

# A Survey on Text Content Sentiment Analysis Levels and Techniques

**Mahesh Patidar**

Research Scholar School of Data Science  
& Forecasting, Devi Ahilya  
Vishwavidyalaya, Indore,  
India

**Dr. V.B. Gupta**

Head School of Data Science &  
Forecasting, Devi Ahilya  
Vishwavidyalaya, Indore,  
India

**Seema Patidar**

Department of Electronics &  
Communication, Oriental Institute of  
Science and Technology Bhopal,  
India

**Abstract** – As Individual convey their thoughts on social media by discussing personal life events or other issues related to politics, science, environment, etc. So this produce large amount of opinion data which can be further mine by utilizing some techniques of pattern, classification, etc. So researcher works on this profitable data, and consequently, semantic annotation of online digital content into respected sentiment was done by utilizing various techniques. Hence this paper gives a detail investigation of different level of sentiment analysis with necessary pre-processing steps. Sentiment analysis techniques adopt by different researcher like supervised, unsupervised, lexicon, etc. were also discuss with their desired outcome. A deep survey of social content sentiment analysis work done by different researcher was also discussed for improving the knowledge of methods used.

**Keywords**– Classification, Sentiment analysis, Ontology, Text Mining, , Un-directed Classification.

## I. INTRODUCTION

As social media attract people towards internet, so number of websites provides different services through which people can share their thoughts, feeling, around the globe. Social media enables people to be linked together and interact with each other anywhere and anytime [1].

Social Media provides a various method for many people to precise and gives the opinion on a current or past event and many other activities around us [2]. More than 500 million people in the world give their opinion and views daily on the web [3]. A large quantity of knowledge is generated from various and different social media in numerous formats, numerous languages in the world. It makes challenges in data analytics to search out the new purpose and extract information from it [4, 5]. The Social Media mining is not exclusively Machine Learning strategies or other intelligent methods to identify and extracts information for sentiment analysis. On other hands, it is essential to know and determine the various domain for which unrelated gathering knowledge from the different place within the word, different time zone, language and separate values to be analyzed from entirely different perspectives [4].

Sentimental analysis primarily focuses on the recognition and categorization of opinions. It can be performed in two ways; using knowledge based approach and the other machine learning techniques [5]. In the first approach,

there is a requirement of the huge database consist of already defined emotions and effective representation of the knowledge. The machine learning approach utilizes the trained and test datasets to design a classifier. Therefore it is easier than the knowledge based approach. There are various types of challenges which were found in the field of sentimental analysis [2]. One of the major challenges is that the word which expresses the opinion may be positive or negative relying upon the circumstances. There are various types of challenges which were found in the field of sentimental analysis. One of the major challenges is to classify the word which expresses the opinion may be positive or negative relying upon the circumstances. Another major challenge is that the way of expressing the opinion of people may not be the same. The relation between textual reviews and the consequences of those reviews can be obtained using opinion mining.

## II. LEVELS OF SENTIMENT ANALYSIS

Analysis of sentiment has been explored primarily at three levels [6]. They are Document level, Sentence level, Entity and Aspect level.

### 1. Document-level Sentiment Analysis

In [7] 159 articles were distributed based on the granularity of sentiment analysis, 73 articles appeared in the document level which makes it the most studied topic in the field. Document-level sentiment classification, as known in the literature, is considered the simplest sentiment analysis task. The task of opinion mining at this

level is to identify opinionated documents and classify them according to their polarities. Authors in [8] mentioned that most of the researchers at this level follow a twostep approach: Topic Relevance Retrieval and Opinion Finding step. The document is considered as a basic information unit which includes multiple sentences. Based on the quintuple introduced in the first section, the task of opinion mining is to determine the overall sentiment of the opinion holder about the entity described in the document. This approach helps the users in decision making by providing a summary of total number of positive and negative documents.

