

TOWARDS AUTOMATIC DETECTION OF SENTIMENTS IN CUSTOMER REVIEWS

ANSHU JAIN¹, DR. SURESH JAIN², MRS. PRAGYA SHUKLA³, HITESH BANDIYA⁴

¹DEPARTMENT OF COMPUTER ENGINEERING, INSTITUTE OF ENGINEERING & TECHNOLOGY
INDORE, INDIA

1jainanshu.85@gmail.com

³pragyashukla_iet@yahoo.co.in

4hiteshbandiya.iet@gmail.com

²Department Of Computer Science, KCB Technical Academy Indore, India

²suresh.jain@rediffmail.com

ABSTRACT

Opinions Play important role in the process of knowledge discovery or information retrieval and can be considered as a sub discipline of Data Mining. A major interest has been received towards the automatic extraction of human opinions from web documents. The sole purpose of Sentiment Analysis is to facilitate online consumers in decision making process of purchasing new products. Opinion Mining deals with searching of sentiments that are expressed by Individuals through on-line reviews, surveys, feedback, personal blogs etc. With the vast increase in the utilization of Internet in today's era a similar increase has been seen in the use of blog's, reviews etc. The person who actually uses these reviews or blog's is mostly a consumer or a manufacturer. As most of the customers of the world are buying & selling product on-line so it becomes company's responsibility to make their product updated. In the current scenario companies are taking product reviews from the customers and on the basis of product reviews they are able to know in which they are lacking or strong this can be accomplished with the help of sentiment analysis. Therefore Our objective of our research is to build a tool which can automatically extract opinion words and find out their polarity by using dictionary, This actually reduces the manual effort of reading these reviews and to evaluate them. The research also illustrates the benefits of using Unstructured text instead of training data which expensive. In this research effort we demonstrate a method which is based on rules where product reviews are extracted from review containing sites and analysis is done, so that a person may know whether a particular product review is positive or negative or neutral. The system will utilize a existing knowledge base for calculate positive and negative scores and on the basis of that decide whether a product is recommended or not. The system will evaluate the utility of Lexical resources over the training data.

KEYWORDS

Opinion Mining, Sentiment Analysis, Sentiment Orientation, Opinion Analysis, Polarity Analysis, Subjectivity Detection, Review Mining

1. INTRODUCTION

Opinion mining deals with the sentimental words which clearly demonstrates ideas about the targeted object at that particular time. It is a type of natural language processing for gathering the information from public about a particular product. Whenever we decide to buy something new,

we always look forward for someone's advice for same. This field has wide range of application due to its increasing popularity among the public. for example, it can help you judge that which product is good and which is bad while both are having the same price, market analysts may use the reviews posted in various blogs for taking major decisions regarding new product launch. Opinion mining actually identifies the author's viewpoint about a subject, rather than simply identifying the subject itself [17]. Opinions are the views that are expressed by an individual on some topic/issue according to his/her own perspective. As these views have been used by several application domains such as business and organization, individual etc., we can say it has become very important to find out efficient ways for extracting opinions. So Opinion mining is the process of studying people's opinions or emotions towards entities, events and their features. In the past few years, it has attracted a great deal of attentions from both academia and industry due to many challenging research problems and a wide range of applications. The another name for opinion mining may be sentiment analysis and subjective analysis. The objective of opinion analysis is to identify emotions from text and determine their polarities. On the basis of determined polarities we can conclude whether the text document represents positive, negative or neutral opinion. The discovered opinions are useful to many practical applications such as opinions in product quality reviews are helpful to potential customers. Users also comment on products in their personal web sites and blogs, which are then aggregated by sites such as Blogstreet.com, AllConsuming.net, and onfocus.com [18]. Meanwhile, opinion mining technique relates indirectly to promote many natural language processing techniques. The sudden growth in the area of opinion mining, which deals with the computational techniques for opinion extraction and understanding created an utmost need to understand and view the internet in a different prospect. In this tool we will use product review extractor interface which extracts the reviews about the various products from system.

