



A Novel Technique to Opinion Classification Problem Solving

Pinki Yadav

Department of CSE

M.Tech Student, GITM Guragon

ABSTRACT

Opinions have become increasingly important in today's world for taking decisions. Web has huge amount of unstructured data in the form of blogs, reviews, forums, social networking websites etc. To classify the information present on the web, we need to develop some method that will automatically help in classifying this information. This can be done by sentiment analysis and Opinion mining. Opinion Mining or Sentiment Analysis is a natural language processing task that mine information from various text forums and classify them on the basis of their polarity as positive, negative or neutral. In this paper, a novel methodology is developed that will help in this classification. Review data is collected for various product domains from micro blogging sites like twitter, face book.

Keywords

Opinions, Sentiment Analysis, Data Mining, Word Net, Social Network.

1. INTRODUCTION

Opinions and emotions are very important in human actions. The access to large quantity of data through internet is not an issue as terabytes of new is produced on the web everyday and this is easily available to every individual. People have been tremendously using social media to advertise their product. Today people not only comment on the existing information, bookmark pages and provide ratings but they also share their ideas, news and knowledge with the community at large.

Sentiment mining, also called opinion analysis, analyzes people's opinions and emotions towards entities such as products, services, organizations, individuals, issues, events, topics, and their attributes. [1]. Opinion analysis helps in retrieving this information from the social media sites, which is mostly in the textual format. So, we need to perform textual analysis. Text based analysis helps to generate the polarity of a textual opinion. Opinions can be divided into positive, negative and neutral categories. Our goal is to extract the opinions and to give a final verdict of whether to accept or reject the product by categorizing the opinions into positive, negative and neutral categories. The technical modeling of opinion analysis in social media is developed with the help of senti word net [2].

The content of the paper discusses firstly, the primary concept about sentiment analysis and some previous works. Next, the various applications of opinion mining are

discussed. The overall architecture is briefly discussed as follows.

The methodology of overall system has been developed which presents a clear picture of getting a final decision of accepting or rejecting a product. This is made possible by first extracting reviews from Amazon, twitter and Flip kart and storing them in the repository. Messages posted as blogosphere are mostly expressed as informal text which require more processing as compared to formal text. Informal text consists of sarcasm, poor grammar, and non dictionary standard words [3]. So preprocessing is required to filter the final words. After preprocessing is done, classification of opinions is carried out using Senti word net dictionary. The final verdict can be taken after this categorization.

The remainder of the paper is organized as follows. Section 2 presents the overview of opinion mining. Section 3 discusses the literature review. In section 4, the overall methodology is proposed. Section 5 discusses the results and section 6 concludes.

2. OPINION MINING

Opinions can be defined as a private state of an individual represented in the form of emotions, sentiments, ideas etc [4]. In other words, these can be simply a positive or negative sentiment, emotion, attitude or appraisal about an entity or an aspect of an entity [5].

Mining refers to as extracting some useful information and representing it in the form of knowledge in an understandable structure as patterns, graphs etc [6]. Opinion mining refers to a sub discipline of computational linguistics that focuses on extracting people's opinion from the web [5]. The Opinion Mining lies at the intersection of IR (Information Retrieval) and CL (Computational Linguistics).

Sentiment analysis on the other hand determines the contextual information, polarity (positive, negative or neutral) and polarity strength (weakly positive, mildly positive, strongly positive) of a document [7]. There are three levels in Opinion mining; Document, sentence and aspect (view) level. They basically help in retrieving the public opinion based on the features. The four pairs of human emotions i.e. joy-sadness, acceptance-disgust, Anticipant-Surprise, Fear-Anger [8] helps in the classification process.

3. LITERATURE

REVIEW

Opinion Mining is the technique of detecting and extracting subjective information in text documents [5].

Karkare, V.Y [9] has proposed the system architecture on evaluation of a product considering user opinions with the holder and the topic. It has an advantage that the method developed does not use any additional resources except the initial opinion lexical analyzer. The paper has a certain disadvantage as no analysis has been done on how to collect and extract features from corpus given by people, no methodology has been illustrated as how to extract and rate features from opinion of products.

Szabo, P. [10] has presented the overall model in the application of opinion classification. The advantages of the paper given are that it has discussed the various problems solving application with respect to the sentiment analysis using dynamic coefficients. But the problems may occur when opinions are expressed implicitly. Ambiguity may arise in the opinion classification method discussed in the work.

Farhan Hassan Khan [11] has presented a novel three way classification algorithm for twitter analysis. The results are taken by performing experiments using random tweets collected which are proved mathematically more accurate removing the limitations of sarcasm and sparsity. The proposed comprehensive TOM (Twitter Opinion Mining) framework has overcome the limitations in the previous sentiment analysis papers. An important feature of the proposed architecture is that the raw data is efficiently processed and mistakes, abbreviation and noise is removed before providing input to the classifier. Preprocessing of data is still a problem with respect to time.