## 2. Sentence-level Sentiment Analysis

Just like document level opinion mining, sentence level opinion mining is also a classification problem. Sentences are regarded as short documents which makes the classification the same for both levels. Most of the researches at document level don't perform a three class classification (positive, negative, and neutral). However, at the sentence level, the neutral class cannot be ignored because sentences may express no opinion or sentiment. Thus, the purpose at this level is to classify each sentence in an opinion document as positive, negative or neutral opinion or sentiment. Sentiment sentence classification is generally performed in two classes of classification problem. The first determines whether the sentence is expressing an opinion (sentiment) or not and the second classify the sentences as positive, negative or neutral. The first step in the process is known in the literature as Subjectivity Classification.

It aims to distinguish opinions (subjective sentences) from facts (objective sentences) [9]. Subjective sentences can express some personal feelings, views judgments or beliefs that might vary from person to person, whereas, objective sentences express factual information which remains valid for all individuals. Because of that some researchers prefer to classify sentences as opinionated or non-opinionated. The second step is called sentence sentiment classification. After classifying the sentences as being subjective (opinionated) or objective (nonopinionated), sentence sentiment classification aims to classify them as positive, negative or neutral. An assumption that generally researchers make at this level of analysis is that a sentence expresses a single sentiment. Thus, sentences that express more than one sentiment are treated differently. More complex sentences (interrogative, comparatives, conditionals and sarcastic sentences) also need advanced techniques.

## 3. Aspect-Level Sentiment Analysis

Polarity classification of opinion text at document and sentence level is helpful in many cases but it does not provide all the necessary details because they do not discover what exactly people liked and did not like. Generally, documents are made of several passages of

opinions of different semantic categories. Thus, classification at coarse-level does not identify sentiments or opinion targets. For example, being positive/negative of the sentiments about an entity in a text document, do not mean that the author is being positive/negative about all the aspects of the expressed entity. Due to the need of a finer grain analysis, aspectlevel sentiment analysis represents a key step. Aspectlevel sentiment analysis (previously called featurebased sentiment analysis) describes that an opinion consists of a sentiment and a target. The objective of the analysis at this level is to discover the specific targets and then specify their sentiment polarities. Using the quintuple definition ( $ei, aij, sijkl, hk, tl$ ), aspect-level sentiment analysis aims to locate the first three components. Therefore the analysis is divided into two tasks: Aspect extraction and Aspect sentiment classification. The first task is also called opinion target extraction [10], because it concentrates on the extraction of both entities and their aspects. Entities appoint to products names, services, events, etc, and aspects, which can be expressed implicitly or explicitly, generally identify the attributes and components of entities. The second step, similar to the identification of the polarity of opinions at coarse granularity, associates a polarity with the various extracted opinion targets. The extraction of remaining components of the quintuple are studied as sub-tasks of aspect-level sentiment analysis called opinion holder extraction and time extraction. The extraction of all quintuples present in a document is helpful to produce a summary of opinions about entities and their aspects.

## III. RELATED WORK

In [5] authors propose a novel strategy to do the sentiment analysis for news article dataset. Based on the content present on online social media inform of text, emojis, etc. for a news incident, a word emotion association network (WEAN) was developed to express its semantic and emotion, which establishes the framework for the news incident sentiment calculation. In light of WEAN, a word emotion calculation is proposed to get the underlying words emotion, which are additionally refined through the standard emotion ontology. With the words emotion close by, author can process each sentence's sentiment. Here paper has found the class of word emotion, while this takes to reduce the sentence sentiment identification because use of words which are not of same class takes to wrong analysis.

