

# Survey on Sentiment Analysis

Suraj D. M., Rohan A. R.,  
Varun A. Prasad, Shirsu Mitra  
Department of Computer Science & Engineering  
B N M Institute of Technology  
Bangalore, India

Dr. Vimuktha Evangeleen Salis  
Associate Professor, Department of Computer Science &  
Engineering  
B N M Institute of Technology  
Bangalore, India

**Abstract**—Emotions play a major role in a person's life as they affect their decision-making process in various kinds of events they face. Computers can be built which can be used to detect emotions but they are often limited to factual information. There are multiple blogs, discussion forums and social network platforms where many users post a large number of informal messages every day. Researchers find emotion detection from text very enthralling because of its expansive range of applications, such as social support, evaluating welfare of a community and even detection and addressing of suicidal tendencies. Algorithms or methods need to be implemented in sentiment analysis to detect an emotion in an informal context. These methods can also be used to identify anomalous or inappropriate utterances that could be construed as threatening to others or self. Nevertheless, existing emotion detection methods or algorithms are used for commercial purposes, developed to evaluate reviews about a product rather than a person's behavior. This survey describes various theories of sentiment analysis that provide a detailed explanation of emotion models.

**Keywords**—Sentiment analysis; emotion detection; human feelings; sentiment; emotions

## I. INTRODUCTION

Sentiment analysis can be defined as a process of detecting or identifying emotion using text analysis, natural language processing and linguistics. The main goal is to observe the attitude of the person and extract his current emotional status. Sentiment detection and analysis can be implemented using unsupervised and supervised learning methods such as support vector machines, artificial neural networks, etc. The result of sentiment analysis is either a positive, negative or neutral opinion of the user based on that subject or topic.

Fig. 1. displays the normal flow of the process involved in sentiment analysis. It involves training a machine learning model using a relevant dataset and then passing inputs to the trained model for sentiment classification. There are dictionary based techniques as well that do not involve training of the model.

Sentiment analysis has many applications in various kinds of fields such as psychology, behavior science and neuroscience as they are an important feature of human behavior. This has attracted the attention of many computer researchers because of its wide range of applications such as social support, evaluating welfare of a community and even prevention of suicidal thoughts. Models can also be used to analyze data from social media, which allows us to gain insights from the general public about a certain product or a topic. Sentiment analysis can also produce interesting results

in various suicide prevention and e-learning platforms. Enthusiastic about its enormous potential, we decided to perform a survey on current systems that detect emotion from text and make it available to the research community.

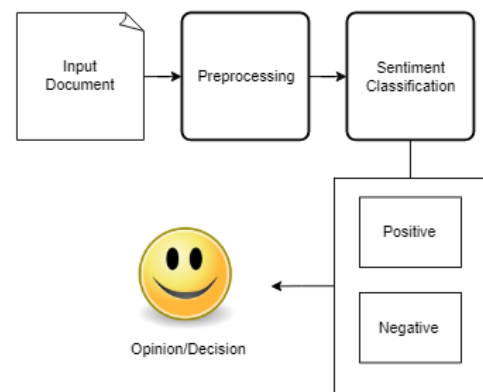


Fig. 1. Sentiment Analysis Flow

This survey depicts recent works in the field of sentiment analysis over textual data. Sentiment analysis is a powerful tool in the field of affective computing whose goal is to enable machines to identify and express emotions. Current implementations of sentiment analysis deal with a person's modality such as a facial expression, text or voice. This survey mainly focuses on reviewing models or approaches that were implemented for detecting emotion in textual data. A statistical comparative analysis was not possible since each work evaluated their systems using different data sets.

## II. COMPUTATIONAL APPROACHES

### A. Machine Learning

Machine learning approaches rely on machine learning algorithms to perform sentiment analysis as a standard classification issue that makes use of syntactic as well as linguistic features. These algorithms use a set of methods to analyze and detect different patterns in the provided dataset to train the models. The trained models can then be used to identify unexplored patterns and classify them based on its emotion. Machine learning algorithms are implemented in two ways i.e. unsupervised learning and supervised learning. Supervised learning is generally performed when the output can be predicted with the given input and unsupervised learning is performed when the output is unpredictable. Various machine learning algorithms that fall under supervised learning algorithms like naive bayes, decision tree learning, etc. or unsupervised learning techniques like neural

networks, cluster analysis methods, etc. can be used to implement sentiment analysis.

*B. Dictionary based approach*

One of the easiest ways to implement sentiment analysis is to use dictionary oriented approaches. WordNet and SentiWordNet are some of the widely used public platforms for sentiment analysis. Dictionaries are usually built with a set of words and its associated sentiment score. These dictionaries can then be used in algorithms for various applications.