1. The tool will summarize the overall polarity of the current review which has been extracted from the web page previously and are stored inside the system and output will generate overall opinion that the review indicates as positive, negative or neutral.
2. The summarization of given review and its features as positive or negative is mainly emphasized which reduces the overall burden & time of manually reading a review and finding out what is good or bad.

Remaining part of the paper has been organized in the following manner:-part 2 describes the necessity for conducting this research .part 3 presents the necessity of proposed automated system, part 4 gives an overview of proposed Technique for opinion summarization, part 5 describes the system design ,part 6 presents the result after conducting experiments with the proposed system and part 7 concludes this paper.

2. NECESSITY

The interest that individual users show in on-line opinions about products and services, and the potential influence is something that vendors of these items are paying more and more attention to. Modern users heavily rely on on-line information, advice and recommendations. This requires a quick, automated and intelligent process to extract quality opinions. Such systems have to cater the needs of on-line information processing along with the provision for missing, confusing, and overwhelming information. Thus, there is a clear need to aid consumers of products and of information by building better information-access systems than are currently in existence. From a few years ago, sentiment mining has become a hot subject among NLP and IR researchers. Large amount of effort has been put into the research on this field that quite a number of papers have

been published and systems for various applications using sentiment mining have been developed and put into commercial use [19] .

The purpose of the project is to develop a tool which will be useful to product manufacturer and customers in reducing their burden of manually reading the reviews and finding out which feature is good or bad. There are two main types (i.e., facts and opinions of information on the web. However, current search engines (e.g., google) are all for facts expressed with topic keywords. There are many factors that make Opinion Mining difficult compared to traditional fact-based text analysis [20].

3. INVOLVED RESOURCES

In this section we describe the necessary resources involved in summarizing the opinions present in product reviews. As our work is concentrated on assigning polarity score to various feature words which are extracted we need to use some existing lexical resource which must be effective and can smartly extract positive and negative score for extracted words. Here comes the utility of SentiWordnet given By A. Esuli and F. Sebastiani,[10],It is described as below:-

(a) SentiWordnet:-it is a lexical resource of sentiment information for terms in the English language designed to assist in opinion mining tasks,where each term is associated with numerical scores for positive and negative sentiment information. It provides a readily available database of term sentiment information for the English language. Each term in SentiWordNet is associated with numerical scores for positive and negative sentiment information.SentiWordNet can be a replacement to the process of manually deriving lists of terms containing sentiment information for opinion mining tasks.

Another resource which is utilized in this research is POS Tagger for part-of-speech tagging purposes which is freely provided by Stanford university for research purposes. This Described as below:-

(b) Stanford POS Tagger:-The Process of assigning different parts of speech tags such as noun,adjective and adverbs to a given text is known as Part-Of-Speech tagging. A Part-Of-Speech Tagger (POS Tagger) is a piece of software that reads text in some language and assigns parts of speech to each word (and other token), such as noun, verb, adjective, etc., although generally computational applications use more fine-grained POS tags like 'noun-plural'. This software is a Java implementation of the log-linear part-of speech taggers. The system requires Java 1.5+ to be installed.

The English taggers use the Penn Tree bank tag set. Here are some tags of the Penn Tree bank English POS tag set:-

Sno.	Tag	Description
1	NNP	noun, proper, singular
2	NNPS	noun, common, plural
3	RB	adverb
4	JJ	adjective or numeral, ordinal
5	JJR	adjective, superlative
6	NN	noun, common, singular
7	RBR	adverb, comparative
8	VB	verb, base form
9	VBD	verb, present participle or gerund
10	WDT	WH-determiner
11	CC	conjunction,coordinating
12	CD	numeral, cardinal
13	DT	determiner

Table 1 Tag set Description

4.PROPOSED TECHNIQUE

Major Tasks

(i) Data Collection:- In this step necessary data which are products reviews here are required to be collected from opinion oriented sites.

(ii) Data Preparation:-Before applying the actual algorithm for the task of opinion mining in any text document which are product reviews here we require some process in order to format the data set as required by classification algorithm.