In [11] authors propose a structure for feature based emotion examination alongside the sentence compression method. Feature based emotion / sentiment investigation is performed dependent on syntactic which represents an opportunity for over regular issue. This sort of issue makes the emotion / sentiment examination too hard to even think about handling the syntactic parsers utilized in

the emotion / sentiment mining strategy. The proposed structure builds up an enhanced sentence compression procedure before the emotion / sentiment examination. For compacting a content for emotion investigation two plans are utilized. That is syntactic compression and extractive compression procedure. Contrasted with extractive compression strategy syntactic is viewed as increasingly proficient on the grounds that it pack the content by evacuating the insignificant words. The proposed method utilizes Aspect-Polarity (A-P) gathering based emotion investigation. The vast majority of the viewpoint put together notion investigation center with respect to the connection between the Features and the extremity words which incredibly influences the effectiveness. To take care of this issue the proposed structure utilizes syntactic examples.

In [12], authors proposed another characteristic type to check its commitment in document level emotion / sentiment investigation. They additionally got outcome on dataset containing 233600 comments with 93.24% exactness. In this paper an exploratory investigation on emotion / sentiment extremity characterization had been led by them. Most importantly, a rating based component had been portrayed which depended on regression model, AND gained from outside autonomous dataset of 233600 film comments. Now commitment of both artificial model and rating based criteria were utilized to accomplish exactness of 91.6% and 89.87% on the datasets from various areas. These outcomes demonstrated that rating based element was increasingly proficient for emotion / sentiment arrangement on extremity comments. Execution could likewise be improved by including bigram and trigram features.

In [13], authors showed content based methodology for the online reviews, film evaluations and so forth utilizing emotion / sentiment examination. These surveys were gathered by administered AI methodologies. For analysis three diverse AI computations were considered for example SVM, ME, NB and these methods evaluation depended on parameters, for example, precision, f-measure and accuracy. In this paper for arranging the film comments of spoiled tomatoes dataset utilizing n gram strategy diverse AI procedures had been proposed. The creators additionally presumed that on examination with other research works their results got better precision.

In [14] authors, developed emotion / sentiment examination on tweets for demonetization. Initially work start with collection of data and afterward changed over it into content documents as inserting dataset. At that point emotion / sentiment examination was performed in the wake of excluding the stop words pursued by deciding the class of the words and ordering the tweets as positive and negative. So another strategy was proposed for emotion / sentiment investigation on demonetization and for this

procedure information cleaning, bigrams, class, emotion scores and graphical strategies were utilized. Here work does not specify the sub-class of the sentiments, as only positive and negative class were identified.

In [15], authors developed comparison examination of various methodologies for emotion / sentiment investigation and theme recognition of Spanish tweets was given arrangement undertakings. For ordering Spanish tweets as per assessment and themes different analyses had been performed. Utilization of stemmers and lemmatizers, ngrams, word types, refutations, valence shifters, interface handling, web indexes, unique Twitter semantics (Hashtags), and distinctive arrangement techniques had been assessed which was spoken to a point by point and complete investigation. The primary end that Anta et. al attracted was that because of their curtness and absence of setting tweets were difficult to manage. These outcomes demonstrated that for examining and arranging the Spanish content it was conceivable to utilize old style techniques. Best exactness that was seen was 58% for subjects and 42% for notions grouping.

In [16], Basha et. al presented that as the presence of E-business item reviews for an item were likewise developing quickly with an exponential factor. To settle on a choice among numerous choice where time and cash were valuable, other individuals emotions would play a significant job. Presently the majority of the associations had emotion / sentiment mining and slant investigation as a piece of their examination. Additionally, pretty much every business was affected by the online networking sites and web journals which drove these organizations to do nostalgic investigation. In this paper they had utilized fluffy standard based frameworks (FRBS) with models, in particular: Mamdani, and Takagi Sugeno Kang (TSK) utilizing FRBS bundle in R. They likewise contrasted these models and other grouping techniques as far as exactness, Recall and F-measure, precision and execution of the strategy. Examinations on the proposed calculation for computing emotions and assessments with respect to the item were directed and furthermore showed the R bundle. Likewise a few instances of the utilization of the bundle and correlation with different bundles had been made.

