# A Survey on Sentimental Analysis in Different Indian Dialects

**Shankar R[1], Shilpa K.M[1], Sridhar Patil[1], Suma Swamy[2]**

Department of Computer Science, BMSIT&M, Bangalore, India[1]

Professor, Department of Computer Science, Sir MVIT, Bangalore, India[2]

**Abstract:** Sentimental Analysis analyses the emotion present behind the speaker/writers reviews or comments. This emotion indeed helps in categorizing a review as positive or negative.In our day-to-day life we are involved in Podcasting, Blogging, Tagging, and Social Network (facebook, twitter, whatsapp etc) where numerous dataset gets generated. In other words, user's feelings are expressed in the form of comments/reviews/feedback. This opinion mining is quite a challenging task and basically everything in this regard is done in English language. If we think of other Indian languages as the resources like reviews and opinions are huge, study on these have not been done. In this paper, we are considering a survey on Sentiment Analysis and Opinion Mining in different languages like Kannada, Hindiand Malayalam. The speaker would be more comfortable in expressing his/her opinions using the native language as it would be more easy and precise.

**Keywords:** Opinion, Bag of words, Naïve Bayesian, corpus, POS tagging.

## I. INTRODUCTION

Sentimental Analysis is a Natural Language Processing and Information Extraction that plans to get author's emotions communicated in positive or negative remarks."What other individuals think" is natural phenomenon. The idea and opinion of others, their reviews have always affected our own opinion.Sentiment analytical systems are being applied in almostevery business and social domains since the result of opinions are central to most human activities and are the key influencers of our behaviours.Social media can provide immediate feedback and it represents a new challenge in communication between a customer and business. Sentiment prediction can be done at the documentlevel, sentence level and phrase level. In documentlevel the sentiment of the entire document issummarized as positive, negative or objective. Sentencelevel prediction classifies individual sentiment bearingsentences. At phrase level phrases in a sentence areclassified according to polarity.

Much of the research work in this regard has been done in English language. In India, we have more than 30 official regional languages where people communicate effectively thus producing huge amount of data. This data can well be addressed to extract numerous and useful information like the buyer's pattern, availability of the goods, product analysis, product feedback, originality etc. The aforementioned research activity is lagging in these languages. So we intend to study the same and it's various effects through machine learning techniques. As an initial start, we are focussing on performing a survey on this topic which has been done in regional languages like Kannada, Hindi and Malayalam.

Numerous forums, blogs, social networks, e-commerce web sites, news reports and additional web resources serve as platforms to express opinions, which can be utilized for understanding the opinions of the general public and consumers on social events, political movements, company strategies, marketing campaigns, product preferences, and monitoring reputations.

## II. LITERATURE SURVEY

A. Sentimental Analysis For Kannada Using Mobile Product Reviews A Case Study:

In this paper they have chosen a lexicon based approach where use of lexicon entity models for aspect extraction and have implemented a simple and statistical Naive Bayes classifier (algorithm) which is used to classify the sentiments of the mobile product reviews in Kannada. Here mobile is major aspect, weekly product reviews have been fed by 'GadgetLoka' from famous Kannada daily newspaper 'Prajavani' by U.B Pavanaja[3]. They have used this to create the corpus in study. The Bayes classifier takes Kannada text as input which is represented in Unicode format. The results indicate that the lexicon based aspect extraction with Naive Bayes sentiment classifier works efficiently for Kannada (Unicode) sentiment analysis. From the results, we can see that the technique is performing with 65 % accuracy, 62.5% precision (PPV), 75% (sensitivity)recall, 55% specificity, 68.75% NPV and 68.2 % F-Measure. The final objective is to apply this method for Kannada text.

B. An Analysis of Sentence Level Text Classification forKannada Language.

The Kannada language is epigraphically about one and half millennia. In this paper the review of positive and negative analysis would help in business development and recommending a system with difference. The size of the feature set has a significant impact on the time required for

classification. This makes dimensionality reduction a requirement for text classification on Kannada. Stop words does not hold the information of class text. In some other cases, sentences might not have sufficient class information as a standalone entity and might use neighbouring sentences to convey the class of information [1]. This can be captured to increase the performance of the classifier. Multi-label classification is appropriate for sentences that can be a part of two or more classes. Also, possible hierarchy among the classes can be explored to support multi label classification.

The work can be extended to online customer reviews in Kannada blogs. It can also be used in extracting opinion in posted Kannada articles online. Fine grained classification, i.e. at sentence level or subsistence level can be achieved using this approach. Depending on the requirements of the application – whether or not it requires high precision/recall, the appropriate methods can be chosen.