(iii) Feature extraction :- The words which typically express the properties of the object.

(iv) Sentiment Classification:- This task determines whether there is opinion on a feature in sentence, and if so, whether it is positive or negative. Existing approaches are based on supervised and semi-supervised methods.

(v)Integration:-Finally we merge all the above tasks.

In this section we describe our proposed algorithm for the effective extraction of useful information from product reviews. This algorithm works on the basis of rule based methodology using lexicon approaches as described in the above sections. For better accomplishment of the steps we need to take some assumptions in advance .so we consider that the product reviews from various review relevant sites have already been extracted with the help of appropriate procedures

and are already stored inside the system. Therefore we skip the task of extracting the reviews and storing them in our system which requires some sort of crawlers to get it done. Then after we take a set of product reviews as input and consider a polarity score along with the relevant feature words as expected outcome. Following are the number of steps required to generate the expected outcome:-

Proposed Algorithm

INPUT:Set of product reviews

OUTPUT:Extracts opinion words with positive and negative polarity

BEGIN

- 1) Collect the reviews R_n for various product P_n
- 2) Each review R_i contains only one sentence in one line
- 3) (a) For each Review R_i in
 - For each Sentence S_i in R
 - For each word W_i in Sentence S_i
 - (b) Check whether the word is in stop word list e.g. A,and,the,for,at etc.
If found then replace it with white space
 - (c) Assign P_no and POS tag to each word.
 - (d) Extract the patterns which matches predefined rules.
 - (e) Assign positive and negative score to each word in extracted pattern. Using Function `assign_score()`
If the pos score>neg score then
Fea_Word score=pos score
Else if pos score<neg score then
Fea_Word score=neg score
- (4) Repeat step 3 (a) and check whether the word matches any word from Conditional dictionary.
Eg negative words such as not,never,neither etc
If found then recalculate its score
Fea_Word score=Fea_Word score+Cond_weight
- 5) Sentence score will be calculated using the

$$i=1$$

$$\text{Sent_Score} = \sum_{n=1}^n \text{Fea_Wordscore}(1) + \text{Fea_Wordscore}(2) + \dots + \text{Fea_Word score}(n)$$
- (6) Finally overall review score will be

$$i=1$$

$$\text{Review_score} = \sum_{n=1}^n \frac{\text{Sent_Score}(n)}{n}$$
- 7) Calculate Accuracy using recall ,precision and F-measure

$$\text{Precision[Positive]} = \text{True Positive} / (\text{True Positives} + \text{False Positives})$$

$$\text{Recall[Positive]} = \text{True Positives} / (\text{True Positives} + \text{False Negatives})$$

Precision-Precision is the fraction of the documents retrieved that are [relevant](#) to the user's information need.

Recall-Recall is the fraction of the documents that are relevant to the query that are successfully retrieved.

F-measure- it is the harmonic mean of precision and recall. It is computed using the formula defined below.

$$F - \text{measure} = 2 * (\text{precision} * \text{recall} / (\text{precision} + \text{recall}))$$

As mentioned in first step of algorithm we retrieve previously extracted reviews and provide them as input to our system. In the next step for accomplishing sentence level sentiment classification we process each sentence in a individual review and then each word in a individual sentence for further processing. Third step comprises of loop control structure for review,sentence and word respectively. Starting at word level we check whether the any word matches the word in stop word list and if does then replace it with white space .The reason behind removing these words is that they contain very less probability of holding any kind of opinion information. After this within the same step we assign position to each word and perform parts of speech tagging of remaining words with the help of existing resources as mentioned in section 3.

