RESEARCH ARTICLE                                                                OPEN ACCESS

# Sentiment Analysis Using Hybrid Approach: A Survey

## Chauhan Ashish P, Dr. K. M. Patel
Computer Engineering School of Engineering, RK University Gujarat, India
Computer Engineering School of Engineering, RK University Gujarat, India

*Abstract*
Sentiment analysis is the process of identifying people's attitude and emotional state's from language. The main objective is realized by identifying a set of potential features in the review and extracting opinion expressions about those features by exploiting their associations. Opinion mining, also known as Sentiment analysis, plays an important role in this process. It is the study of emotions i.e. Sentiments, Expressions that are stated in natural language. Natural language techniques are applied to extract emotions from unstructured data. There are several techniques which can be used to analysis such type of data. Here, we are categorizing these techniques broadly as "supervised learning", "unsupervised learning" and "hybrid techniques". The objective of this paper is to provide the overview of Sentiment Analysis, their challenges and a comparative analysis of it's techniques in the field of Natural Language Processing.
*Keywords*—Sentiment Analysis, opining mining, SVM (Support Vector Machine), Rainforest, Hybrid Approach

## I. INTRODUCTION

Sentiment analysis aims to uncover the attitude of the author on a particular topic from the written text. Other terms used to denote this research area include opinion mining and subjectivity detection. Main aim of Sentiment analysis is to minimize the gap between the human beings and machine. Because sentiment analysis proposes the feature which provides an approx accurate opinion on any object or person based on different techniques and methods. So it is the collection of human and electronic intelligences for gaining the opinion on text mining. The user generated content is available in various forms such as web logs, reviews, news, discussion forums. Web 2.0 3.0 has provided a platform to share the feelings and views about the products and services. Basically, this problem can be explained better using the following user review about a cell phone.I bought an sony Phone a few days ago. it was really cool and light weight.the speaker quality was also really clear and sound.but the battery life was too weak but that is fine for me . but my cousin was angry with me as i did not tell her before i brought this sony phone.the phone is too expensive. my cousin wanted me to exchange this Sony phone The above text contains total seven sentences.It is Contains both kind of sentiments; i.e. It contains two issues, quality of product (a positive sentiment is attached with it) and cost issues (negative sentiment is attached with this issue), so decision of more important sentiment is also a problem. There are several challenges in Sentiment analysis. The first is an opinion word that is considered to be positive in one situation may be considered negative in another situation. A second challenge is that people don't always express opinions in a same way. The third one is mixed Review in that we got mixes review from the user like this phone is good but price is so costly thats kind of review called Multi-theme documents. In such type of document problem statement does not always remain so clear. This paper is organized as follows: Section- II gives literature survey, section-III gives describes analysis of some techniques proposed in literature and section-IV will conclude the paper.

## II. DATA SOURCES

User's opinion is a major criterion for the improvement of the Sentiment analyis. Blogs, review sites, data and micro blogs provide a good understanding of the reception level of the products and services.

- **Blogs**
With an increasing usage of the internet, blogging and blog pages are growing rapidly. Blog pages have become the most popular means to express one's personal opinions.

- **Review sites**
With an increasing usage of the internet, blogging and blog pages are growing rapidly. Blog pages have become the most popular means to express one's personal opinions.The reviews for products or services are usually based on opinions expressed in much unstructured format. The reviewers data used in most of the sentiment classification studies are collected from the e-commerce websites like www.amazon.com (product

reviews), www.yelp.com (restaurant reviews), www.CNET download.com (product reviews) and www.reviewcentre.com .

- **DataSet**

Most of the work in the field uses movie reviews data for classification. Movie review datas are available as dataset (http:// www.cs.cornell.edu/People/pabo/movie-review-data). Other dataset which is available online is multi-domain sentiment (MDS) dataset. (http:// www.cs.jhu.edu/mdredze/datasets/sentiment). The MDS dataset contains four different types of product reviews extracted from Amazon.com including Books, DVDs, Electronics and Kitchen appliances, with 1000 positive and 1000 negative reviews for each domain. other review dataset is available at : http://www.simafore.com/download-the-dataset-and- market-basket-analysis-process.

