# Comparative Study of Available Technique for Detection in Sentiment Analysis

## Miss. Siddhi Patni[1] ,Prof. Avinash Wadhe[2]

[1,]*M.E (CSE) 2nd Semester G.H.Raisoni College of Engineering&Management Amravati, India*
[2,] *ME (CSE) G.H.Raisoni College of Engineering   & Management Amravati, India*

### ABSTRACT:
*Our day-to-day life has always been influenced by what people think. Ideas and opinions of others have always affected our own opinions. As the Web plays an increasingly significant role in people's social lives, it contains more and more information concerning their opinions and sentiments. The distillation of knowledge from this huge amount of unstructured information, also known as opinion mining and sentiment analysis. It has recently raised growing interest for purposes such as customer service, financial market prediction, public security monitoring, election investigation, health related quality of life measure, etc.  Sentiment Analyzer  (SA) that extracts sentiment (or opinion) about a subject from online text documents. Instead of classifying the sentiment of an entire document about a subject, SA detects all references to the given subject, and determines sentiment in each of the references using natural language processing (NLP) techniques. There are various machine learning algorithms that attempt to predict the sentiment or opinions of some documents or information of particular data and organize data, such as finding positive and negative reviews while diminishing the need for human effort to classify the information. This paper compares the NLP and machine learning methods of sentiment analysis and determines which one is better.*

*Keywords: Sentiment Analyzer, Natural Language Processing, Machine Learning.*

## I.      INTRODUCTION

A vital part of the information era has been to find out the opinions of other people. In the pre-web era, it was customary for an individual to ask his or her friends and relatives for opinions before making a decision. Organizations conducted opinion polls, surveys to understand the sentiment and opinion of the general public towards its products or services. In the past few years, web documents are receiving great attention as a new medium that describes individual experiences and opinions. With explosion of Web 2.0 [1] applications such as blogs, forums and social networks came. The rise of blogs and social networks has fueled a market in personal opinion, reviews, ratings, recommendations and other forms of online expression. For computer scientists, this fast-growing mountain of data is opening a tantalizing window onto the collective consciousness of Internet users. An emerging field known as sentiment analysis is taking shape around one of the computer world's unexplored frontiers, translating the vagaries of human emotion into hard data. This is more than just an interesting programming exercise. For many businesses, online opinion has turned into a kind of virtual currency that can make or break a product in the marketplace. Therefore organizations have evolved and now look at review sites to know the public opinion about their products instead of conducting surveys. However, gathering all this online information manually is time consuming. Therefore automatic sentiment analysis is important. To do so, the main task is to extract the opinions, facts and sentiments expressed in these reviews. Sentiment analyzer that extracts sentiment  about a given topic using NLP techniques consists of 1) a topic specific feature term extraction, 2) sentiment extraction, and 3) (subject, sentiment) association by relationship analysis. SA utilizes two linguistic resources for the analysis: the sentiment lexicon and the sentiment pattern database.

Using machine learning algorithm for sentiment analysis is to extract human emotions from text or documents. Metrics such as accuracy of prediction and precision/recall are presented to gauge the success of these different algorithms.  A system is there to process the documents and to predict human reactions, as well as provide results.

## II.      LITERATURE SURVEY

Sentiment Analysis is the computational study of opinions, sentiments and emotions expressed in text. Lui [2] mathematically represented an opinion as a quintuple (o, f, so, h, t), where o is an object; f is a feature of the object o; so is the polarity of the opinion on feature of object o; h is an opinion holder; t is the time when the opinion is expressed. The goal of sentiment analysis is to detect subjective information contained in various sources and determine the mind-set of an author towards a text or document. The research field of sentiment analysis is rapid in progress due to the rich and diverse data provided by Web 2.0 applications. Blogs, review sites, forums, micro blogging sites, wikis and social networks are all the sources used for sentiment analysis. Today, a huge amount of information is available   from this application of Web 2.0, among these types of information available, one useful type is the sentiment, or opinions people express towards a subject that *is* either a topic of interest or a feature of the topic .There has been extensive research on automatic text analysis for sentiment, such as sentiment classifiers [3], affect analysis, automatic survey analysis, opinion extraction [4], or recommender systems. These methods typically try to extract the overall sentiment revealed in a document, positive, negative or neutral .Two challenging aspects of sentiment analysis are: First, although the overall opinion about a topic is useful, it is only a part of the information of interest. Document level sentiment classification fails to detect sentiment about individual aspects of the topic. In reality, for example, though one could be generally happy about his car, he might be dissatisfied by the engine noise. To the manufacturers, these individual weaknesses and strengths are equally important to know, or even more valuable than the overall satisfaction level of customers. Second, the association of the extracted sentiment to a specific topic is difficult. Most statistical opinion extraction algorithms show    some shortcomings and hence developed sentiment Analyzer (SA) that extracts topic-specific features, extracts sentiment of each sentiment-bearing phrase, makes (topic/ feature), sentiment association.