## IV. CLASSIFICATION TECHNIQUES

### 1. Lexicon-based approach

The main procedure that can be utilized for emotion / sentiment analysis is the dictionary based strategy. It utilizes a dictionary that comprises of terms with separate emotion / sentiment scores to each term. The term can be related with a solitary word, expression or figure of speech [10]. The emotion / sentiment is characterized dependent on the nearness or nonappearance of terms in the vocabulary. The vocabulary based methodology

incorporates corpus-based methodology and word reference based methodology that is examined further.

**Dictionary-based approach:** The principle thought behind the word reference based methodology is to utilize lexical databases with assessment words to remove emotion / sentiment from the report. In view of [17], a lot of seed estimation words (for example great, terrible) with their polarities are gathered by hand. Toward the start, this underlying set does not need to be huge, 30 sentiment words is sufficient [18]. Subsequent stage is to utilize the polar words to enhance a set by searching up for separate equivalent words and antonyms in a lexical database. The look-into technique is iterative. At every emphasis the calculation takes refreshed arrangement of words (extended set) and searches again until there will be no new words to incorporate. At last, a lot of emotion words can be inspected with a motivation behind erasing mistakes.

## 2. Corpus-based approach

In [19] Bing Liu shows that corpus-based methodology can be connected in two cases. First case is a recognizable proof of emotion words and their polarities in the area corpus utilizing a given arrangement of emotion / sentiment words. The subsequent case is for structure another vocabulary inside the specific area from another dictionary utilizing a space corpus. The discoveries propose that regardless of whether emotion / sentiment words are space subordinate it can happen that a similar word will have inverse direction relying upon setting. The examination directed by Hazivassiloglou and McKeown [20] is noticeable in the writing about corpus-based system. Creators proposed a technique that concentrates semantic direction of conjoined descriptive words from the corpus. The procedure depends on the use of literary corpora and seed emotion / sentiment words (descriptors). Uncommon etymological standards are connected to the corpora so as to find emotion / sentiment words with relating polarities. Creators accept that modifiers have a similar extremity on the off chance that they are joined by the combination "and". In any case, the combination "yet" is utilized for connecting descriptive words with inverse polarities. Moreover such conjunctions as "or", "either-or", "not one or the other nor" are utilized. Some of the time these standards don't appropriate. In this way, creators likewise foresee the polarities of the conjoined descriptive words to check whether the polarities are the equivalent or not, for this reason log-direct regression model is utilized. After forecast arrange, the chart is gotten that gives interfaces between descriptive words. At that point bunching is completed on the chart to isolate descriptors into positive and negative subsets. To finish up, Hazivassiloglou and McKeown had the option to accomplish 90% exactness.

## 3. Decision tree

It is another approach to perform characterization. Decision tree [21] is a classifier that is displayed as various hierarchical decay of data space. The tree structure contains two sorts of nodes: leaf node (contains the estimation of the objective quality, for example positive or negative mark in twofold order assignment) and choice node (contains a condition on one of the properties for space division). The division of the information space is done recursively in hierarchical structure of decision tree.

## 4. Supervised machine learning:

These techniques accept the labeled data that are utilized for the learning procedure. Once training data is pass than obtained output is compared with desired one if class match than new data is use for training otherwise updation of weight were done. As training informational collection, marked reports must be utilized. Typically, bag of words method in [24] is utilized to speak to a report as a feature vector  $d = (w_1, w_2, \dots, w_i, \dots, w_N)$ , where N is set of all the one of a kind terms in the preparation dataset and  $w_i$  is weight of the I-th term. To change over training dataset to an element vector, dictionary with N exceptional words must be made from the input training dataset.

