorporated into an office's correspondence and effort methodology.

1.2.1 Characteristics of Social Media: They are usually (conventional) unbind services.

- It provides a venue for online, interpersonal or mass communication.
- It permits individuals from all edges of the world the chance to communicate with each other.
- To the fullest extent of these locales offer easy to use includes features that constantly updated themselves.
- They can specialize in a solitary reoccurring topic or conceal a huge number of various themes.
- It permits individuals to advance themselves while making an online fingerprint 'of who they are, in addition to what they may aspire to be.
- It provides new freeway for communication and give-and-take ideas.

Permanent availability of posted data on social media destinations makes it a persistence platform.

- It provides the facility of replicability, accessibility and search ability.

1.2.2 Types of Social Media: A wide range of types of online networking is everywhere throughout the Internet. Users have various motivations to utilize these online networking outlets. Different types of social media are as follows [7]

1. Blogs: A blog is a dedicated streamline for personal diary or journal. It provides a platform to express thoughts and passions to the world.

2. Wikis: It is an aggregate site that allows everyone to create or modify a page using her Web browser. Actually, a wiki is a blend of a CGI script and a gathering of plain content documents that permits clients to make Web pages.

3. Social bookmarking: This is the demonstration of storing bookmarks with their keywords to an open Web site. It is especially valuable when gathering an arrangement of assets which are to be imparted on others.

*4. Social network sites:* The Social network sites abbreviated as SNS is the expression used to portray any website that permits people to grab the benefits of social networking services.

- *Status-update services:* Also called micro blogging administrations, for example, Twitter allows users to post short tweets regarding people, events and to see updates created by others.
- *Virtual world content:* This is innovation which moves clients to an alternate, envisioned and genuine, platform. This can be achieved through either 360 video catch or CGI generation. Imaginary universe is one of the examples of virtual world in which users create avatars interacting with each other.
- *Media-sharing sites:* It permits clients to post recordings or photographs on YouTube, Pinterest etc. Users can then share that media with the world or just a select group of friends.

### 1.3 Types of Sentiment Analysis

There are three different types of sentiment analysis

*1. Document-Level Sentiment Analysis:* In this type of sentiment analysis is performed on the entity on which and identify the positive negative view on the single entity by using the documents [10]

*2. Comparative Sentiment Analysis:* In many cases users express their views by comparing with the similar product or entity. The main goal is to identify the opinion from the comparative sentence. For e.g.: "I drove the Verna, it does not handle better that Honda city superb."

*3. Aspect-Based Sentiment Analysis:* Document level and sentence level analysis gives good results when they are used on single entity but when we want to analyze the multiple entity then we need aspect based sentiment analysis. For e.g.: "I am a Samsung phone lover. I like the look of the phone. The screen is big and clear. The camera is fantastic. But, In any case, there are a couple of drawbacks as well; the life of battery isn't up-to the check and access to what's app is troublesome." Classifying the positive and negatives of this review hides the valuable information about the product .In the aspect based sentiment analysis it allows us to analyze the positive and negative aspect of an item. This type of analysis is mainly domain specific. In this analysis firstly aspects are identified and then location of the aspect in the review. After this polarity of the view defines the aspect.

### 1.4 Methodology of Sentiment Analysis

Below steps presents the methodology of sentiment analysis.

**1.4.1 Pre-processing**: This process was done before the tweet-based usage of feature extractor in order to design or build the feature vector. The process is initialized using the following steps. Such steps convert the plain tweet text into the elements of processing nature with an additional

information utilized by the feature extractor. The tools of third party were used for all the steps specifically handling tweeted text unique nature [14].