C.  HOMS: Hindi Opinion Mining System.
In this paper, the Hindi movie reviews are analysed. The document could be positive or negative or neutral. They are classified by performing opinion mining at the document level by using Machine Learning and Part of Speech Tagging (POS)[6]. In Machine learning, Naive Bayes Classifier is used. In POS tagging, adjectives are considered as opinion words and based on their count, the document is classified.The sentences for polarity detectionare done using two methods: Supervised using Naive Bayes, Unsupervised method by POS tagger.First it extracts the data from the online website,and then text mining is applied for classification. Second, it extracts adjectives and counts them.

D. A Practical Approach to Sentimental Analysis of Hindi Tweets.
Twitter is enormously utilized as a part of communicating in the form of tweets.Twitter is used to express user's opinion about the particular entity which has limited character of 140. These opinions and reviews are very helpful for any other person who is interested in that particular entity. The technique is to extract an opinion from the tweets and defining them as positive, negative or neutral. The work proposed here is sentimental analysis on Hindi tweets[7]. Every few minutes the tweet is extracted by making connection with twitter using Google Script. Pre-processing is done by deleting the extra symbols in tweets and stop words are also removed. A Senti-word is created which contains an adjective. The proposed algorithm uses subjective lexicon method.

E. Domain Specific Sentence Level Mood Extraction from Malayalam Text.
In this paper initial step is to gather corpus from Malayalam novels. In order to avoid grammatical mistake, sentences are manually typed. As much as 100 sentences are collected initially. Words present inside the sentences are tagged manually as adjectives and adverbs. These adjectives and adverbs present in sentence define the emotion of the sentence. The SO-PMI-IR classifier

classifies input into two types classes-desirable or no desirable. This results to classification of input to –Joy, Sorrow, Anger or Neutral. In this paper instead of Naïve Bayes, Maximum Entropy etc, the SVM machine learning is used to provide more accuracy [8].

## III. DIFFERENT LEVELS OF SENTIMENTAL ANALYSIS

A.  Document Level Sentiment Analysis
The basic information unit could be a single document of opinion text. During this document level classification, one review about a single topic is taken into account. However withthe case of forums or blogs, comparative sentences could appear. Customers could compare one product with another which has similar characteristics. The challenge within the document level is classification of sentence into a document and this document might not be relevant in expressing the opinion regarding an entity. So subjectivity/objectivity classification is extremely necessary during this kind of classification. The irrelevant sentences should be eliminated from the process works. Supervised and Unsupervised learning ways can be used for the document level classification. Any supervised learning algorithm like Naïve Bayesian, Support Vector Machine, is used to train the system.

For training and testing knowledge, the reviewer rating (in the shape of 1-5 stars), is used[13]. The categories which will be used for the machine learning are term frequency, adjectives from part of speech tagging, Opinion words and phrases, negations, dependencies etc. Labeling the polarities of the document manually is time consuming. The unattended learning is done by extracting the opinion words within the document. The point-wise mutual information is created to seek out the linguistics of the extracted words. So the document level sentiment classification has its own benefits and drawbacks. Advantage is that we tend to get associated the overall polarity of opinion within a couple of specific entities. Disadvantage is that a completely different document withvarious emotionswould not get processed and categorized as Positive or Negative polarity.

B. Sentence Level Sentimental Analysis.
In the sentence level sentiment analysis, the polarity of every sentence is calculated. Identical document level classification strategies are often applied to the sentence level for classification problem. The subjective sentences contain opinion words that facilitate in determining the sentiment regarding the entity. Just in case of easy sentences, one sentence bears one opinion regarding an entity. However there can be some sentences with difference of opinionated texts [13]. In such cases, sentence level sentiment classification isn't desirable. Knowing that a sentence is positive or negative is of lesser use than knowing the polarity of a specific feature of a product. The advantage of sentence level analysis lies within the subjectivity/ objectivity classification. The traditional algorithms are often used for the training processes.

## C. Phrase Level Sentimental Analysis.

The phrase level sentiment classification may be a way more pinpointed approach to opinion mining. The phrases that contain opinion words are noticed and a phrase level classification is done. This may be advantage or disadvantage. In some cases, the precise opinion regarding an entity is often properly extracted [13]. However in another cases, wherever discourse polarity also matters, the result might not be absolutely correct. Negation of words may occur. In such cases, this level of sentiment analysis handles the negation in a more or less precise way. However if there are sentences with negating words that occur with the exception of the opinion words, phrase level analysis isn't desirable. Conjointly, varied dependencies aren't considered here. The words that seem terribly around each other are thought-about to be in a phrase.