## 5. Unsupervised machine learning methods

Un-supervised learning methodology utilizes unlabeled datasets so as to find the structure and locate the comparable patterns from the information. This technique is generally utilized when a gathering of dependable clarified dataset is not known, yet collection of unlabeled information is simpler. It doesn't cause any troubles when new class information must be recovered. Turney [23] utilizes un-supervised approach for the comments characterization. Comments are ordered into prescribed (approval) and not suggested (disapproval) classes. The author recovers states that comprise of two words dependent on labels designs. The examples are planned so that they need to catch slant phrases. Each expression is a mix of descriptor/qualifier and action word/thing (by and large, 5 examples are proposed). Grammatical feature tagger is utilized to the archive so as to choose which expressions must be recovered.

## V. TEXT PREPROCESSING

As article is group of sections. Passages (paragraph) are accumulation of sentences. While sentences are group of words. So entire preprocessing center around word in the document with no punctuation. So in pre-processing of data there are two regular techniques initially is stop word removal, and second is stem word evacuation [8] .

**Stop List Removals:** As sentence is group with number of words yet a portion of those words are simply use to develop correct sentence in spite of the fact that it doesn't



make any data in the sentence. So reorganization of those words at that point expelling is term as Stop word evacuation. So a vector of words is store by the analyst which help in recognizing of stop words. This evacuation of stop words help in decrease the execution time of the calculation, simultaneously words which not give any productive data is additionally evacuated. Stop words resemble {a, the, for, an, of, and, etc.}. So content archive is change into group of words which is then contrast and these words and afterward each match word is expelled from the document.

In order to understand this assume an sentence {India is a great country in the world} then after pre-processing it become {India, great, country, world} while stop words {is, a, in, the} in the sentence are removed.

**Stem Word Removal** In this technique words which are practically same in prefix can be replace by single word. This can be said accumulation of words share same word is term as stem. So there event in the article make same impact yet, while training in content mining calculation it make difference, so update each word from the group into single word is done in this stem word evacuation pre-handling step. This can be understand by {play, plays, playing} then replace each with word {play}.

## VI. EVALUATION PARAMETER

In order to evaluate results there are many parameter such as accuracy, precision, recall, F-score, etc. Obtaining values can be put in the mention parameter formula to get results.

$$\text{Precision} = \frac{\text{True\_Positive}}{\text{True\_Positive} + \text{False\_Positive}}$$

$$\text{Recall} = \frac{\text{True\_Positive}}{\text{True\_Positive} + \text{False\_Negative}}$$

$$F\_Score = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Accuracy} = \frac{\text{Correct\_Classification}}{\text{Correct\_Classification} + \text{Incorrect\_Classification}}$$

In above true positive value is obtain by the system when the ranked tweet / comment is in favor of user query and system also says that tweet / comment is in favor of the user query. While in case of false positive value it is obtain by the system when the input tweet / comment is in

favour of user query and system do not rank that tweet / comment in their list.

## VII. CONCLUSIONS

This study examines the significance and impacts of content based sentiment examination challenges. Here paper has brief a significant issue of content sentiment / emotion identification from un-organized data. Different methods with their required features were discussed, where neural network gives better opportunity for content classification. Here paper related work of sentiment analysts shows that most of researcher work on two class identification of sentiment either positive or negative, while multiclass need more work for improving the accuracy. Hence in future un-supervised technique should be developed by researcher which should be accurate and less time taken.