The machine learning approach applicable to sentiment analysis mostly belongs to supervised classification in general and text classification techniques in particular. Thus, it is called "supervised learning". In a machine learning based classification, two sets of documents are required: training and a test set. A training set is used by an automatic classifier to learn the differentiating characteristics of documents, and a test set is used to validate the performance of the automatic classifier. A number of machine learning techniques have been adopted to classify the reviews like Naive Bayes (NB), maximum entropy (ME), and support vector machines (SVM) have achieved great success in text categorization.

## III. SENTIMENT ANALYSIS METHODS

### A.  *Natural Language Processing Approach*

In the natural language processing method of sentiment analysis this paper extract the opinion or sentiment by using sentiment analyzer [5] that includes feature term extraction a feature term of a topic is a term that satisfies one of the following relationships:

a)   a part-of relationship with the given topic

b)   an attribute-of relationship with the given topic.

c)   an attribute-of relationship with a known feature of the given topic.

There are two   linguistic resources used by sentiment analysis sentiment lexicon and sentiment pattern database.   The sentiment lexicon contains the sentiment definition of individual words in the following form:
<lexical_entry> <POS> <sent_category>
Lexical_entry is a (possibly multi-word) term that has sentimental connotation, POS is the required POS tag of lexical entry, sentiment_category: + / -
The following is an example of the lexicon entry:
"Excellent" JJ +.
Sentiment pattern database contains sentiment extraction patterns for sentence predicates. The database entry is defined in the following form:
<predicate>  <sent_category>  <target>
 predicate: typically a verb, sent_category: + / - / [˜] source is a sentence component (SP/OP/CP/PP) whose sentiment is transferred to the target. SP, OP, CP, and PP represent subject, object, complement (or adjective), and prepositional phrases, respectively. The opposite sentiment polarity of source is assigned to the target, if ˜ is specified in front of source target is a sentence component (SP/OP/PP) the sentiment is directed to.  As a preprocessing step to sentiment analysis, we extract sentences from input documents containing mentions of subject terms of interest.  After parsing each input sentence by a syntactic parser, SA identifies sentiment phrases from subject, object, adjective, and prepositional phrases of the sentence. Within the phrase, we identify

all sentiment adjectives defined in the sentiment lexicon. For example, vibrant is positive sentiment phrase for the sentence

"The colors are vibrant."

Extract all base noun consist of at least one sentiment word. The sentiment of the phrase is determined by the sentiment words in the phrase. For example, excellent pictures are a positive sentiment phrase because excellent is a positive sentiment word. For a sentiment phrase with a word with negative meaning, such as not, no, never, hardly, seldom, or little, the polarity of the sentiment is reversed. SA extracts T- and B-expressions in order to make (subject, sentiment) association. From a T-expression, sentiment of the verb (for sentiment verbs) or source (for trans verb), and from a B-expression, sentiment of the adjective, is assigned to the target.

### B. Machine Learning approach

For machine learning there is a system [6] consisted of first processing the confessions in order to extract a feature set, before passing the data into a supervised learning algorithm.

**1) Parser Method:** In order to refine our data and improve the feature set, we removed all HTML tags using a Python parser. This was essential towards refining our dataset because HTML tags do not convey emotions and would skew our feature vector by including phrases that have no semantic meaning (e.g. '&nbsp ;'). Emoticons, on the other hand, are an excellent way of conveying emotions through text because it captures the emotion of the writer by including a facial expression. Therefore; we captured this unique feature set and used it to improve our feature vector.