Krzysztof Jędrzejewski [12] has discussed the properties of social networks associated with opinion mining. A new opinion classification method has been explored which a variant of semantic orientation and presents the results by testing the algorithm for accuracy on data sets from real world. The future work in this paper has concentrated on the response of opinions in text and improving the performance of the algorithm by creating an active learning strategy.

Anna Stavrianou [13] has presented a model based on opinion based graph which has focused towards content oriented domain. The comparison between the existing user based graph approach and the proposed opinion based graph has been illustrated. The paper has various advantages. The proposed model has given a better technique of handling knowledge extracted from the discussion. The mining of the discussion has reduced the dimension space of the data. The limitation of the paper is that the opinion based graph is temporally dependent. No work has been done to explore how opinion changes over time.

Bing Liu [6] has explained the heuristic and rule based methods discussing the overall description of what opinion mining is, techniques used in sentiment classification and

how opinion summarization can be performed. The limitation of the work is that more work needs to be done on probabilistic methods.

4. PROPOSED WORK

In search engines and in data mining systems, Web crawler serve as the one of the most important module [14]. Generally, the crawler downloads html source code of the web pages, and uses regular expression to find all URLs in the downloaded pages, and use these URLs as seeds for next turn of download. But to retrieve specific information efficiently, we need to perform some panellised form of crawling. We can use some of the specific tags from the popular social networking sites and extract the opinions.

The generic framework of our work is shown in Fig 1.

First, the user posts a query according to his her interest. Then, the query will be preprocessed in order to remove the unwanted content. The query is searched within various social platforms and WWW will be searched for the particular HTML Page and parsing is performed, where Opinion Retriever will extract the relevant opinions and store them in a repository.

Thereafter the opinions collected are passed over to another module. In the Opinion Classification module, the opinions are identified and classified by using Senti word dictionary for the sentiment analysis.

The algorithm used in our work is as follows.:

Input: User query which will be searched from social networking sites.

Online reviews will be extracted.

Output: Sentiment Analysis

Procedure:

Data Preprocessing()

Polarity Identification()

Opinion Orientation()

End

The reviews are collected in the Opinion repository. The data preprocessing is performed. Tasks like Tokenization will split the reviews into each tokens. The part-of-speech tagging will categorize the English grammar in nouns, verbs, adjectives, pronouns, prepositions, conjunctions and interjections. Each word (token) will be labeled with its appropriate part of speech. Thereafter, Stemming works on reducing the root words. And finally, Stop word removal will remove all the unwanted words in each review sentence which can be checked against stop word list. It can be of types like pronouns, prepositions, conjunctions etc.

Polarity Identification is used to categorize the tokens into positive, negative and neutral categorizes with the help of Senti word net.

Opinion Orientation will identify the opinions as to give the final verdict of whether to recommend a product or to reject a product by looking into more of positive or negative tokens.

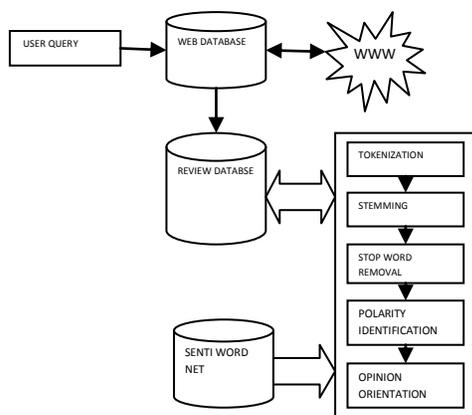


Fig 1: Generic Framework

5. RESULTS

The proposed system uses customer review dataset about a product effectively. The social networking sites like Amazon, Flip-Kart, Twitter are used to extract opinions. Each review is websites is assigned with a different rating like 0-5 stars, a review label and date, a reviewer name and the review content [15]. Sony camera product reviews are used in our system. The query is posted which consist of product name. The review text is retrieved. Reviews are split into individual sentences which are further split into tokens with the sentence id and token id. The details of the dataset used in the proposed system are shown in Table 1 as follows.

Table 1: Corpus

S No.	Corpus	Sony Camera
1	Reviews	90
2	Total sentences	350
3	Positive Sentences	175
4	Negative Sentences	132
5	Total Opinion Sentences	307

Precision, Recall and F-Measure are used for evaluation of our proposed framework and comparison is done[16].

Precision is defined as the ratio of relevant opinions retrieved to the total number of opinions retrieved (relevant and irrelevant opinions retrieved). Mathematically,

$$\text{Precision} = \frac{RTO}{RTO + RWO}$$

Where RTT is the relevant opinions retrieved and RWT is the irrelevant opinions retrieved.