*1. Step1: Tokenization:* It is a process of text conversion as a string into elements process-able known as tokens. In terms of tweets methodology, such elements can be emoticons, words, links, punctuations or hashtags. As shown in figure.3, "an insanely awsum time..." text was busted into "an", "insanely", "awsum", "time". The elements here get separated by some space whereas the sentence based punctuation ending such as full stop or exclamation mark get separated more often by a space. The hastags along with symbol "#" that precedes the tag is required to get retained as the symbol "#" may suggest distinct sentiments than the word to be used regularly in text. Therefore, the Twitter-based particular form of tokenizer helps to extract tokens.
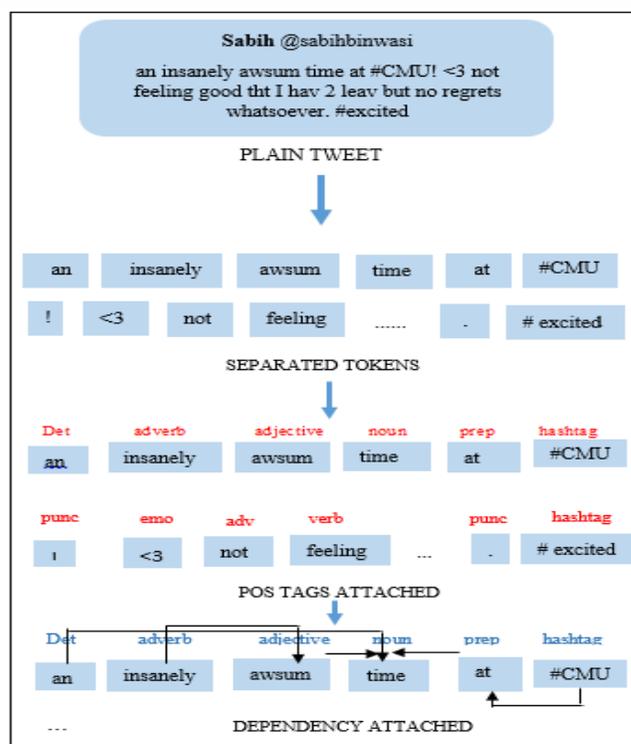


Figure.2 Pre-processing stages

*2. Step 2: Parts of Speech Tags (POS):* The POS tags are generally characterized by sentence based words dependent over the different categories of grammar in context to a word. Such an information or data is necessary for the process of sentiment analysis as words contain distinct values of sentiment that are basically POS dependent. Considering an example, the word "good" acting as a noun has no sentiment whereas "good" in its adjective form reflects a sentiment in a positive way. Fig.2 above shows that each token gets

extracted in its last step and gets assigned a POS tag. The accuracy maintained by a POS tagger is about 93%.

*3. Step 2: Dependency Parsing:* It represents the relationship extraction among sentence-based words. Such a method is very useful for identification of relationship between "good" and "bad" in form of phrases like "not really good" where there is only adjacent word relationship. It explains the parent-child relation between tweet tokens as shown in fig.2. The accuracy maintained by dependency parsing is about 80%.

*1.4.2 Feature Extraction*

This is a process of designing a feature-vector from a known tweet. The entry (each) in case of feature vector is an integer that contributes on the attributing a class of sentiment to a tweet class. Such a contribution strong to negligible form. Here, the strong class represents feature based value entry influencing sentiment true class whereas the negligible class presents no such relation between sentiment class and feature value  [13] [14]. The algorithm here identifies the strength dependency between classes and the features using strong correlated form of features and preventing noisy-feature usage. *(a) Feature Set-Bag of Words:* These represents "bag or words" called unigrams as a set of features where the token frequency is considered to be a feature vector. This feature set was unanimously selected by the practitioners. The feature vector based entry is assigned to each of the specific tokens that were found in the trained labelled set. If such a token occurs in tweet, then it gets assigned by 1 as a binary value, otherwise it is considered to be zero. As shown in figure above, the order of token-based sequence or the structure of grammar is not at all preserved. For instance, the token "awsum" forms a part of tweet hence it is marked or assigned as "1" whereas the token "hater" does not occur in tweet, hence it is labelled as "0". The token "hater" has a column for it in a feature vector. Hence, it would take place in some tweet in the trained data-set form.

It was further analyzed that the indication of only word presence yields good performance than the word-based frequency. In such case, an entry is also assigned for specific ordered token pairs termed as bi-grams. The token pair "insanely awsum" where it is assigned or labelled as "1" if it forms a part of tweet, otherwise it is considered to be "0". It indicates that the system is equipped not only to indicate the token presence but also indicate its context. *(b) Feature Selection:* It is expected to add unigrams, bigrams and trigrams adding large entries to a vector feature which can make the space of vector highly dimensional resulting in a more complex and hard task to identify the relation among each feature. This represents the major issue of text-based classification and is popularly known as Curse of Dimensionality [1] [2]. The process notices that some of the features are not relevant for the operation of sentiment analysis, hence these un-necessary features are required to be

