



## Sentiment Analysis and Opinion Mining: A Concept

Puja M. Dadhe<sup>1</sup>, Dilip S. Sadhankar<sup>2</sup> and R.N. Jugele<sup>3</sup>

Department of Computer Science, Shivaji Science College, Nagpur.<sup>1,3</sup>

Department of Computer Science, SFS college, Nagpur<sup>2</sup>

poojadadhe@gmail.com<sup>1</sup> dileep.sadhankar@gmail.com<sup>2</sup> rn\_jugele@yahoo.com<sup>3</sup>

### Abstract-

An important part of our information-gathering behavior has always been to find out what other people think. With the growing availability and popularity of opinion-rich resources such as online review sites and personal blogs, new opportunities and challenges arise as people now can and do, actively use information technologies to seek out and understand the opinions of others. The sudden eruption of activity in the area of opinion mining and sentiment analysis, which deals with the computational treatment of opinion, sentiment and subjectivity in text, has thus occurred at least in part as a direct response to the surge of interest in new systems that deal directly with opinions as a first-class object. In marketing and advertising domains Opinion Mining is being larger domain. Advertiser needs to analyze performance/ popularity of ads that posted on site. Star rating based mechanism may go fraud, because of robots or automatic responders. So, current system needs to be analyzed using comments & natural language processing.

**Keywords-** Sentiment, Opinion, machine learning, mood, reception, tweets, lexicon.

### I. INTRODUCTION

Sentiment analysis is a type of natural language processing for tracking the mood of the public about a particular product or topic. Sentiment analysis, which is also called opinion mining, involves in building a system to collect and examine opinions about the product made in blog posts, comments, reviews or tweets. Sentiment analysis can be useful in several ways. For example, in marketing it helps in judging the success of an ad campaign or new product launch, determine which versions of a product or service are popular and even identify which demographics like or dislike particular features [1]. Sentiment analysis or opinion mining is the computational study of peoples' opinions, appraisals, attitudes and emotions toward entities, issues, events, topics and their attributes. The task is technically challenging and practically very useful. For example, businesses always want to find public or consumer opinions about their products and services. Potential customers also want to know the opinions of existing users before they use a service or purchase a product. With the explosive growth of social media (i.e., reviews, forum discussions, blogs and social networks) on the Web, individuals and organizations are increasingly using public opinions in these media for their decision making. However, finding and monitoring opinion sites on the Web and distilling the information contained in them remains a formidable task because of the proliferation of diverse sites. Each site typically contains a huge volume of opinionated text that is not always easily deciphered in long forum postings and blogs. Opinion Mining has become very important after the advent of various social networking sites.

### II. DATA SOURCE

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

#### 2.1. 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. Bloggers record the daily events in their lives and express their opinions, feelings, and emotions in a blog [ii]. Many of these blogs contain reviews on many products, issues, etc. Blogs are used as a source of opinion in many of the studies related to sentiment analysis [4].

#### 2.2. Review sites

For any user in making a purchasing decision, the opinions of others can be an important factor. A large and growing body of user-generated reviews is available on the Internet. The reviews for products or services are usually based on opinions expressed in much unstructured format. The reviewer's 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.CNETdownload.com (product reviews) and www.reviewcentre.com, which hosts millions of product reviews by consumers. Other than these the available are professional review sites such as www.dpreview.com, www.zdnet.com and consumer opinion sites on broad topics and products [6].

## 2.4. Micro-blogging

Twitter is a popular micro blogging service where users create status messages called "tweets". These tweets sometimes express opinions about different topics. Twitter messages are also used as data source for classifying sentiment.

### III. TERMINOLOGY

With the advent of social media, data is captured from different sources, such as mobile devices and web browsers and it is stored in various data formats. As the social media content is unstructured with respect to traditional storage system (such as RDBMS) tools are needed that processes and analyze this disparate data. Big data technology is made to handle the different sources and different formats of the structured and unstructured data. Sentiment Analysis is done on three levels.