## IV. POLARITY

Polarity detection is an element of sentimental analysis where sentences or documents are classified as negative or positive, [2] machine learning algorithms are often applied.Opinion words like stunning, lovely, nice represent positive polarity and unhealthy, ugly, atrocious represent negative polarity. In the table, few words are listed as positive and negative words that define polarity. Collection of opinion words or phrases employed in sentiment classification is named as opinion lexicon. A corpus is extracted automatically from twitter and other resources mentioned above and a classifier is made to work out positive, negative and neutral sentiments. A lexicon based approach is used to extract sentiment from text. An example tool referred to as SentiStrengthuses a lexicon based sentiment analysis of short text. Every word within the lexicon is assigned a score between +5(very positive) and -5(very negative)[2]. Sentiment Analysis in Indian Languages may be a new topic for analysis. WordNet, POS tagger, SentiWordNet and alternative resources are out there for only a few languages.

TABLE I classification of words

| Positive words | Negative words |
| --- | --- |
| Good | Bad |
| Awesome | Worst |
| Superb | Terrible |
| Constantly | Struggled |
| Thrilling | Unbelievable |
| Funny | Boring |
| Sleep | Miserable |
| Genius | Profanity |
| Witty | Disliked |
| Funny | Skip |
| Wait | Dull |
| Crafted | Annoyed |
| Amazing | Sadly |
| Complex | Unlinked |
| Enjoys | Cliché |
| Wow | Selfish |

| | |
| --- | --- |
| Roller | Unrealistic |
| Unfold | Implausible |
| Love | Sorry |
| Beautiful | Contrived |
| Thrilling | Worse |
| Favor | Poorly |
| Turner | Wasted |
| Genius | Ridiculous |
| Terrific | Dumb |
| Fascinating | Depressing |
| Immediately | Terrible |
| Woven | Silly |
| Pieces | Unlikeable |
| Coaster | Quit |
| Insightful | Mess |
| Twice | Hated |
| Edge | Implausible |
| Blown | Struggled |
| Prime | Bothered |
| Brilliant | Negative |
| Perfect | Mistake |
| Love | Death |

A SentiWordNet development for Indian Languages like Bengali, Hindi and Telugu is tried using WordNet, and Corpus based ways in [7]. In [6], a Hindi sentiment lexicon known as Hindi-SentiWordNet (H-SWN) was developed using English SentiWordNet (SWN) and English-Hindi WordNet linking. The sentiment analysis of documents was performed using sentiment annotated corpora, machine translation and resource based sentiment analysis ways.Stop words are words that don't hold information regarding the category of the text. The meaningless words of a language are typically known as stop words.

## V. APPLICATIONS AND CHALLENGES IN SENTIMENTAL ANALYSIS

There is a large explosion of 'sentiments' available from social media at the side of Twitter, Facebook, message boards, blogs, and user forums. These firms help other people to look at their reputation and notice timely feedback relating to their product and actions. Sentiment analysis offers these organizations the facility to look at the various social media sites in real time for the development of business. Selling managers, PR firms, campaign managers, politicians, and even equity investors and web consumers are the direct beneficiaries of sentiment analysis technology. Twitter and Facebook iscentre of attention of the various sentiment analysis applications. One application that performs analysis of tweets that contain a given term is tweetfeel (http://www.tweetfeel.com)[4]. Another important domain for sentiment analysis is that of the stock markets. A sentiment analysis system can use these various sources to look out an article that discusses the businesses and utilises an automatic commerce system. One such system isthat the Stock instrument (http://www.thestocks onar.com). This technique (developed by Digital Trowel)

shows graphically the daily positive and negative sentiment concerning each stock aboard the graph of the worth of the stock.

Challenges in sentimental analysis are as follows. Positive or negative words may have same meaning. We can write a sentence like "this phone sucks". This is negative sentence but "this vacuum cleaner sucks" is a positive sentence. A sentence may not contain any sentiment like "can u tell me which phone is good" and "if I find a good camera in the shop I will buy it" here both sentences do not express positive nor negative, so in order to classify, it becomes difficult.

Sarcasm sentence with or without sentiment word are difficult to handle: example- "what a great car! But stopped working in two days". Sarcasm are not so common in consumer reviews about product and services, but are very common in political discussion, it becomes difficult to deal with the sentences.

## VI. SPAM IN SENTIMENTAL ANALYSIS

As opinions are necessary for several applications, it's no surprise that individuals have began to cheat the system. Opinion spam refers to faux or unreal opinions that try and mislead readers or machine-controlled systems by giving unworthy positive opinions to some target objects so as to market the objects and/or by giving malicious negative opinions to another object so as to wreck their reputations. Detection of such spam is incredibly necessary for applications. Mechanical assignment of values to opinions is beneficial as opinions will then be graded with their utility values. With the ranking, the reader will target those quality opinions. We must always note, however, that spam and malicious entries are totally different ideas, detection is a component of sentimental analysis wherever sentences or documents are classified as negative or positive. For this, Machine learning algorithmic rule will be applied [11].