removed. So, different researchers conducted the study of selecting the features as per the requirement. A feature attribution evaluation was conducted to analyses the impact of features. A very popular method named Chi-Squared Feature Selection was used which carried a classification algorithm evaluating the feature-based dependency value and the dependency of each of the class. Thus, if the feature has a high correlational dependence, then it is assigned with a very high rank. The method of Chi-Squared Feature Selection outperform well for text classification.
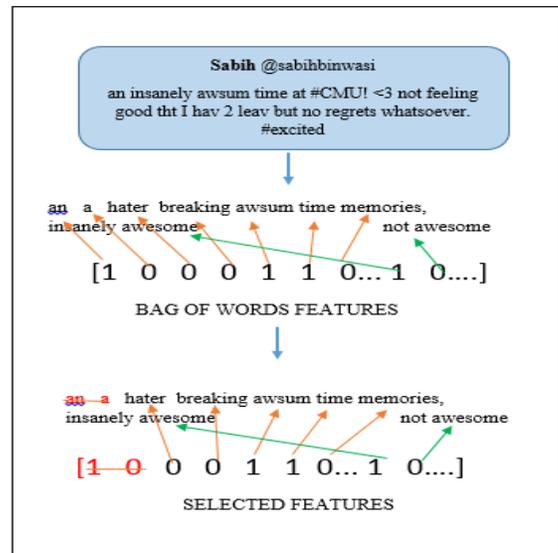


Figure.3 Feature Extraction

*(c) Social Media Feature Set:* Emoticons are used as symbols basically to express the gesture or feelings using language characters and punctuation. In fig. the emoticon "<3" depicts heart representing love symbol [16]. These emoticons strongly analyses the positive or negative sentiment. The figure below shows a tweet sample illustrating the social media feature set representing a number of positive emoticons labelled as "1" as per the dictionary while no such emoticons are found. In feature-based vector, the presence of hashtag, URL and mention ("@username") is also included. Such types of features are mostly used in the showing the diffrenciating
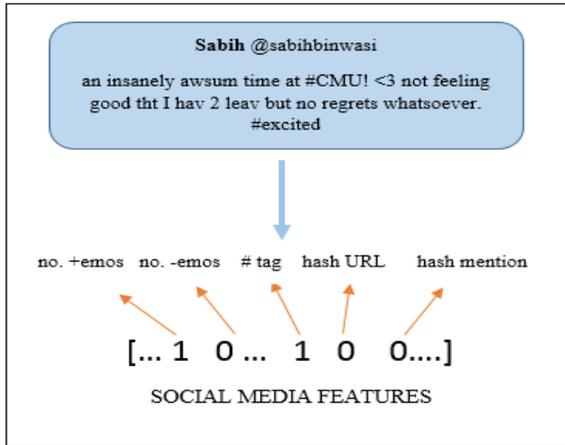
Figure.4 Social Media Feature Set

impact between the polar i.e. positive and negative and non-polar i.e. neutral class. These type of indicators helps to identify the relation between sentiment class and the indicators itself. For, instance, a tweet with "@username1" is considered as a rare form of token in the trained data set, the algorithm of classification would fail in identifying a relation. With "@" considered to be as the binary feature, the chances for relation formation increases.

*(d) Lexical Feature Set:* These are basically driven by the use if lexicons. The task of sentiment lexicon analysis maps the n-grams or tokens to score-based polarity. These have been reported successfully with an ability to locate issues based on classification methods of sentiment analysis. After the process of bag-of-words features.
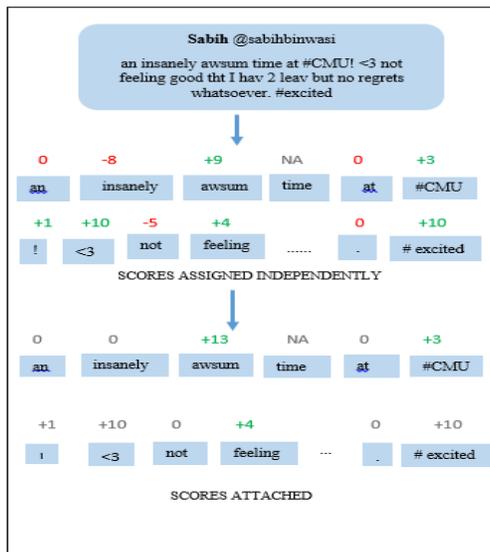


Figure.5 Lexical Feature Set

Mohammad et al., 2013 analyzed such a feature to be most successful one, solving the issue of identification of full-text sentiment, token sentiment using the lexicon-based feature set which makes the task very easy. We take an example of *AFINN* which is an affective lexicon discovered by Finn Årup Nielsen. Such type of lexicon was constructed manually for the mapping of the most frequently used words on Twitter having a ranges from -5 to +5, where, -5 is labelled as a negative type of token and +5 is considered as the most positive type of token. The tokens which does not contain any kind of sentiment are labelled as "0". Each of the tweet token is labelled independently in a polarity score. For instance, "time" is not labelled any type of score.