- **Document level-** Document Level Sentiment Analysis is performed for the whole document and then decide whether the document express positive or negative sentiment.
- **Sentence level-** Sentence level Sentiment Analysis is related to find sentiment from sentences whether each sentence expressed a positive, negative or neutral sentiment and is closely related to subjectivity classification. Many of the statements about entities are factual in nature and still carry sentiment. Current Sentiment Analysis approaches express the sentiment of subjective statements and neglect such objective statements that carry sentiment.
- **Aspect or Entity level-** Entity or Aspect Level Sentiment Analysis performs fine grained analysis. The goal of it is to find sentiment on entities and/or aspect of those entities. For example consider a statement "My HTC Wildfire S phone has good picture quality but it has low phone memory storage." so sentiment on HTC camera and display quality is positive but the sentiment on its phone memory storage is negative.

Sentiment analysis is usually conducted between two levels

- A coarse level

A fine level Coarse level sentiment analysis deals with determining the sentiment of an entire document and Fine level deals with attribute level sentiment analysis. As per technical perspective there are two main approaches for sentimental analysis. They are given below:

**A. Symbolic Techniques:** In July 2013, Neethu M S and Rajasree R [5] proposed that Symbolic techniques also known as knowledge

based approach. In this technique, available lexical resources are used. In this sentiment analysis approach, bag-of-words approach is used. The BOW model focuses on the words list or says string of words and it cannot check the context of the sentence. This model contains a list of words that have own value when found in the given text, totally focuses on the words and take care nothing about the language fundamentals. The difficulty in using a Knowledge base approach is that it requires a large lexical database. This has become harder and harder to provide as the language of social networks is so trend dependent and changeable that lexicon datasets cannot keep up. Therefore, Knowledge based approaches to sentiment analysis are not as popular as they are used to be.

**B. Machine Learning Techniques:** In contrast to Knowledge based approaches, Machine Learning techniques are not using any lexicon resources list, instead a training set and a test set is used in order to classify them. Training set contains input vectors and corresponding class labels for training the network. After that test set is used to validate the given model by checking the class labels to unknown feature vectors. There are different machine learning techniques like SVM, maximum entropy and Naïve Bayes etc. This allows the algorithm to remain dynamic in the face of ever changing social network language lexicons. In this methodology, a classification model is developed using a training set which tries to classify the input feature vectors into corresponding class labels. It uses the results from the knowledge based techniques and those of the machine learning techniques to ensure a thorough analysis of the dataset.

### IV. METHODOLOGIES PROPOSED

Opinion features such as reviews on a particular product are typically domain-specific. The feature appears frequently in the given review domain, and which are outside the domain is domain-independent corpus about product [7]. Domain-specific opinion features are mentioned more frequently in the domain corpus of reviews, as compared to a domain-independent corpus. A domain-dependent review corpus and a domain-independent corpus is observed.

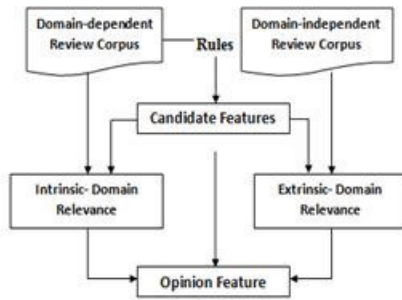


Figure 1: IEDR Workflow

Figure 1 show that, first extract a list of candidate features from the review corpus by defining manually syntactic rules. Each extracted candidate feature, will estimate its IDR, which represents the statistical association of the candidate to the given domain corpus, and extrinsic-domain relevance, will reflects the statistical relevance of the candidate to the domain-independent corpus. Only candidates with IDR scores more exceeding a predefined intrinsic relevance threshold and EDR scores less than another extrinsic relevance threshold are extracted as valid opinion features. In short, this paper identifies opinion features that are domain-specific and at the same time domain-independent corpus are removed and ignored.

**3.1 Candidate Feature Extraction**

Opinion features appear as the subject or object of a review sentence are generally nouns or noun phrases. In the dependence grammar, the subject opinion feature has a syntactic relationship of type subject verb with the sentence predicate. The object opinion feature has a dependence relationship of verb-object on the predicate [2]. It also has a dependence relationship of preposition-object on the prepositional word in the sentence.

**3.2 Opinion Feature Extraction**

Domain relevance characterizes how much a term is related to a particular corpus based on two kinds of statistics, dispersion and deviation. Dispersion identifies how significantly a term is mentioned in overall documents by measuring the distributional significance of the term across different documents in the entire domain. Deviation results about how frequently a term is mentioned in a particular document by measuring its distributional significance in the document. Both dispersion and deviation are calculated using the frequency-inverse document frequency term weights which is a well known technique.