- **Twitter Dataset**

    Twitter is a social news website .The Twitter Senti- ment Analysis Dataset contains classified tweets, each row is marked as 1 for positive sentiment and 0 for negative sentiment.and twitter dataset is available at
:http://socialcomputing.asu.edu/datasets/Twitter

- **Movie Review Dataset**

    This is a dataset for binary sentiment classifica- tion containing substantially more data than previ- ous benchmark datasets. it provide a set of 25,000 highly polar movie reviews for training, and 25,000 for testing. There is additional unlabeled data for use as well. Raw text and already processed bag of words formats are provided. Movie dataset is available at: http:// ai.stanford.edu/amaas/data/sentiment/

### III. LITERATURE SURVEY

    Mainly there are three categories of learning algorithms.
    This impact category can be divided into several sub- categories.

#### A. *Supervised Learning Algorithms*

    Some of the most predominant Supervised Learning tech- niques in Sentiment Analysis have been SVM, Nave Bayesian Classifiers and other Decision Trees.

- Nave Bayes

    A Nave Bayes classifier is a simple probabilistic model based on the Bayes rule along with a strong independence assumption [2]. The Nave Bayes model involves a simplifying conditional independence as- sumption. That is a given class (where people don't express opinions

in the same way; they use opinion words as positive or negative comments), the words are conditionally independent of each other. This assumption does not affect much the accuracy in text classification but makes really fast classification algorithms applicable for the problems.

- Maximum Entropy

    It is conditional exponential classifier. It maps each pair of feature set and its label to a vector[3]. It is also called as loglinear classifier because they work by extracting some set of features from the input, combining them linearly (each feature is multiplied by its weight and added up) and using this sum as exponent. It is parameterized by a set of weights that are used to combine the joint-features that are generated from a set of features by an encoding.

- Decision Tree

    It is a tree in which internal nodes are represented by features, edges represent tests to be done at feature weights and leaf nodes represent categories which results from above tests[5]. It categorizes a document by starting at the tree root and moving successfully downward via the branches (whose conditions are satisfied by the document) until a leaf node is reached. The document is then classified in the category that labels the leaf node. Decision Trees have been used in many applications in speech and language processing.

- Rainforest

    It is one of the most attractive classification models Developing a unifying framework that can be applied to most decision tree algorithms.Its based on Top- down approach .Its use for Fast Decision TreeCon- struction of Large Datasets AVCgroup Constructin.

- Support Vector Machines

    In comparisons, SVM has outperformed other classi- fiers such as Nave Bayes. While SVM has become a dominant technique for text classification, other algorithms such as Winnow and AdaBoost have also been used in previous sentiment classification studies. SVM gives highest accuracy results in text classifi- cation problems. SVM represents example as points in a space which are mapped to a high dimensional space where mapped examples of separate classes are divided by as wide as possible tangential possible distance to the hyper plane. New examples are mapped into this same space and depending on which side of the hyper plane they are positioned, they are predicted to belong to a certain class. SVM hyper planes are

fully determined by a relatively small subset of the training instances, which are called support vectors. The rest of the training data have no influence on the trained classifier. SVMs have been employed success- fully in text classification and in a variety of sequence processing applications[6].

### B. Unsupervised Learning Algorithms

These are also known as lexicon based techniques. This involve learning patterns in the input when no specific output values are supplied, this means that the learner only receives an unlabelled set of examples.Unsupervised methods can also be used to label a corpus that can later be used for supervised learning. An agent purely based on unsupervised Sentiment analysis approaches Supervised learning Unsupervised learn- ing

- k-mean

K-Means tries to find the natural clusters in the data, by calculating the distance from the centers of the clusters[7]. The position of centers is iteratively changed until the distances between all the points are minimal. The centers are initially randomly assigned. K-Means can find only local maximum, and the final label assignment can be suboptimal. Common practice is to repeat the algorithm on the same data multiple times, and to report the best result. We have repeated the procedure 10 times in our experiments. We have used Euclidean distance as dissimilarity metric be- tween feature vectors. We use B3 measure to evaluate the performance of the classifiers.