### II. MECHANISM OF DATA MINING

This process of data mining refers to the core of Knowledge Discovery Process (KDD). It is comprised of few steps which lead to new knowledge from raw data collections.

- *Data integration:* At this step, numerous heterogeneous information sources collect in a common source and collection of relevant prior knowledge is done. It required to formulating goals which the process have to achieve. This process may lead to initial data preparation [2] [4].

- *Data selection:* During this step, the data pertinent to the analysis is settled and recovered by information gathering. Information can be gathered from the data warehouse or from a transitory storage system.

- *Data pre-processing:* Data cleansing is another name for this process, it is a stage in which noisy information, anomalies, missing values and



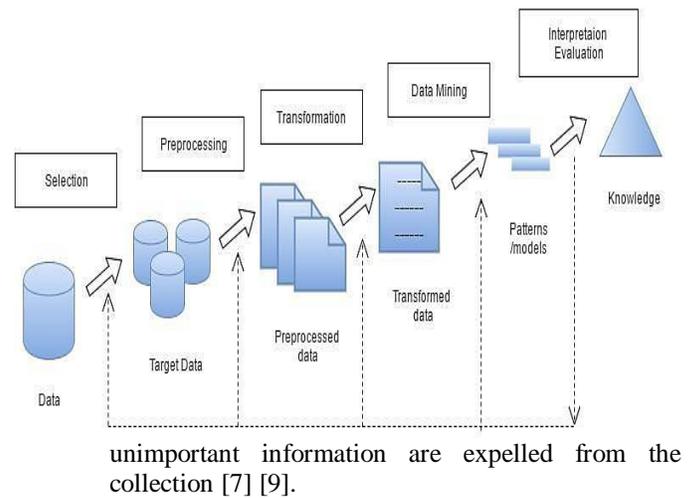unimportant information are expelled from the collection [7] [9].

Figure.6 Mechanism of Data Mining

- *Data transformation:* It is also called data combination; a stage which produces appropriate

forms of data is gained from the selected. This step leads to feature selection and reduces dimensions.

- *Data mining:* It is an important step which is used for extraction of potentially useful patterns by applying clever strategies. At this stage, different mining algorithms like Associations, sequences, classification, clustering, etc. are applied here.
- *Pattern Interpretation and evaluation:* In this progression, strictly interesting patterns in light of given measures representing knowledge are recognized.
- *Knowledge representation:* At this last stage, the discovered knowledge is outwardly shown to the client. This crucial step utilizes visualization techniques to help clients comprehend and decipher the obtained results.

### 2.1 An Architecture for Data Mining

The vital elements of a classic data mining system [17] are as follows and shown in Figure 1.2.

*1. Knowledge Base:* For entire data mining process, the knowledge base is very helpful. It may be search or evaluate the patterns for results. It may have data from experiences of users which can be utilized for the data mining process. To generate more reliable and accurate results the inputs can be gained by the data mining engine from the knowledge base. The modules of pattern evaluation interface with this phase to furthermore update it.

*2. Graphical User Interface:* This module intercommunicates between client and data mini framework. This module provides the better utilization of the framework effectively and productively extracting the real complexity of the procedure. At the point when the user client determines an inquiry, this module displays the outcome in an easily understandable manner by interacting with the mining framework.

*3. Pattern Evaluation Modules:* By utilizing a threshold value the measurements of interestingness of the pattern are done by pattern evaluation module. To focusing on the search for finding interesting patterns this module associates with the data mining engine [8].

*4. Data Mining Engine:* It is the core part of any data mining system. Numerous functions of data mining additionally classification, clustering, time-series analysis, association, characterization, prediction and so forth are performed by various modules of mining engine.

*5. Database or Data Warehouse Server:* It contains the realistic data that is to be processed. Hence, considering the data mining request of the user then relevant data is obtained by a server to process the request.