In step 3 (d) we extract useful patterns which are the collection of most probable words those are considered to be good polarity indicators. So the useful patterns which we used for our methodology are as below:-

Pattern	First word	Second Word	Third word
P1	Adjective(good)	Noun(Image)	Noun(Quality)
P2	Verb(Recommended)	Noun(camera)	-----
P3	Adverb(very)	Adjective(short)	Noun(life)
P4	Adverb(not)	Adverb(very)	Adjective(good)

Table 2 Description of predefined Pattern rules

After successful extraction of predefined rules the resulting words will be assigned to the positive and negative score with help of knowledge dictionary and Function assign_score() will decide whether a word is assigned with positive or negative score. Each word in the pattern will get the positive or negative score. Next step 4 again performs the matching each word in input review with the conditional dictionary if found then recalculation of polarity score is done .Conditional dictionary contains contextual word and due their presence polarity will get affected so we perform dual checking to detect these words. Rest of the steps 5,6,7 Contains the formula for calculating scores for a word,sentence and review respectively.

5. SYSTEM DESIGN

Figure 1 represents our proposed system design,where task of opinion mining is divided in to two parts, one is data preprocessing and other is data main processing. In data preprocessing we firstly remove stop words from input text so as to be more focused on required data and then parts of speech tagging for further processing of remaining text. In data main processing we firstly extract the predefined phrases and then polarity score is assigned to them .Nextly we evaluate our technique By calculating recall and precision over obtained results.

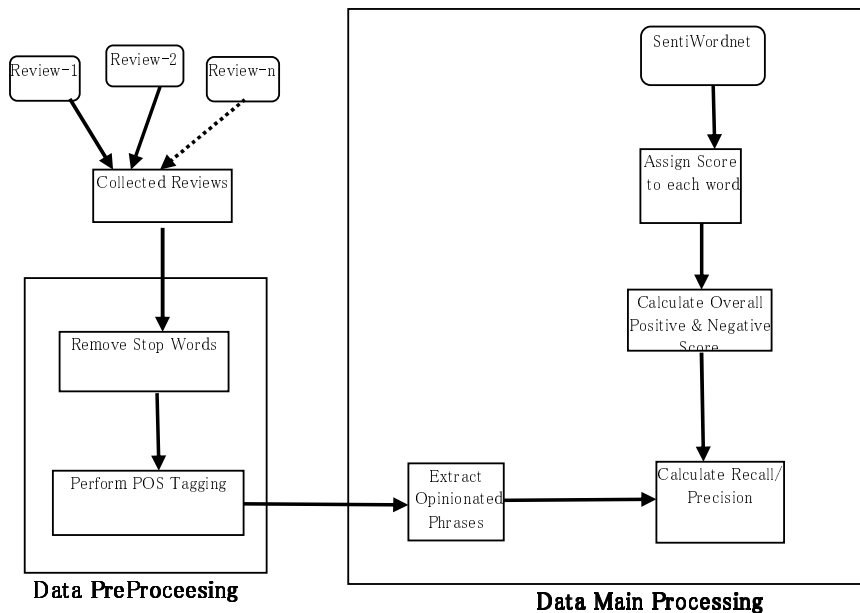


Figure 1 Proposed System design

The tool named as Opminer Works in two Modules which are as follows:-

Data Preprocessing

We processed only the description part of each review, split the review in to sentences and assigning position to each words.

- (i) Removing stop words:-We remove words like digits, prepositions, articles etc,for this create a list of words, which are to remove from the review. It helps better extraction of opinion phrases/words.
- (ii) Not all the words in review sentences are useful for identifying product features and orientations of the discussed product. As Hu and Liu point out, nouns and noun phrases in the sentences are likely to be the features that customers comment on, while adjectives are often used to express opinions and feelings.

Perform Part-Of-Speech (POS) tagging with the help of Stanford POS tagger to assign POS tag to each word (whether the word is a noun, verb, adjective,etc). POS tagging is important as it allows us to generate general language patterns.

Data Main Processing

- (i) Consider the noun phrases as feature words in a sentence and adjective phrases as product features ,this replacement is necessary because different products have different features. The replacement ensures that we can find general language patterns which can be used for any product.
- (ii) Now we Implement a matching algorithm which will extract the positive/negative score of product features extracted in the previous step with the help of knowledge base named as SentiWordnet.