## REFERENCES

- [1]. P Reza Zafarani, Mohammad Ali Abbasi, Huan Liu, Social Media Mining – An Introduction, Cambridge University Press, Publisher Location, 2014.
- [2]. Kavanaugh, A. L., Fox, E. A., Sheetz, S. D., Yang, S.Li, L. T., Shoemaker, D. J., ...Xie, L. Social media use by the government: from the routine to the critical. Government Information Quarterly, 29(4), 2012.
- [3]. Twitter Usage Statistics – Internet Live Stats. (n.d.). Retrieved October 22, 2018, from <http://www.internetlivesstats.com/twitter-statistics>.
- [4]. Liu, B., Sentiment analysis and opinion mining. Synthesis lectures on human language technologies, 2012.
- [5]. Dandan Jiang<sup>1</sup>, Xiangfeng Luo<sup>1</sup>, Junyu Xuan, And Zheng Xu “Sentiment Computing for the News Event Based. on the Social Media Big Data”. Digital Object Identifier 10.1109/ACCESS.2016.2607218 IEEE Acss 2017.
- [6]. Priyanka Patil, Pratibha Yalagi, Sentiment Analysis Levels and Techniques: A Survey in International Journal of Innovations in Engineering and technology, April 2016.
- [7]. Ravi, K. and Ravi, V. (2015). A survey on opinion mining and sentiment analysis. Know.-Based Syst., 89(C):14– 46
- [8]. Missen, M. M. S., Boughanem, M., and Cabanac, G. (2012). Opinion mining: reviewed from word to document level. Social Network Analysis and Mining, 3:107–125.
- [9]. Chaturvedi, I., Cambria, E., Welsch, R. E., and Herrera, F. (2018). Distinguishing between facts and opinions for sentiment analysis: Survey and challenges. Information Fusion, 44:65 – 77.
- [10]. Liu, B. (2012). Sentiment Analysis and Opinion Mining. Morgan & Claypool Publishers.

- [11]. Wanxiang Che, Yanyan Zhao, Honglei Guo, Zhong Su, and Ting Liu, "Sentence Compression for Spect-Based Sentiment Analysis" IEEE/ACM transactions on audio, speech, and language processing, vol. 23, no. 12, December 2015
- [12]. Nguyen, D. Q., Nguyen, D. Q., Vu, T., & Pham, S. B. (2014). Sentiment classification on polarity reviews: an empirical study using rating-based features. In Proceedings of the 5th workshop on computational approaches to subjectivity, sentiment and social media analysis (pp. 128-135).
- [13]. Tiwari, P., Mishra, B. K., Kumar, S., & Kumar, V. (2017). Implementation of n-gram methodology for rotten tomatoes reviews dataset sentiment analysis. International Journal of Knowledge Discovery in Bioinformatics (IJKDB), 7 (1), 30-41.
- [14]. Arun, K., Srinagesh, A., & Ramesh, M. (2017). Twitter Sentiment Analysis on Demonetization tweets in India Using R language. International Journal of Computer Engineering in Research Trends, 4 (6), 252-258.
- [15]. Anta, A. F., Chiroque, L. N., Morere, P., & Santos, A. (2013). Sentiment analysis and topic detection of Spanish tweets: A comparative study of NLP techniques. Procesamiento del lenguaje natural, 50, 45-52.
- [16]. Basha, S. M., Zhenning, Y., Rajput, D. S., Iyengar, N., & Caytiles, D. R. (2017). Weighted Fuzzy Rule Based Sentiment Prediction Analysis on Tweets. International Journal of Grid and Distributed Computing, 10 (6), 41-54.
- [17]. Hailong, Z., Wenyan, G., & Bo, J. (2014, September). Machine learning and lexicon based methods for sentiment classification: A survey. In Web Information System and Application Conference (WISA), 2014 11th (pp. 262-265)
- [18]. Hu, M., & Liu, B. (2004, August). Mining and summarizing customer reviews. In Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 168-177).
- [19]. Liu, B. (2012). Sentiment analysis and opinion mining. Synthesis lectures on human language technologies, 5(1), 1-167.
- [20]. Hatzivassiloglou, V., & McKeown, K. R. (1997, July). Predicting the semantic orientation of adjectives. In Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics (pp. 174-181)
- [21]. Aggarwal, C. C., & Zhai, C. (2012). A survey of text classification algorithms. In mining text data (pp. 163-222). Springer US
- [22]. Tang, B., Kay, S., & He, H. (2016). Toward optimal feature selection in naive Bayes for text categorization. IEEE Transactions on Knowledge and Data Engineering, 28(9), 2508-2521
- [23]. Turney, P. D. (2002, July). Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews. In Proceedings of the 40th annual meeting on association for computational linguistics (pp. 417-424).