**2) Spell Checking:** There are many spelling errors. In order to reduce problems of over fitting as a result of having too many unique spellings, rather than raw data through a spell checker and corrected all the spelling errors.

**3) Features:** In this paper, there are three features considered: bag of words, WordNet2 synsets, and sentiment lexicons.

3.1)   *Bag of Words (BoW):*  It treats each unique word token as a separate feature.

3.2)   *Word Net Synsets:* In order to further improve the quality of the feature set and decrease over fitting, we used WordNet to map the words in the confessions onto their synonym set (synset).

3.3)   *Sentiment Lexicons:* Sentiment lexicons are groupings of words into emotion.



Fig. 1 Model  Diagram.

**4) TF-IDF:** Function words such as 'and', 'the', 'he', 'she' occur very often across all confessions. Therefore, it makes little sense to put a lot of weight on such words when using bag of words to classify the documents. One common approach is to remove all words found in a list of high frequency stop words. A better approach is to consider each word's Term Frequency-Inverse Document Frequency (TF-IDF) weight. The intuition is that a frequent word that appears in only a few confessions conveys a lot of information, while an infrequent word that appears in many confessions conveys very little information.

## IV.   THEORETICAL ANALYSIS

From our initial experience with sentiment detection, we have identified a few areas of potentially substantial improvements. We expect full parsing will provide better sentence structure analysis, thus better relationship analysis. Second, more advanced sentiment patterns currently require a fair amount of manual validation. Although some amount of human expert involvement may be in-evitable in the validation to handle the semantics accurately, we plan on more research on increasing the level of automation systems understand and manipulate natural languages to perform the desired tasks.

At the core of any NLP task there is the important issue of natural language understanding. The process of building computer programs that understand natural language involves three major problems: the first one relates to the thought process, the second one to the representation and meaning of the linguistic input, and the third one to the world knowledge. Thus, an NLP system may begin at the word level  to determine the morphological structure, nature (such as part-of-speech, meaning) etc. of the word – and then may move on to the sentence level – to determine the word order, grammar, meaning of the entire sentence, etc. Then to the context and the overall  environment or domain. A given word or a sentence may have a specific meaning or connotation in a given context or domain, and may be related to many other words and/or sentences in the given context.

In machine learning, the broad field of Artificial Intelligence, which aims to mimic intelligent abilities of humans by machines. In the field of Machine Learning [6] one considers the important question of how to make machines able to "learn". Learning in this context is understood as inductive inference, where one observes examples that represent incomplete information about some "statistical phenomenon". In unsupervised learning one typically tries to uncover hidden regularities (e.g. clusters) or to detect anomalies in the data (for instance some unusual machine function or a network intrusion). In supervised learning, there is a label associated with each example. It is supposed to be the answer to a question about the example. Based on these examples (including the labels), one is particularly interested to predict the answer for other cases before they are explicitly observed. Hence, learning is not only a question of remembering but also of generalization to unseen cases.

## V.   COMPARATIVE STUDY

In this paper there is a comparative study between NLP and ML approach. There are several parameters [7] considered:

### A.  *Keyword Selection*

Topic based classification usually uses a set of keywords to classify texts in different classes. In sentiment analysis we have to classify the text in to two classes (positive and negative) which are so different from each other. But coming up with a right set of keyword is not a petty task. This is because sentiment can often be expressed in a delicate manner making it tricky to be identified when a term in a sentence or document is considered in isolation.

### B.   *Sentiment is Domain Specific*

Sentiment is domain specific and the meaning of words changes depending on the context they are used in.  Consider an example:  go read the book
The example has a positive sentiment in the book domain but a negative sentiment in the movie domain, it suggests that the book is preferred over the movie, and thus have an opposite result [8].

### C.  *Multiple Opinions in a Sentence*

Single sentence can contain multiple opinions along with subjective and factual portions. It is helpful to isolate such clauses. It is also important to estimate the strength of opinions in these clauses so that we can find the overall sentiment in the sentence, e.g.: The picture quality of this camera is amazing and so is the battery life, but the viewfinder is too small for such a great camera‖. It expresses both positive and negative opinions [8].