*6. Data Cleaning, Integration, and Selection:* To pass the content to data warehouse server, the data is required to clean, integrate and select. The data available in different sources and formats cannot be utilized straight forwardly for the process of data mining the reason behind it is that the unavailability of complete and reliable data. Along these lines, first data required cleaning and integration. After that, a larger number of data more than required gathered from various data sources from which the data of interest should be chosen and passed to the server. As a component of cleaning, integration and selection of data various techniques may be performed [1] [17].

*7. Database, Data Warehouse, World Wide Web, or other Information Repositories:* These are the real sources of data. For effective data mining, we require large volumes of historical data. A large amount of data is usually stored by organizations in databases. Warehouses may contain many databases, content records, and different sorts of data archives spreadsheets. Sometimes, even plain content records or spreadsheets are saved. A Huge source of data example includes www or the internet [13] [16].

### III. ALGORITHM USED

***1. Naïve Bayes classifier***: It is based on the Bayesian theorem of probability. It is a Supervised Learning algorithm which is used for classification. It solves the problem in continuous as well as categorical value attributes.

- It is mainly used in text classification and spam filtering.
- Recommendation based systems also used this classifier.

Bayes's theorem is the base of the classifier on which it is developed. In this theorem conditional probability that an event y belongs to class n can be calculated from the conditional probabilities of finding particular event in each class and the unconditional probability of the event in each class. In given data y $\varepsilon$ Y and C classes where y is a random variable. Conditional probability that an event y belongs to class n is computed by using following formula (1)

$$P(C_n|y) = P(C_n)\frac{P(y|c_n)}{P(y)} \dots\dots\dots\dots\dots\dots (2)$$

This equation is used for pattern classification and finds the probability of the given data y belongs to class n and gives the optimum class with high probability among all classes.

For Statistically Independent Features perform the given equation:

$$P(y|C_n) = \prod_{i=0}^{p} P(y|C_n) \dots\dots\dots\dots\dots\dots\dots\dots\dots$$
(2)

Here, y is a p-dimensional vector data y= $(y_1, y_2, \dots\dots y_n)$
Summarized the Naïve Bayes result:

$$n = argmax_n\, P(c_n) \prod_{i=0}^{n} P(y_i|C_n) \dots\dots\dots\dots\dots\dots$$
(3)

***2. Random Forest:*** It uses decision tree models for classification and regression. It gives the more accuracy and

information in variables results. In this classifiers features are selected randomly from the training data set while maintain the class distribution. It splits the data set by using the Gain Index method. This method also estimated the missing data and run very efficiently on large size data set. Random forest method is used for feature selection and support vector machine for the classification of these features. Accuracy in result is 85% in identification.

If classification is done by using random forest method then following steps are performed on image. In this consider the feature space of M-dimensional $\{A_1, A_2 \ldots \ldots \ldots A_m\}$. Here we calculate the weights $\{W_1, W_2 \ldots \ldots W_m\}$ for every feature in the space. Improved algorithm is used to grow each decision tree in random forest using weights.

*3. Feature Weight Computation:* To register the feature weight, it is utilized to gauge the usefulness of each information feature A as its connection to the class Y. The high estimation of the weight shows the question in preparing information corresponds with the estimations of features. Accordingly, A is useful to the class names of new objects.

*4. Support Vector Machine:* Vapnik introduced for support vector machine, and is popular tool for supervised machines learning methods which are based on the minimization of the structural risk. The SVM basic characteristics is the original non-linear data into data class and the separation margin among itself is maximized and typing points nearer from the support vectors.

The training sample is
$$n = \{(u_i, v_i)|i = 1, 2, \ldots., m\}$$
Where
m←Sample no.
$\{u_i\} \in r_k$←Input vector set
$v \in \{-1,1\}$←Desired corresponding input vector
Then, optimal classification of existing hyper-plane has following condition to meet:
$$\begin{cases} \omega^t u_i + B \le 1, v_i = 1 \\ \omega^t u_i + B \le -1, v_i = -1 \end{cases}$$
Where
$\omega^\tau$←Super plane omega vector,
B←offset quality
Then, the decision function is classified as:
$$F(u_i) = sgn(\omega^t u_i + B)$$
SVM classification model is described with optimization model $min_{\omega,\xi,B} P(\omega, \xi)$

$$min_{\omega,\xi,B} P(\omega, \xi_i) = \frac{1}{2}\omega^t\omega + \frac{1}{2}\gamma \sum_{i=1}^{m} \xi_i^2$$
$$v_i[\omega^t\phi(u_i) + B = 1 - \xi_i, i = 1,2, \ldots., m$$
$$\xi = (\xi_1, \xi_2, \ldots. \xi_m)$$
Where
$\xi_i$←slack variable

