

Opinion Mining and Sentiment Analysis Based On Natural Language Processing

Ms. Swati N. Manke

Computer Engineering Department
BSIOTR, Wagholi
Pune, India

e-mail: engg.swatimanke@gmail.com

Prof. Nitin Sivale

Computer Engineering Department
BSIOTR, Wagholi
Pune, India

e-mail: nitinrajni3@gmail.com

Abstract— In marketing and advertising domains Opinion Mining is being larger domain. Advertiser as well as customer needs to analyze performance and popularity of product. Till now Star rating based mechanism is being used to analyze the performance and popularity of the product. The star rating mechanism uses the number of star ratings obtained by the product which may go fraud because of robots or automatic responders. So, the system needs to be analyzed using comments natural language processing. The proposed system collects the comments written by the customer about the product relevant with respect to opinion mining and by using Naive Bayes algorithm the popularity of the product is analyzed. False positive and false negative comments can be removed by using irrelevant comment removal mechanism. This system presents basic definitions used in opinion mining area which is based on natural language processing. The results obtained using the proposed system are so accurate and soundly support and overcome the problems in the existing system. The system can be used the by any online marketing website as well as in any field where the feedback from the customers can be collected.

Keywords- *Opinion, Sentiment, Domain-dependent, Domain-independent, Corpus, Relevance.*

I. INTRODUCTION

In the past few years due to many challenging research problems and practical applications opinion mining (or sentiment analysis) has attracted a great deal of attention from researchers of natural language processing and data mining. Opinion lexicon expansion and opinion target extraction are two fundamental problems in opinion mining. An opinion lexicon is a list of opinion words such as good, excellent, poor, and bad used to indicate positive or negative sentiments. Sentence and document sentiment classification and feature-based opinion summarization forms the foundation of many opinion mining tasks. Opinion targets are related to the topics on which opinions are expressed. They are important because without knowing the targets, the opinions expressed in a sentence or document are of limited use. The computational treatment of opinion, sentiment, and subjectivity has recently attracted a great deal of attention, in part because of its potential applications. It also has proven useful for companies, recommender systems, and editorial sites to create summaries of people's experiences and opinions that consist of subjective expressions extracted from reviews or even just a review's polarity positive or negative. Data-driven methods, resisting traditional text categorization techniques are document polarity classification poses a significant challenge. Opinion Mining is being larger domain of the marketing and advertising domains. Advertiser needs to analyze performance popularity of advertise that a person 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 and natural language processing. Fraud comments could be removed by using irrelevant comment removal mechanism suggested in proposed system.

In the proposed system the role and importance of social networks as preferred environments for opinion mining and sentiment analysis especially. Briefly, describing selected properties of social networks that are relevant with respect to opinion mining and

outline the general relationships between the two disciplines. In the proposed system the related work and provide basic definitions used in opinion mining is presented.

Then, introduce the original method of opinion classification and test the presented algorithm on real world data sets acquired from popular Polish social networks, reporting on the results. The results are promising and soundly support the main thesis of the paper, namely, that social networks exhibit properties that make them very suitable for opinion mining activities.

II. RELATED WORK

The document, sentence, or even phrase (word) levels can be generally analyzed Opinions and sentiments expressed in text reviews. The document-level (sentence-level) opinion mining is used to classify the overall subjectivity or sentiment expressed in an individual review document (sentence). Wei Jin and Hung Hay Ho [2] proposed a novel and robust machine learning approach for web opinion mining and extraction. This model provides solutions for server problems that have been not provided by previous approaches. This system can self learn new vocabularies based on the pattern its has learned, which is used in text and web mining. A novel bootstrapping approach is used to handle situations in which collecting a large training set could be expensive and difficult to accomplices. In the proposed system the effectiveness of proposed approach in web opinion mining and extraction from product review are demonstrated in result. Guang Qiu, Bing Liu, Jiajun Bu and Chun Chen [4] focuses on to important tasks in opinion mining that are opinion lexicon expansion and target extraction. In the proposed system, a propagation approach to extract opinion words and targets iteratively given only a seed opinion lexicon of small size. The identified relations between opinion words and targets are used for extraction in the proposed system a novel method for new opinion words polarity assignment and noisy targate priming are proposed. The new approach is compared with others on standard testing data set. The result of this

paper shows that this approach out performance other state-of-the-art methods. Bo Pang and Lillian Lee [6] examine the relation between subjectivity detection and polarity classification. The subjectivity detection can compress reviews in shorter extracts that still retains polarity information at a level comparable to that of the full review. By using Naive Bayes polarity classifier the subjectivity extract are shown to be more effective input than the originating document. The paper shows that the minimum-cut framework results in the development of efficient algorithm for sentiment analysis. Via this framework contextual information can lead to statistically significant improvement in polarity classification accuracy. Niklas Jacob and Iryna Gurevych [3] show how a CRF-based approach for opinion target extraction in a single and cross domain setting. In the proposed system a comparative evaluation of this approach on data set from four different domains are presented. The CRF-based approach out performance a supervised baseline on all dataset in the single domain setting. The CRF-based approach also yields promising results in the cross domain setting. Maniquing Hu and Bing Liu [7] proposed a number of techniques for mining features from product reviews based on data mining and natural language processing method. To produce feature based summary of a large number of customer review of a product sold online is the objective of this paper. Opinion mining becomes increasingly important as more people are buying and expressing their opinion on the web. Experimental result of this paper indicate that the propose techniques are effective in performing tasks. Opinion mining based on document, sentence, or phrase (word) level does not represent what exactly people like or dislike.