### C. Hybrid Techniques

In Hybrid Techniques both combination of machine learn- ing and lexicon base approaches ate used. Researchers have proved that this combination gives improved performance of classification [8]. Mudinas et al. proposed a concept-level sentiment analysis system, called pSenti, which is developed by combining lexicon based and learning-based approaches. The main advantage of their hybrid approach using a lexi-con/learning symbiosis is to obtain the best of both worlds- stability as well as readability from a carefully designed lexicon, and the high accuracy from a powerful supervised learning algorithm. Their system uses a lexicon from public resources for initial sentiment detection. They used sentiment words as features in machine learning method. The weight of such a feature is the sum of the sentiment value in the given review. For those adjectives which are not in sentiment lexicon, their occurring frequencies are used as their initial values. Their hybrid approach pSenti achieved 82.30

- Method one

In this method we have use the Stock Sentiment Analysis (SSA)hybrid system which integrates 3 com- plementing modules: a Dictionary Based SA,a Pattern Based SA, and a Semantic Event Based SA [10].

- method two

In this method hybrid classification method is pro- posed based on coupling classification methods using arcing classifier and their performances are analyzed in terms of accuracy. A Classifier ensemble was de- signed using Naive Bayes (NB), Genetic Algorithm (GA)[12].

## IV. CONCLUSION

Thus, I examined that all of the above algorithms and techniques, which is used in sentiment analysis, are not ac- curate 100 percent. All of these have their own merits and de-merits. But yes, Supervised machine learning techniques have shown more better performance in compare to other once i.e. unsupervised machine learning techniques. However, the unsupervised methods is important too because supervised methods demand large amounts of labeled training data that are very expensive. But Support Vector Machines (SVM) has high accuracy than other algorithms. In most of case SVM gives higher efficiency and accuracy. but here I have examined that The accuracy of hybrid classifier methods can be increased by increasing the number of classifiers.In future we need to enhance the classification accuracy.

## V. ACKNOWLEDGMENT

## REFERENCES
[1] Huosong Xia ,Min Tao, Yi Wang,"*Sentiment Text Classification of Cus- tomers Reviews on the Web Based on SVM*" Department of economics and management, Business Intelligence and Data Mining Lab Wuhan University of Science and Engineering Wuhan , China ,2010
[2] Bing Liu,"*Sentiment Analysis and Opinion Mining Morgan Claypool Publishers*, May 2012.
[3] Kechaou, Z ,Ben Ammar, M. Alimi "*Improving e-learning with sentiment analysis of users' opinions,*" Global Engineering Education Conference (EDUCON), 2011 IEEE , vol., no., pp.1032- 1038, 4-6 April 2011
[4] Sowmya Kamath S, Anusha Bagalkotkar, Ashesh Khandelwal, Shivam Pandey,

Kumari Poornima, *Sentiment Analysis Based Approaches for Understanding User Context in Web Content*, 978- 0-7695-4958-3/13, 2013 IEEE.

[5] V.K. Singh, R. Piryani, A. Uddin,"*Sentiment Analysis of Textual Reviews*" Department of Computer Science South Asian University New Delhi, India 2013

[6] Michelle Annett and Grzegorz Kondrak , "*A Comparison of Sentiment Analysis Techniques: Polarizing Movie Blogs*"

[7] B. Liu Web Data Mining: Exploring hyperlinks, contents, and usage data," Opinion Mining. Springer, 2007

[8] A. Mudinas, D. Zhang, M. Levene, *Combining lexicon and learning based approaches for conceptlevel sentiment analysis*, Proceedings of the First International Workshop on Issues of Sentiment Discovery and Opinion Mining, ACM, New York, NY, USA, Article 5, pp. 1-8, 2012.