B←offset
$\omega$←support vector
$\gamma$←classification parameter for balancing the model complexity and fitness error.
Transforming the optimization problem into dual space and for solving it, Lagrange function is introduced:

$$l(B, \omega, \alpha, \xi) = \frac{1}{2}\omega^t\omega$$
$$+ \frac{1}{2}\gamma \sum_{i=1}^{m} \xi_i^2$$
$$- \sum_{A=1}^{m} \alpha_i\{v_i[\omega^t\phi(u_A) + B] - 1 + \xi_i\}$$

Where
$\alpha_i$←Lagrange multiplier
Then, describing the classification decision function:
$$F(x_i) = sgn(\sum_{i=1}^{m} \alpha_i v_i A(u, u_i) + B)$$

## III. RELATED WORK

Chen et.al. [1] proposed a visualization approach called TagNet which is used for sentiment analysis. This approach combines the tag clouds with improved node-link diagrams which presents the time-varying heterogeneous information. The proposed algorithm improves the scalability of large dataset. Krishna, et al. [2] proposed a model on fuzzy logic for feature based opinion mining and sentiment analysis. The opinions are used for making decisions to choose a product or any interesting topic. The proposed approach is used to extract the features from the tweets. This is completed by using machine learning and fuzzy approach. The classification of sentiment analysis and review is effectively done by this approach. Shidaganti et.al. [3] proposed a technique which is basically a combination of data mining and machine learning. The proposed work is done on the tweeter data for analysis of tweet for gather the user's opinion regarding any topic of issue. Basically the tweeter platform is used by peoples to express their view in short message related to brands, celebrities, products and also on political issues. In this work TF-IDF and clustering algorithm is discussed with their efficiency. Rout et al. [4] worked upon the unstructured data of social media like blogs, tweeter for sentiment and emotion analysis. This work is performed on supervised and unsupervised approach on different databases. The unsupervised approach is used for automatic identification of sentiments for the tweets. Different algorithm of machine learning like SVM and maximum entropy are used to identify the sentiments. The unigram, bigram and POS features are used effective feature form the tweeter. Mumtaz et al. [5] proposed an approach which is a combination of machine

learning and lexical based approach. The proposed hybrid approach gives high accuracy than the classical lexical method and provides the enhanced redundancy than machine learning approach. The proposed approach is used for opinion mining by using the concept of natural language processing which extract the sentiments from the text related to an entity. Sharma, et al. [6] outlines the features of the various establishments; achievement is marked after so many years of independence. Many people are living far from financial inclusion. The PMJDY has made the mark of success and try to overcome the loopholes of the initiatives which are in continuous sequence. PMJDY goal is to cover all households of the country and provide the banking facilities along with inbuilt insurance coverage. The research has been surveyed in the rural areas of Jaipur district. In this to find the facts and figures both the primary and secondary sector data has been collected and try to find the correlation between them. It is also seen that nearness of banks increase the likelihood inclusion. In the coming time it will see that how much awareness is spread by the direct and indirect channels, & how the policy makers help to increase the interests. Jones, et al. [7] focused on the implication of PMJDY scheme for the development all over the India. With this step of Prime Minister, India is digitizing even the rural areas are conscious about their bank accounts and to enjoy benefits provided by the government. The various initiatives are taken Jan Dhan Yojna, Swachh bharat. Demonetization of currency all these measures lead towards the progress in the structure of Indian economy. This scheme ensures the better quality of life in the country and helps in improve the living status of the people. India is the only nation where 50% of people is in the working age group. Balachandran et al [8] outlined the today's scenario about choosing the right institute is the most challenging task. To get the estimated idea about the particular institute every student surf over the social network sites for the reviews, ratings about the particular institution. But it is difficult to analyze the statistical aspect from the reviews. In this Aspect based Sentiment Analysis is directly implemented on the reviews which gives us negative and the positive reviews of the particular institution. The various techniques are used for the aspect identification such as NLP-based technique, Machine Learning based (ML), unsupervised approach, Dictionary based, Corpus based. The best possible result analytic is given by the NLP and the ML classifier to classify each aspect into their respective category. Al-Smadi et al. [9] worked on aspect based sentiment analysis on the reviews of hotel's in Arabic by using the concept of Long-short term memory neural networks. It is implemented in two levels that are character level and aspect based with random field classifier and polarity classification based. The proposed approach gives better results by the enhancement of 39%. Zahrotun et al. [10] has discussed about the clustering