**V. PROPOSED SYSTEM ARCHITECTURE**

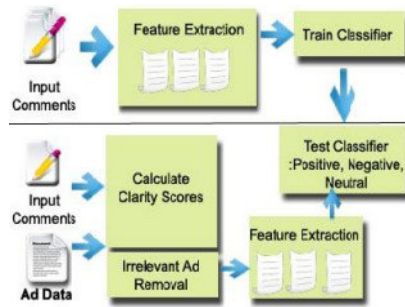


Figure 2: System Architecture

The system Architecture of proposed system is as shown in the figure 2. In first part of the system, it is shown that input will be collected from various online shopping websites such as amazon, flipcart, snapdeal, Jabong etc[3]. The comments which are written by the customers about any product in textual format in natural language is collected and those comments are used for feature extraction. After feature extraction the comments are passed to the trainer classifier for finding the patterns of comments. To identify patterns, techniques like N-Gram Extraction and part of speech Extraction are used by the trainer classifier. Collection comments and identifying patterns of comments is the *online* process of this system.

In second part of figure of the system *offline* process is shown. From the collected comments which are in natural language textual format the irrelevant comments are removed and clarity score is calculated. To remove irrelevant comments K-L Divergence algorithm is used and clarity score is also calculated using a threshold value. In this process the domain-dependent features and domain-independent comments are separated for feature extraction. In feature extraction NER tagger, Naive Bays classifier and porter streaming algorithms are used. With the help of trainer classifier and feature extraction the test classifier gives the feedback about a specified product as positive, negative and neutral.

**VI. CONCLUSION**

In Future, the social networks can give perfect solution to the problem of opinion acquisition and dissemination and perceived as natural enablers for opinion mining applications. In this paper, concept presented a proof of, examples of analysis that aim at gathering user opinions in two different application areas. Both experiments suggest that the networks fuelling the websites in question provide relevant context for opinion mining. The system aware of the fact that has not

utilized the information from the social network directly in the opinion mining algorithm. Merely, the system has tested the ability to attain high accuracy and quality of sentiment prediction using the data harvested from a social network site. It includes the user's reception of opinions contained in the text and further improvements of the presented all expect to attain the improvement of classification performance due to the utilization of information derived from the social networks, namely, the information on relationships and connections between users. We also intend to develop an active learning strategy for this type of classification task. The system suggested by this paper can be used in online marketing field as well as advertising field. The same system can be also used in any field where feedback about service can be collected. For example in hotels, railway services, about teacher. It can be also implemented for different languages.

#### REFERENCES

- [1] B. Liu, 2012, "Sentiment Analysis and Opinion Mining," Synthesis Lectures on Human Language Technologies, vol. 5, no. 1, pp. 1-167.
- [2] BoPang and Lillian Lee, 2008, "Opinion mining and Sentiment analysis:", Foundation and Trends in Information Retrieval, VOL.2.No. 1-2.
- [3] G. Vinodhini, RM.Chandrasekaran, 2012, "Sentiment Analysis and Opinion Mining: A Survey", International Journal of Advanced

Research in Computer Science and Software Engineering, Volume2, Issue 6

- [4] Martin, J. (2005). Blogging for dollars. Fortune Small Business, 15(10), 88-92.
- [5] Neethu M S and Rajasree R, July 4 - 6, 2013, "Sentiment Analysis in Twitter using Machine Learning Techniques", fourth ICCCNT international conference, Tiruchengode, India.
- [6] Popescu, A. M., Etzioni, O.: Extracting Product Features and Opinions from Reviews, In Proc. Conf. Human Language Technology and Empirical Methods in Natural Language Processing, Vancouver, British Columbia, 2005, 339-346.
- [7] Swati N. Manke, Nitin Shivare, 2015, "A Review on: Opinion Mining and Sentiment Analysis based on Natural Language Process", International Journal of Computer Applications(0975-8887) Volume 109-No.4.

#### VII. REFERENCE BOOK

- [i] A Survey of Opinion Mining and Sentiment analysis, Book by Bing Liu and Lei Zhang, University of Illinois, Chicago.
- [ii] Using Web Mining & Social Network analysis to study the emergence of cyber communities in blogs by Micheal Chau and Jennifer Xu.