III. RELATED WORK

*A. Machine Learning*

[1] The authors of this paper propose an approach to sentiment analysis which makes use of support vector machines (SVMs) to combine diverse sources of potentially pertinent information. Support vector machine is a popular and a very powerful tool for classifying vectors of real-valued attributes.

This model uses kernel, a function that maps a huge space of data points and is generally used on data that is not linearly separable. Models that integrate the introduced features with unigram models have demonstrated better performance, according to the authors. They have demonstrated that these hybrid SVMs showed superior performance, and also produced the best published results over this data, when they combined SVMs based on real-valued favorability measure with feature-based unigram-style SVMs in experiments involving movie review data from Epinions.com. Further experiments on a smaller music reviews dataset, using a feature set annotated for topic, are also reported. The results obtained from this experiment have shown that including topic data into such models may yield better performance.

[3] The authors of this paper analyses the problem in document classification not by topic, but by its overall sentiment, e.g., classifying negative and positive reviews from a movie dataset. Three machine learning models (Support vector machines, naive bayes and maximum entropy classification) were implemented on the movie reviews dataset. These models performed poorly on traditional topic-based categorization compared to the overall sentiment approach. The machine learning algorithms that this paper employed are:

- *Naïve Bayes*

The main objective of this approach is to classify a document *d* based on its class *c* using the Bayes rule. The equation (1) and (2) are used for classification.

$$P(c | d) = \frac{P(c)P(d | c)}{P(d)}, \tag{1}$$

$$P_{NB}(c | d) := \frac{P(c) \left( \prod_{i=1}^m P(f_i | c)^{n_i(d)} \right)}{P(d)}. \tag{2}$$

- *Maximum Entropy*

Maximum Entropy (ME, or MaxEnt, for short) is an alternative method which has been proven to be more effective in applications that uses natural language processing (NLP).

We can estimate  $P(c|d)$  using the equation (3).  $Z(d)$  is the normalization function.

$$P_{ME}(c | d) := \frac{1}{Z(d)} \exp \left( \sum_i \lambda_{i,c} F_{i,c}(d, c) \right), \tag{3}$$

- *Support Vector Machines*

Support vector machine is a powerful tool and is widely used for classifying vectors of real-valued attributes. This algorithm uses kernel, a function that maps a huge space of data points and is generally used on data that is not linearly separable.

After performing various tests, the authors conclude that support vector machines give better results in traditional text categorization, outperforming maximum entropy and Naïve Bayes. The results of the aforementioned techniques are shown in Table I, which presents average three-fold cross-validation accuracies, in percentage. The statistics in boldface represent the model that showed the highest accuracy for a particular setting.

TABLE I. ACCURACY OF MODELS

Features <sup>a</sup>	# of features	Frequency or Presence	NB	ME	SVM
U	1615	freq.	<b>78.7</b>	N/A	72.8
U	1615	pres.	81	80.4	<b>82.9</b>
U + B	32330	pres.	80.6	80.8	<b>82.7</b>
B	16165	pres.	77.3	<b>77.4</b>	77.1
U + POS	16695	pres.	81.5	80.4	<b>81.9</b>
Adjectives	2633	pres.	77	<b>77.7</b>	75.1
U + P	22430	pres.	81	80.1	<b>81.6</b>

<sup>a</sup> U: Unigram; B: Bigram; P: Position

*B. Dictionary based approach*

[1] This paper attempts to perform sentiment analysis through a simple and effective approach by using semantic orientation with pointwise mutual information (PMI). A real number measure of the positive or negative sentiments expressed by a word or a phrase is called semantic orientation (SO). This method can also be used for multiple word phrases. The terms “value phrases” and “SO phrases” are used by the authors to refer to these phrases. The technique deployed is to calculate PMI for different phrases and use a fixed PMI as reference for calculating the relative value. PMI is calculated using the equation (4).  $PMI(w_1, w_2)$  in equation (4) is the probability of co-occurrence of  $w_1$  and  $w_2$ .

$$PMI(w_1, w_2) = \log_2 \left( \frac{p(w_1 \& w_2)}{p(w_1)p(w_2)} \right) \tag{4}$$

The difference between PMIs with the word “excellent” and the word “poor” for a phrase is called it’s SO. SO is calculated using equation (5).

$$SO(phrase) = \log_2 \left( \frac{\text{hits}(phrase \text{ NEAR } \text{“excellent”}) \text{hits}(\text{“poor”})}{\text{hits}(phrase \text{ NEAR } \text{“poor”}) \text{hits}(\text{“excellent”})} \right) \tag{5}$$

[7] This paper attempts to classify reviews as recommended (thumbs up) or non-recommended (thumbs

down) using a simple unsupervised learning algorithm. Reviews are classified based on the average semantic orientation (SO) of the phrases that consist of adverbs or adjectives in the review. The input to this algorithm is a written review, and the output is the classification. The initial step described by the authors is to use a part-of-the-speech tagger for identification of phrases in the input text that contain the adjective or adverb. The next step is to compute the semantic orientation of the extracted phrases. The PMI-IR algorithm can be used to estimate the SO of phrases that consist of adverbs or adjectives. PMI-IR uses Information Retrieval (IR) and Pointwise Mutual Information (PMI) to calculate the similitude between pairs of phrases or words.