technique that are used in the data mining. Clustering technique is the grouping of meaningful data that belongs to the same class or comes under the same group. Jaccard and cosine similarity, and the combination of both are used to get the optimum value similarity. Cosine similarity that gives the measures of similarity within two non-zero vectors of inner product space that measures cosine angle between the product space. But from the combination of these it is expected it increases the values by two titles. This work is carried out over the practical work in the form of title document in the Department of Informatics Engineering University of Ahmad Dahlan. The results obtained from this study are cosine similarity gives the good approximations than the Jaccard and from the combination of both. Agarwal.et al [11] has discussed about the documents that we get in the form of unstructured, semi-structured and structured data. To assemble the bulk of data and to classify them is very difficult according to their respective domain. To overcome this problem of domain specific clustering two algorithms are used Jaccard and cosine similarity algorithm techniques to find the similarity within two documents. Cosine similarity within two documents gives the fast result because the cluster generation is steady in comparison to the Jaccard coefficient. Jaccard coefficient uses more complex mathematical formulae for computing the similarity between two documents. So Cosine similarity gives us more accurate, reliable results. Virmani.et al [12] has discussed the collaboration of sentiment analysis with the opinion extraction, summarization and maintain the record of each and every student. To get the enhanced and collaborated opinion about the student modifies the existing algorithm in the paper. To analyze the opinion a database of sentiment word has been used. Firstly set the score to each and every sentiment word in the database. Whenever a sentiment word is encounter in the sentence, it matches with the database and set the score accordingly. Then from these scores the cumulative opinion value is evaluated. The algorithm gives us numerical value for the opinion. If the numerical score is high shows the positive remark and if the numerical score is

low shows negative remarks. If the remarks given by the two teacher are very high but at the same time remarks given by the one teacher is low then on collaboration it will gives average score. The overall performance depends upon the teacher remarks, the sentiment word used by the teacher is not match with word used in database it effects the overall score. Medhat, et al. [13] has discussed about the various Sentiment Analysis (SA) applications, recent enhancements in the algorithm that are investigated and present in the paper briefly. Recent articles are discussed in the proposed paper that gathers the interest of the readers in SA (transfer learning, emotion detection, building up of the resources). The surveys are conducted on various SA algorithms gives refined

categorization. Emotion detection algorithm is mainly used to enhance and analyze the emotions either it could be explicit or implicit. In this the various algorithm are used to represent the sentiments and emotions some are Pont-wise Mutual information (PMI), Chi-square, Latent Semantic Indexing (LSI). The sentiment classification techniques are segmented into machine learning (ML), hybrid and lexicon based approach. SC and FS algorithms are area of research in the proposed paper. The most commonly used ML algorithms are mostly used for solving SC problems. Bifet et.al. [14] This

paper focused on the challenges that Twitter information faces, concentrating on order issues, and afterward consider these streams for supposition mining and sentiment analysis. To manage gushing unequal classes, author proposed a sliding window Kappa measurement for assessment in time-changing information streams. Utilizing this measurement we play out an investigation on Twitter information utilizing learning calculations for information streams.

Table.1 Existing Scheduling Model

| Author's Name | Year | Methodology Used | Proposed Work |
| --- | --- | --- | --- |
| Shidaganti et al. | 2018 | Machine Learning | The proposed work is done on the tweeter data for analysis of tweet for gather the user's opinion regarding any topic of issue. Basically the tweeter platform is used by peoples to express their view in short message related to brands, celebrities, products and also on political issues. |
| Al-Smadi | 2018 | Deep Neural Network | It is implemented in two levels that are character level and aspect based with random field classifier and polarity classification based. |
| Sharma et al. | 2017 | Survey | In this to find the facts and figures both the primary and secondary sector data has been collected and try to find the correlation between them. It is also seen that nearness of banks increase the likelihood inclusion. In the coming time it will see that how much awareness is spread by the direct and indirect channels, & how the policy makers help to increase the interests. |
| Bifet et.al. | 2010 | Kappa Measurement | Focused on the challenges that Twitter information faces, concentrating on order issues, and afterward consider these streams for supposition mining and sentiment analysis |
| Zahrotun et al. | 2016 | Jaccard and Cosine Similarity | Discussed about the clustering technique that are used in the data mining. |

### IV. CONCLUSION

In this study, the concept of Support Vector Machines (SVM) is used for classification of algorithm with binary classification process. Such type of method helps in analyzing different feature vectors with an assigned class in order to identify the relation dependency between a sentiment and each of the feature. Here, each of the vector is considered as a point of data in vector dimensional space that equals to the size of feature-set. The SVM helps in identifying the vector dimension based hyperplane which divides the class into two types. One is the considered as "best" i.e. defined as

a good type of separation gained by the hyperplane having the large distance to the point nearest to the training data type of any kind of class known as functional margin. In general, if the margin is large then the classifier error gets reduced.