A. Extraction of Opinion Features

A sub problem of opinion mining is opinion feature extraction, with the vast majority of existing work. Previous approaches of opinion mining are classified into two supervised and unsupervised. Supervised learning models including hidden Markov models and conditional random fields have been used to tag features or aspects of commented entities. Supervised models can be carefully tuned to perform well on a particular domain, but need extensive retraining when applied to a different domain. A decent-sized set of labeled data is generally needed for model learning on every domain. In review sentences, an unsupervised NLP approaches extract opinion features by mining syntactic patterns of features are implied. The approaches attempt to discover syntactic relations among feature terms and opinion words in sentences by using carefully crafted syntactic rules, or semantic role labeling. Syntactic relations identified by the methods helps to locate features associated with opinion words, but could also inadvertently extract large number of invalid features due to the colloquial nature of online reviews. The results of statistical analysis on a given corpus to understand the distributional characteristics of opinion features, unsupervised corpus statistics approaches used. This approaches resistance to the colloquial nature of online reviews given a suitably large review corpus. Hu and Liu proposed [8] an association rule mining (ARM) approach to mine frequent item sets as potential opinion features, which are nouns and noun phrases with high sentence-level frequency (or support). ARM, relies on the frequency of item sets, has the some limitations for the task of feature identification, 1) frequent but invalid features are extracted incorrectly, and 2) rare but valid

features may be overlooked. To address feature-based opinion mining problems, a mutual reinforcement clustering (MRC) introduced by Su et al. [13] this approach is used to mine the associations between feature categories and opinion word groups, based on a co-occurrence weight matrix generated from the given review corpus. MRC is able to extract infrequent features, provided that the mutual relationships between feature and opinion groups found during the clustering phase is accurate which is unlike to other methods. Due to the difficulty in obtaining good clusters on real-life reviews, MRC's precision is low. The existing approaches to feature extraction only use the knowledge or patterns mined from a given single review corpus, by completely ignoring the possible variations present in a different domain-independent corpus. In the proposed system, IEDR approach utilizes the fact that word distribution characteristics vary across different types of corpora, in particular domain- specific versus domain-independent, to derive powerful hints that help discriminate valid features from the invalid ones. In the first step of this approach, some syntactic dependence rules to extract candidate features, similar to NLP approaches are defined. In the second step, employ the IEDR measures to identify the desired domain-specific opinion features. The key difference between IEDR compared to existing methods lies in its smart fusion of domain-dependent and domain-independent information sources.

III. IMPLEMENTATION DETAILS

A. Problem Statement

In the advertising and marketing field the current system is based on star rating system. The people also make their opinion about a particular product on star rating. As the star rating can be increase by the false users as well as some robots can be also designed to make the fake ratings the existing system is unable to give the accurate result. So there is a need of developing such a system which is able to analyze the popularity on the basis of comments wrote by the users. The proposed system will able to mine users' intent from comments. Then Irrelevant comments removal will increase opinion mining performance of system and False positive and false negative rates may reduced. The system can also Resistant to fake opinion postings. According to above discussion the existing system having some disadvantages as-

1. Star rating systems are easy to attack.
2. By visiting and like the web page of a product again and again by the same person the star rating can be increase.
3. False positive and false negative rates are more.

The above problems can be overcome by the system suggested by this paper. So the problem statement of this paper is as below.

1. To classify opinion about any product as Positive, Negative and neutral.
2. False positive and false negative rates can be reduced.
3. Fake opinions can be reduced.

B. Overview

An opinion feature such as "screen" in cell phone reviews is typically domain-specific. That is, the feature appears frequently in the given review domain, and rarely outside the domain such as in a domain-independent corpus about Culture. As such, domain-specific opinion features will be mentioned more frequently in the domain corpus of reviews, compared to a domain-independent corpus. Given a domain-dependent review corpus and a domain independent corpus, we first extract a list of candidate features

from the review corpus via manually defined syntactic rules (denoted "Rules" in the figure). For each extracted candidate feature, we estimate its IDR, which represents the statistical association of the candidate to the given domain corpus, and extrinsic-domain relevance, which reflects the statistical relevance of the candidate to the domain-independent corpus. Only candidates with IDR scores exceeding a predefined intrinsic relevance threshold and EDR scores less than another extrinsic relevance threshold are confirmed as valid opinion features. In short, we identify opinion features that are domain-specific and at the same time not overly generic (domain-independent) via the inter corpus statistics IEDR criterion.

IV. METHODOLOGIES FOR OPINION MINING

An opinion feature such as reviews on a particular product is typically domain-specific. The feature appears frequently in the given review domain, and which are outside the domain is domain-independent corpus about product. 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 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 reflect 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.