[9] Kurt Junshean Espinosa, Kevin Llaguno, Jaime Caro,"*Sentiment analysis of Facebook statuses using Naive Bayes classifier for language learning*" Department of Computer Science University of the Philippines Cebu city, Philippines 2010

[10] Ronen Feldman ,Benjamin Rozenfeld, Micha Y. Breakstone,"*SSA A Hybrid Approach to Sentiment Analysis of Stocks*" School of Business Administration The Hebrew University of Jerusalem Jerusalem, ISRAEL

[11] Amandeep Kaur,Vishal Gupta,"*A Survey on Sentiment Analysis and Opinion Mining Techniquesg*" University Institute of Engineering and Technology, Chandigarh, India 2013

[12] M.Govindarajan,"*Sentiment Analysis of Movie Reviews using Hybrid Method of Naive Bayes and Genetic Algorithm*" Department of Com- puter Science and Engineering, Annamalai University, Annamalai Nagar, Tamil Nadu,India 2013

[13] K. Revathy,Dr. B. Sathiyabhama,"*A Hybrid Approach for Supervised Twitter Sentiment Classification*" Department of Computer Science and Engineering, Sona College of Technology, Salem, India 2013

[14] Balaji Jagtap, Virendrakumar Dhotre,"*SVM and HMM Based Hybrid Approach of Sentiment Analysis for Teacher Feedback Assessment*" Dept. of Information Technology, Maharashtra Institute of Technology, Pune ,india 2014

[15] Basant Agarwal,Namita Mittal,"*A Hybrid Approach for Supervised Twitter Sentiment Classification*" Department of Computer Engineering Malaviya National Institute Technology Jaipur, India 2014

[16] Malhar Anjaria, Ram Mahana Reddy Guddeti,"*A Hybrid Approach for Supervised Twitter Sentiment Classification*" Department of Information Technology National Institute of Technology Karnataka, Surathkal, Man- galore ,india 2014

[17] Neethu M S,Rajasree R,"*Sentiment Analysis in Twitter using Machine Learning Techniques*" Department of Computer Science and Engineering College of Engineering Trivandrum, India 2013

[18] Khin Phyu Phyu Shein, Thi Thi Soe Nyunt,"*Sentiment Classification based on Ontology and SVM Classifier*" Computer Software Technology Department University of Computer Studies, Yangon. Yangon, Myanmar. 2010

[19] Khin Phyu Phyu Shein, Thi Thi Soe Nyunt,Kurt Junshean Espinosa, Kevin Llaguno, Jaime Caro,"*Sentiment analysis of Facebook statuses using Naive Bayes classifier for language learning*" Department of Informatics University of Piraeus Piraeus, Greece 2010

[20] Wang Zuhui ,Jiang Wei,"*Online Reviews Sentiment Analysis Applying Mutual Information*" Research center of Information Management and Information System Harbin institute of technology Harbin, China 2012

[21] Dae-Ki Kang,"*Effective Sentiment Analysis based on Term Evaluation by Bayesian Model Selection Criteria*" Division of Computer and Infor- mation Engineering, Dongseo University, Busan, South Korea 2013

[22] Dae-Ki Kang,"*Effective Sentiment Analysis based on Term Evaluation by Bayesian Model Selection Criteria*" Division of Computer and Infor- mation Engineering, Dongseo University, Busan, South Korea 2013

[23] Zied Kechaou, Mohamed Ben Ammar, Adel. M Alimi, "*Improving e- learning with sentiment analysis of users opinions*" REsearch Group on Intelligent Machines, University of Sfax, National Engineering School of Sfax (ENIS), BP 1173, Sfax, 3038, Tunisia 2011

[24] Lei Shi, Bai Sun, Liang Kong and Yan Zhangz,"Web Forum Sentiment Analysis based on Topics" Department of Machine Intelligence Peking University Beijing 100871, China 2009

[25] Pranav Waila, V.K. Singh , M. K. Singh,"Blog Text Analysis Using Topic Modeling, Named Entity Recognition and Sentiment Classifier Combine" Department of Computer Science, South Asian University,
New Delhi, India 2013

[26] Wei Wei , Jon Atle Gulla,"Sentiment Analysis In a Hybrid Hierarchical Classification Process" Department of Computer and Information Science Norwegian University of Science and Technology Sem Sffilands vei, NO- 7491 Trondheim, Norway 2012