Recall is defined as the ratio of relevant opinions retrieved to the manually retrieved opinions by the classifier (relevant opinions retrieved and relevant opinions not retrieved). Mathematically,

$$\text{Recall} = \frac{RTO}{RTO + RNO}$$

Where RTT is number of relevant opinions retrieved and RNT are relevant opinions not retrieved.

F-Measure is the harmonic mean of both the precision and recall. Mathematically,

$$\text{F-measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

It has been observed that Precision of crawling is high i.e. ranges from 81.26% to 93.7% and Recall of crawling process is also high i.e. ranges from 83.75% to 92.4%.

6. CONCLUSIONS

Online Reviews have become the prime means of communication. In our research, we have discussed the sentiment classification process using Senti Wordnet. A novel methodology is presented with the algorithm in retrieving the opinions and to give a final verdict of recommendation by categorizing the opinions into positive and negative category. The evaluation is done with the information retrieval search strategies and results are presented. Performance of sentiment classification needs to be enhanced by combining different types of techniques in order to overcome their individual drawbacks and benefit from each other's merits. In future, we can apply the supervised and unsupervised learning techniques on the prepared data to compare the individual performances.

REFERENCES

[1] Rabby, F., Masud. M.A., Sayful Islam, G.M., Billah, M., 2013. "Sentiment Mining in Social Network Using Textual Opinion", NCICIT 2013: 1stNational Conference on Intelligent Computing and Information Technology, November 21, CUET, Chittagong-4349, Bangladesh

[2] <http://wordnet-online.freedicts.com/>

[3] Bahrainian, S.-A., Dengel, A., 2013. "Sentimen Analysis and Summarization of Twitter Data", Computational Science and Engineering (CSE), 2013 IEEE

16th International Conference
on DOI: 10.1109/CSE.2013.44, Page(s): 227 – 234, IEEE
Xplore.

[4] K. Khan, B. Baharudin, A. Khan, A. Ullah,
2012, "Mining opinion components from unstructured
reviews: A review", 1319-1578 2014.

[5] Bhatia, S., Sharma, M., Bhatia, K., 2015. " Strategies
for Mining Opinions: A Survey", International Conference
on "Computing for Sustainable Global Development",
IEEE Xplore.

[6] Liu, B., 2012. "Sentiment Analysis and Opinion
Mining", Morgan & Claypool Publishers

[7] Osimo, D., Francesco, M., Anderson, C.,
2008. "Research Challenge on Opinion Mining and
Sentiment Analysis", Wired Magazine, 16(7), 16–07.

[8] Kamath, Bagalkotkar, S.S., Kandelwal, A., Pandey, A.,
Poornima, S., 2013. " Sentiment Analysis Based
Approaches for Understanding User Context in Web
Content", Communication Systems and Network Technologi
es (CSNT), International Conference
on DOI: 10.1109/CSNT.130, Page(s): 607 – 611, IEEE
Xplore.

[9] Karkare, V.Y., Gupta, S.R., 2014. "Product
Evaluation Using Mining and Rating Opinions of Product
Features", Electronic Systems, Signal Processing and
Computing Technologies (ICFSC). 2014 International
Conference. DOI: 10.1109/ICFSC.2014.72 Publication
Year: Page(s): 382 – 385, IEEE Xplore.

[10] Szabo, P., Machova, K., 2012. "Various approaches to
the opinion classification problems solving", Applied
Machine Intelligence and Informatics (SAMI), 2012 IEEE
10th International Symposium

on DOI: 10.1109/SAMI.2012.6208929, Page(s): 59 - 62,
IEEE Xplore.

[11] Khan, F.H., Bashir, S., Qamar, U., 2014. "TOM:
Twitter opinion mining framework using hybrid
classification scheme", Decision Support Systems, Volume
57, January, Pages 245-257, Science Direct.

[12] Jędrzejewski, K., Morzy, M., 2011. "Opinion Mining
and Social Networks : A Promising Match", International
Conference on Advances in Social Networks Analysis and
Mining.

[13] Stavrianou, A., Velcin J., Chauchat J-H., 2009. "A
combination of opinion mining and social network
techniques for discussion analysis".

[14] Gao, Y., Peng, C., 2013. "Design and Implementation
of Distributed Crawler System for Opinion Mining",
Proceedings of the 2nd International Conference on
Computer Science and Electronics Engineering (ICCSEE
2013).

[15] Jeyapriya, A., and Kanimozhi Selvi., 2015. "Extracting
aspects and mining opinions in product reviews using
supervised learning algorithm." Electronics and
Communication Systems (ICECS), 2015 2nd International
Conference on. IEEE.

[16] Bhatia, S., Sharma, M., Bhatia, K.K., 2015. "Sentiment
Knowledge Discovery using Machine Learning
Algorithms", Journal of Network Communications and
Emerging Technologies (JNCET) Volume 5, Special Issue
2, December (2015).