A. Candidate Feature Extraction

Opinion features appear as the subject or object of a review sentences 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. It also has a dependence relationship of preposition-object on the prepositional word in the sentence.

B. 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

The system Architecture of proposed system is as shown in the following figure. 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. 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 off line 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 in 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.

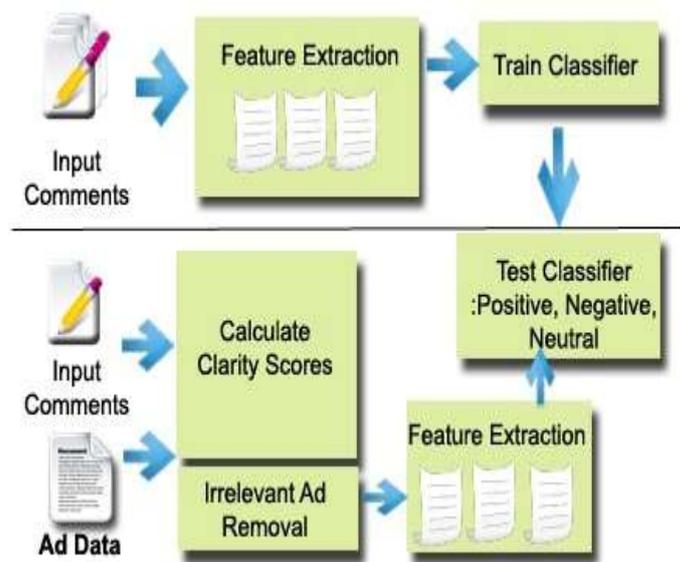


Fig 1. System Architecture

VI. EXPERIMENTAL RESULTS

As the existing system is based on star rating it can be easily attacked. The proposed system is designed to recover the problems present in the existing system as discussed in the section 3. The following graph is showing the comparison between existing system and proposed system.

VII. CONCLUSION

In Future by implementing this system 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 proposed system, 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 web sites in question provide relevant context for opinion mining. The system aware of the fact that has not



Fig 2. Experimental Result Comparison

utilized the information from the social network directly in the opinion mining algorithm. Merely, this 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] Zhen Hai, Kuiyu Chang, Jung-Jae Kim, and Christopher C. Yang "Identifying Features in Opinion Mining via Intrinsic and Extrinsic Domain Relevance" VOL. 26, NO. 3, MARCH 2014.
- [2] W. Jin and H.H. Ho, "A Novel Lexicalized HMM-Based Learning Framework for Web Opinion Mining," Proc. 26th Ann. Int'l Conf. Machine Learning, pp. 465-472, 2009.
- [3] N. Jakob and I. Gurevych, "Extracting Opinion Targets in a Single and Cross-Domain Setting with Conditional Random Fields," Proc. Conf. Empirical Methods in Natural Language Processing, pp. 1035-1045, 2010CP.
- [4] G. Qiu, B. Liu, J. Bu, and C. Chen, "Opinion Word Expansion and Target Extraction through Double Propagation," Computational Linguistics, vol. 37, pp. 9-27, 2011.
- [5] S. M. Kim and E. Hovy, "Extracting Opinions, Opinion Holders, and Topics Expressed in Online News Media Text," Proc. ACL/COLING Workshop Sentiment and Subjectivity in Text, 2006.
- [6] B. Pang and L. Lee, "A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts," Proc. 42nd Ann. Meeting on Assoc. for Computational Linguistics.
- [7] M. Hu and B. Liu, "Mining and Summarizing Customer Reviews," Proc. 10th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 168-177, 2004.
- [8] B. Liu, "Sentiment Analysis and Opinion Mining," Synthesis Lectures on Human Language Technologies, vol. 5, no. 1, pp. 1-167, May 2012.
- [9] Y. Jo and A.H. Oh, "Aspect and Sentiment Unification Model for Online Review Analysis," Proc. Fourth ACM Int'l Conf. Web Search and Data Mining, pp. 815-824, 2011.
- [10] G. Qiu, C. Wang, J. Bu, K. Liu, and C. Chen, "Incorporate the Syntactic Knowledge in Opinion Mining in User-Generated Content," Proc. WWW 2008 Workshop NLP Challenges in the Information Explosion Era, 2008.
- [11] A. Popescu and O. Etzioni, "Extracting Product Features and Opinions from Reviews," Proc. Human Language Technology Conf. and Conf. Empirical Methods in Natural Language Processing, pp. 339- 346, 2005.
- [12] V. Hatzivassiloglou and J.M. Wiebe, "Effects of Adjective Orientation and Gradability on Sentence Subjectivity," Proc. 18th Conf. Computational Linguistics, pp. 299-305, 2000.
- [13] Q. Su, X. Xu, H. Guo, Z. Guo, X. Wu, X. Zhang, B. Swen, and Z. Su, "Hidden Sentiment Association in Chinese Web Opinion Mining," Proc. 17th Int'l Conf. World Wide Web, pp. 959-968, 2008. <http://nll.vnunet.com/news/1116995>. [Accessed: Sept. 12, 2004]. (General Internet site)