The numerical rating of a phrase can be calculated by subtracting the mutual information between the word “poor” and the given phrase from the mutual information between the word “excellent” and the given phrase. The proposed classification algorithm’s performance is evaluated over 410 reviews from Opinions, sampled randomly over four different domains: reviews of travel destinations, automobiles, movies and banks. The SO (phrase) is calculated using equation (4) PMI is calculated using the equation (5).

The final conclusion obtained shows the result of the reviews whether it’s a positive or negative overall with accuracy. The thumbs up represents the number of positive reviews found and the thumbs down represents the negative reviews and the sum shows the total reviews found in a document. The results obtained from the experiment is shown in Table II.

TABLE II. CLASSIFICATION RESULTS

Average Semantic Orientation	Author's Classification		
	Thumbs Up	Thumbs down	Sum
Positive	28.33%	12.05%	40.83%
Negative	21.67%	37.50%	59.17%
Sum	50.00%	50.00%	100.00%

[2] This paper discusses the topic of emotion detection and its progress over time. The approaches described by the authors are:

- Keyword analysis and VOS viewer

Keyword analysis is an effortless approach. All the keywords are first indexed and is then processed to minimize the space. This data is used to produce results in the form of tables and graphs. VOS viewer is a program for constructing and viewing bibliometric maps. There are 3 steps in this approach: Similarity matrix, VOS mapping, translation, reflection and rotation, and the last step is the optimization step.

- Latent Dirichlet Allocation

Latent Dirichlet Allocation is a simple model that could be used for analyzing text documents. This model has the ability to capture multi-topic characteristics from the text. LDA presumes that data has patterns and is structured even if these characteristics cannot be examined without difficulty. This approach is based on the concept that the data in the document

consists of information related to a random combination of topics which are characterized by a particular sequence of words.

The authors conclude that keyword analysis approach is reliable and a quick way to analyze various text documents, as the publicists and authors themselves provide the keywords for this process. LDA is a very powerful tool, it is also capable of giving deeper insights as it perceives information beyond keywords and looks at words that have been used to provide a description for a topic.

[11] The authors of this paper propose a dictionary based technique for domain specific sentiment analysis. The authors used sentiwordnet (SWN), a lexicon which is a publically available dictionary that includes adjectives, adverbs, and verbs. The dataset chosen by the authors for this approach is the movie reviews dataset. The method proposed in the paper performs analysis using linguistic features ranging from adverb + adjective to adverb + adjective + verb combinations at the document level. For any kind of aspect based sentiment analysis, the author follow these 3 steps:

1. Identification of aspect from the review.
2. Locating the aspect.
3. Detecting the sentiment polarity of review.

The results of the lexicon-based approach proposed by the authors is compared with the performance of the Alchemy API. It is observed that the results of the sentiwordnet has higher accuracy than the Alchemy API.

The paper explores various aggregation, linguistic feature selection and weighting schemes. For sentiment classification at a document level, they used the SWN (AAC) and SWN (AAAVC) schemes and obtained results on the data using both these schemes. The author presents the values of performance metrics obtained for their implementations and compares them with the results of the Alchemy API.

It is observed that out of total of 760 actual positive reviews, Alchemy API labels 634, SWN (AAAVC) labels 688 and SWN (AAC) labels 678 as positive. Similarly, the three schemes label 140, 99 and 98 reviews as negative out of a total number of 240 actual negative reviews.

[12] This paper discusses an unsupervised learning approach to intra-sentence discourse relations for polarity classification at the sentence level. The authors have presented a conversation scheme based on practical results. Then, an unsupervised model was presented to analyze data ranging from a small set of phrase-cue-based patterns to a large set of common semantic sequential representations. The models performance was further enhanced by incorporating semantic sequential representations as features in supervised techniques. The results obtained from this experiment have shown that the implemented methods not only identified discourse relations effectively but have also achieved significant improvement ( $p < 0.01$ ) in the classification of sentence level polarity.

Although this paper proposes a method for text in Chinese, the same unsupervised model can be implemented to classify sentence level polarity in other languages.