### D.  *Negation Handling*

Handling negation can be tricky in sentiment analysis. For example:
I like this dress‖ and I don't like this dress‖
These sentences differ from each other by only one token but consequently are to be assigned to different and opposite classes. Negation words are called polarity reversers.

| Parameter of Sentiment Analysis | NLP Approach | ML Approach |
|---|---|---|
| Keyword Selection | Not Efficient. | Most Effective. |
| Sentiment is Domain Specific | Efficient to check grammar in specific sentiment statement. | Not more effective and require specific training to statement. |
| Multiple Opinions in a Sentence | Not efficient for frequently changing opinion. | Efficient for differ opinion statement by using ML agent. |
| Negation Handling | More efficient for negation statement. | Equally Efficient for negation statement. |

**Table I Comparative Study Of Nlp And Ml Approach.**

In this paper , two different approaches are considered and compared with the help of different parameters so from this table it can be noticed that NLP approach is much efficient in keyword selection **,** efficient to check grammar in specific sentiment statement , not efficient for frequently changing opinion , more efficient for negation statement and for ML approach it has noticed that it is more efficient in keyword selection, not more effective in domain specific sentiment and require specific training to statement, efficient for differ opinion statement by using ML agent and equally efficient for negation statement. So therefore if there is combination of two approaches then analysis of sentiment will be more effective.

## VI. CONCLUSION

Sentiment Analyzer (SA) consistently demonstrated high quality results of for the general web pages. Although some amount of human expert involvement may be inevitable in the validation to handle the semantics accurately, plan on more research on increasing the level of automation. Nonetheless, the synset and sentiment lexicons, used are better suited to more formal styles of writing. An alternative approach is to replace our synsets and lexicons with "slang" versions or even the automatic generation of sentiment lexicons on a slang corpus. Another area of interest is the difficulty in correlating topics with sentiment. Intuition says that topics themselves should portray different sentiments, and so should be useful for sentiment analysis. This method turns out to be fairly crude, as sometimes topics may be too neutral or too general. Thus, it is concluded that hybrid approach that is combination of NLP and ML approach can strengthen analysis of sentiment or opinions on different parameters and can give a better result than applying individual approach .

## REFERENCES

[1] Tim O'Reilly, Web 2.0 Compact Definition: Trying Again (O'Reilly Media, Sebastopol), http://radar.oreilly.com/archives/2006/12/web_20_compact.html. Accessed 22 Mar 2007
[2] Liu B. Sentiment Analysis and Subjectivity. Handbook of Natural Language Processing, Second edition, 2010
[3] B. Pang, L. Lee, and S. Vaithyanathan. Thumbs up? Sentiment classification using machine learning techniques. *In Proc. of the 2002 ACL EMNLP Conf., pages 79–86, 2002.*
[4] S. Morinaga, K. Yamanishi, K. Teteishi, and T. Fukushima. Mining product reputations on the web. *In Proc. of the 8th ACM SIGKDD Conf., 2002.*
[5] Jeonghee Yi, Tetsuya Nasukawa, Razvan Bunescu, Wayne Niblack*,"Sentiment Analyzer: Extracting Sentiments about a Given Topic using Natural Language Processing Techniques",* Proceedings of the Third IEEE International Conference on Data Mining ,2003 .
[6] Raymond Hsu, Bozhi See, Alan Wu,*''Machine Learning for Sentiment Analysis on the Experience Project",* 2010.
[7] Akshi Kumar ,Teeja Mary Sebastian, Sentiment Analysis: "*A Perspective on its Past, Present and Future",*2012 .
[8] Pang, B and Lee L. Opinion mining and sentiment analysis. Foundations and Trends in Information Retrieva, l 2008, (1-2), 1–135

**Author's Profile:**
**Miss. Siddhi S. Patni** is doing M.E (CSE) from G.H Raisoni College of Engineering and Management, Amravati and has done B.E in Information Technology from SGBAU, Amravati.

**Prof. Avinash P. Wadhe:** Received the B.E and from SGBAU Amravati university and M-Tech (CSE) From G.H Raisoni College of Engineering, Nagpur (an Autonomous Institute). He currently an Assistant Professor with the G.H Raisoni College of Engineering and Management, Amravati SGBAU Amravati University. His research interest include Network Security, Data mining and Fuzzy system .He has contributmore than 20 research paper. He had awarded with young investigator award in international conference.