E-mail spam and internet spam are quite familiar to most people. E-mail spam refers to unsought industrial e-mails marketing products and services [5], whereas internet spam refers to the use of "illegitimate means" to boost the search rank positions of target websites.

The reason for spam is especially because of political economy. As an example, within the internet context, the economic and/or publicity price of the rank position of a page came back by a search engine is of great importance. If somebody searches for a product that your computing device sells, however the product page of your website is stratified terribly low (e.g., on the far side the highest 20) by a search engine, then the possibility that the person can move to your page is very low, coupled with to buy the product from your site. This can be actually unhealthy for the business. There are currently several firms that are within the business of serving to others improve their page ranking by exploiting the characteristics and weaknesses of current search ranking algorithms. These firms are referred to as search engine optimization (SEO) firms.

Due to the explosive growth of the user-generated content, it's become a typical observation for people to seek out and to scan opinions on the net for several purposes. For instance, someone plans to shop for a camera. Most likely, he or she's going to attend a merchant or reviewer web site (e.g., amazon.com) to scan the reviews of some cameras. If he or she finds that the majority reviews are positive, he or she is probably for sure to shop for the camera. However, if most reviews are negative, he or she's going to actually select another camera. Positive opinions may end up in important financial gains and/or fames for organizations and people. This, sadly, also gives good incentives for opinion spam, which refers to human activities[5] writing spam reviews in order to deliberately mislead readers or automated opinion mining systems.

## VII. TOOLS IN SENTIMENTAL ANALYSIS

In order to extract emotions, certain tools are available. These tools are essential for data processing classification of images and text. Few tools are listed below.

WEKA: Machine learning algorithm for Data Mining which does Tagging, Parsing, provides lexical resources such as WordNet.

STANFORD CORENLP: helps in POS tagging, Named entity recognizer, Parsing, Co-reference resolution system, and bootstrapped pattern learning.

Robust Accurate Statistical Parsing: provides Statistical Parser, Tokenization, POS tagging, N-gram search.

NLTK: Classification, Tokenization, Stemming,

LingPipe: Entity extraction, POS tagging, Clustering, Lemmatization and Parsing.

GATE: Tokenizer, Gazetteer, Sentence splitter, POS Data pre-processing, Classification, Regression, Clustering, Association rules, Visualization.

APACHE OPENNLP: Tokenization, Sentence segmentation, Part-of-speech tagging, Named entity extraction, Chunking, Parsing, Co reference resolution.

## VIII. CONCLUSION

In this paper, we had a walk-through on various aspects of sentimental analysis in Indian Dialects. The discussion of tagging, classification of the sentences or reviews in to positive or negative polarities was wholesome. The challenges faced in all these languages are quite intriguing and essentially it gives the researchers an edge in performing opinion mining task in these regional languages as the information is abundant.

## ACKNOWLEDGMENT

# REFERENCES

[1] Jayashree R, Srikanta Murthy K, "An Analysis of Sentence level Text Classification for the Kannada Language"2011IEEE website.

[2] Deepamala. N, Dr. Ramakanth Kumar. P,"Polarity Detection of Kannada documents"2015IEEE website.

[3] YashaswiniHegde, S.K. Padma," Sentiment Analysis for Kannada using MobileProduct Reviews A Case Study" IEEEwebsite2015.

[4] HaseenaRahmath P "Opinion Mining and Sentiment Analysis - Challenges and Applications" international journal 2014.

[5] Chandni,"Sentiment Analysis and its Challenges "international journal 2015.

[6] VandanaJh, Manjunath N,P Deepa Shenoy," HOMS: Hindi Opinion Mining System" , 2015.

[7] Yakshi Sharma,VeenuMangat,Mandeep Kaur," A Practical Approach toSentimentAnalysis of Hindi Tweets"2015.

[8] Neethu Mohandas, Janardhanan PS Nair, Govindaru V," Domain Specific Sentence Level Mood Extraction from Malayalam Text".

[9] Bo Pang1 and Lillian Lee2, "Opinion Mining and Sentiment Analysis",Foundations and TrendsR_ in Information Retrieval Vol. 2, Nos. 1–2 (2008) 1–135.

[10] Subhabrata Mukherjee," Sentiment Analysis A Literature Survey"2012.

[11] Pravin Jambhulkar "A Survey Paper on Cross-Domain Sentiment Analysis".

[12] Raisa Varghese1, "a survey on sentiment analysis and opinion mining"

[13] Dr. Ritu Sindhu "A Novel Approach for Sentiment Analysis and Opinion Mining"