#### IV. CONCLUSION

Sentiment analysis helps in the identification of a person’s attitude and emotional status. A human’s emotion is a

complicated metric that can be conveyed in negative or positive ways. The current emotional status of a human can contribute a lot in the advancement of computer-human communication. Humans can provide inputs to the computer in different forms, but in the era of Web 2.0, textual information is the most common form of input that is used to communicate with a computer. Therefore, emotion detection from textual data should be focused as an important research affair to enhance computer-human interactions in the future.

This paper demonstrates various theories that have been used in the development of models in the field of sentiment analysis. We have presented a literature review of the recent works in this field. This survey reviews models and approaches that were implemented for detecting emotion in textual data and also talks about its challenges and applications.

We hope that researchers in related fields may find the survey presented in this paper useful and advance scientific discussions regarding emotion detection from texts.

#### REFERENCES

- [1] Tony Mullen and Nigel Collier, National Institute of Informatics (NII), Hitotsubashi 2-1-2, Chiyoda-ku, Tokyo 101-8430, Japan "Sentiment analysis using support vector machines with diverse information sources".
- [2] Oskar Ahlgren "Research On Sentiment Analysis: The First Decade" 2016 IEEE 16th International Conference on Data Mining Workshops.
- [3] Bo Pang and Lillian Lee, Department of Computer Science, Cornell University Ithaca, NY 14853 USA. Shivakumar Vaithyanathan, IBM Almaden Research Center, 650 Harry Rd. San Jose, CA 95120 USA "Thumbs up? Sentiment Classification using Machine Learning Techniques" Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), Philadelphia, July 2002, pp. 79-86.
- [4] Priya Nambisan, PhD University of Wisconsin- Milwaukee. Zhihui Luo, PhD University of Wisconsin – Milwaukee. Akshat Kapoor, M.S. University of Wisconsin – Milwaukee. Timothy B Patrick, PhD University of Wisconsin – Milwaukee. Ron A Cisler, PhD University of Wisconsin – Milwaukee. "Social Media, Big Data and Public Health Informatics: Ruminating behavior of depression revealed through Twitter" 2015 48th Hawaii International Conference on System Sciences.
- [5] Dongkeon Lee, Kyo-Joong Oh, Ho-Jin Choi School of Computing, Korea Advanced Institute of Science and Technology (KAIST), Daejeon, Korea "The ChatBot Feels You – A Counseling Service Using Emotional Response Generation".
- [6] Kyo-Joong Oh, DongKun Lee, ByungSoo Ko, Ho-Jin Choi School of Computing, Korea Advanced Institute of Science and Technology (KAIST), Daejeon, Republic of Korea "A Chatbot for Psychiatric Counseling in Mental Healthcare Service Based on Emotional Dialogue Analysis and Sentence Generation".
- [7] Peter D. Turney, "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews", Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL).
- [8] Apoorv Agarwal, Boyi Xie, Ilia Vovsha, Owen Rambow, Rebecca Passonneau, "Sentiment Analysis on Twitter Data", Proceedings of the Workshop on Language in Social Media (LSM 2011), pages 30–38.
- [9] Ms.A.M.Abirami Dept. of Information Technology, Thiagarajar College of Engineering, Madurai, Tamil Nadu, India. Ms.V.Gayathri Dept. of Information Technology Thiagarajar College of Engineering Madurai, Tamil Nadu, India "A SURVEY ON SENTIMENT ANALYSIS METHODS AND APPROACH" 2016 IEEE Eighth International Conference on Advanced Computing (ICoAC).
- [10] H. N. Io1, C. B. Lee, Department of Accounting and Information Management, University of Macau, China "Chatbots and Conversational Agents: A Bibliometric Analysis".
- [11] V.K. Singh, R. Piryani, A. Uddin, P. Waila, "Sentiment Analysis of Movie Reviews", conference on International Multi-Conference on Automation, Computing, Communication, Control and Compressed Sensing, IEEE-2013.
- [12] Lanjun Zhou, Binyang Li, Wei Gao, Zhongyu Wei, Kam-Fai Wong, Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong, Shatin, NT, Hong Kong, China "Unsupervised Discovery of Discourse Relations for Eliminating Intra-sentence Polarity Ambiguities" Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, pages 162–171, Edinburgh, Scotland, UK, July 27–31, 2011.
- [13] Michael Gamon. 2004. "Sentiment classification on customer feedback data: noisy data, large feature vectors, and the role of linguistic analysis". Proceedings of the 20th international conference on Computational Linguistics.
- [14] Alec Go, Richa Bhayani, and Lei Huang. 2009. Twitter sentiment classification using distant supervision. Technical report, Stanford.
- [15] Apoorv Agarwal, Fadi Biadisy, and Kathleen Mckeown. 2009. Contextual phrase-level polarity analysis using lexical affect scoring and syntactic n-grams. Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009), pages 24–32.