### V. REFERENCES

[1] Chen, Yang. "TagNet: Toward Tag-Based Sentiment Analysis of Large Social Media Data." In *2018 IEEE Pacific Visualization Symposium (PacificVis)*, pp. 190-194. IEEE, 2018

[2] Krishna, Kalpesh, Preethi Jyothi, and Mohit Iyyer. "Revisiting the Importance of Encoding Logic Rules in

Sentiment Classification." *arXiv preprint arXiv: 1808.07733* (2018).

[3] Shidaganti, Ganeshayya, Rameshwari Gopal Hulkund, and S. Prakash. "Analysis and Exploitation of Twitter Data Using Machine Learning Techniques." In *International Proceedings on Advances in Soft Computing, Intelligent Systems and Applications*, pp. 135-146. Springer, Singapore, 2018.

[4] J. K. Rout, K. K. R. Choo, A. K. Dash, S. Bakshi, S. K. Jena, and K. L. Williams, "A model for sentiment and emotion analysis of unstructured social media text," *Electron. Commer. Res.*, vol. 18, no. 1, pp. 181–199, 2018.

[5] D. Mumtaz and B. Ahuja, *A Lexical and Machine Learning-Based Hybrid System for Sentiment Analysis*, vol. 7. Springer, Singapore, 2018.

[6] N. Sharma, "Pradhan Mantri Jan Dhan Yojana ( PMJDY ) - A Conceptual Study," *Int. J. Res. Granthaalayah*, vol. 5, pp. 143–152, 2017.

[7] T. M. Jones, "A Study on the Implications of Pradhan Manthri Jan Dhan Yojana on the Growth of Indian Economy," vol. 06, no. 03, pp. 461–466, 2017.

[8] L. Balachandran and A. Kirupananda, "Online reviews evaluation system for higher education institution: An aspect based sentiment analysis tool," *2017 11th Int. Conf. Software, Knowledge, Inf. Manag. Appl.*, pp. 1–7, 2017

[9] M. Al-Smadi, M. Al-Ayyoub, H. Al-Sarhan, and Y. Jararwell, "An aspect-based sentiment analysis approach to evaluating Arabic news affect on readers," *J. Univers. Comput. Sci.*, vol. 22, no. 5, pp. 630–649, 2016.

[10] L. Zahrotun, "Comparison Jaccard similarity, Cosine Similarity and Combined Both of the Data Clustering With Shared Nearest Neighbor Method," *Comput. Eng. Appl.*, vol. 5, no. 11, pp. 2252–4274, 2016.

[11] N. Agarwal, M. Rawat, and M. Vijay, "Comparative Analysis Of Jaccard Coefficient and Cosine Similarity for Web Document Similarity Measure," *Int. J. Adv. Res. Eng. Technol.*, vol. 2, no. 5, pp. 18–21, 2014.

[12] D. Virmani, V. Malhotra, and R. Tyagi, "Sentiment Analysis Using Collaborated Opinion Mining," *arXiv Prepr. arXiv1401.2618*, no. January, 2014.

[13] W. Medhat, A. Hassan, and H. Korashy, "Sentiment analysis algorithms and applications : A survey," *Ain Shams Eng. J.*, vol. 5, no. 4, pp. 1093–1113, 2014.

[14] Bifet, Albert, and Eibe Frank. "Sentiment knowledge discovery in twitter streaming data." In *International conference on discovery science*, pp. 1-15. Springer, Berlin, Heidelberg, 2010.

[16] Michell's, "What is Social Media?" [Online]. Available: https://www.tes.com/lessons/OioCS800E23GLQ/is-society-becoming-addicted-to-social-media. [Accessed: 29-May-2018].

[17] H. Ashida and T. Morita, "Architecture of Data Mining Server : DATAFRONTBerver." [Online]. Available: https://data-flair.training/blogs/data-mining-architecture/. [Accessed: 15-May-2018].