6. RESULTS

We practically conduct experiment in order to test our proposed technique. For this we collected variety of reviews about various products from review containing sites such as Amazon.com, Eopinions.com etc and store them inside our system. After this We input them in to our system for further processing and then we calculate recall & precision with the help above mentioned formulas. we have achieved 68% precision and 56% recall as the proposed model is in development phase.

7. CONCLUSION

Opinion Mining is a emerging technology applying to but yet need get improvement to achieve more convenient and it is emerging and rapidly growing field .however the previous research work mainly focused on training dataset for analysing the sentiments present in user generated content, they succeed in providing good level of accuracy but due to lack of availability of annotated corpus for same purposes, we emphasis on the utility of lexical resources .In this paper we showed the usability of these sources over supervised learning . Opinion Mining has become Important for all types of fields in extracting likes and dislikes and intensity of likes and dislikes. So the research conducted here is hopefully of significant use to help the researchers in knowing about the past and recent trends in the same area. Opinion Mining faces many problems and issues as still it has ambiguity of the kind that a word which is positive in one context may be negative in another context, secondly identification of opinionated text from the given document is itself a tough task, nextly identification of positive or negative opinion words and overall sentiment in a document is also a challenging. Opinionated text may be fake, irrelevant or may be ambiguous. So, the research scope of this emerging field is very big and vast and needs to be explored properly for effective opinion mining systems.

REFERENCES

- [1] B. Liu. "Sentiment Analysis and Subjectivity". Handbook of Natural Language Processing, Second Edition, (editors: N. Indurkha and F. J. Damerau), 2010.
- [2] B. Liu, "Sentiment Analysis: A Multi-Faceted Problem IEEE Intelligent Systems", 2010
- [3] B. Liu, "Opinion Mining", invited contribution to Encyclopedia of Database Systems, 2008
- [4] Hu and Bing Liu, "Mining Opinion Features in Customer Reviews, Mining", American Association for Artificial Intelligence 2004.
- [5] Bo Pang and Lillian Lee, "Opinion Mining and Sentiment Analysis", Foundations and Trends in Information Retrieval Vol. 2, Nos. 1-2 (2008) 1-135 2008.
- [6] Mingqing Hu and Bing Liu, "Mining and Summarizing Customer Reviews", KDD'04, August 22-25, 2004, Seattle, Washington, USA, 2004.
- [7] Ruifeng Xu^{1,2}, Kam-Fai Wong², Qin Lu¹, Yunqing Xia³, Wenjie Li¹ "Learning Knowledge from Relevant Webpage for Opinion Analysis", 2008 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, 2008
- [8] A. Esuli, "Opinion Mining", Language and Intelligence Reading Group, June 14, 2006, Pisa, Italy. Martin, J. R. and White, P. R. R. (2005). The Language of Evaluation: Appraisal in English.
- [9] A. Esuli and F. Sebastiani, "SentiWordNet: A publicly available lexical resource for opinion mining," in Proceedings of Language Resources and Evaluation (LREC), 2006.
- [10] A. Esuli and F. Sebastiani, "PageRanking WordNet ,synsets: An application to opinion mining," .
- [11] Q. L. Miao, Q. D. Li, and R.W. Dai, A sentiment mining and retrieval system. Expert Systems with Applications: An International Journal, 2009, 36 (3): 7192-7198.
- [12] H. Yu, and V. Hatzivassiloglou, Towards answering opinion questions: Separating facts from opinions and identifying the polarity of opinion sentences. In: Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing. Morristown, NJ, USA: Association for Computational Linguistics, 2003, pp. 129 - 136.

- [13] Xiaojun Li, Lin Dai, Hanxiao Shi, "Opinion Mining of Camera Reviews Based on Semantic Role Labeling", 2010 Seventh International Conference on Fuzzy Systems and Knowledge Discovery (FSKD 2010).
- [14] Turney, P. (2002). Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews. In Proceedings of ACL-02, 40th Annual Meeting of the Association for Computational Linguistics, pages 417–424. Association for Computational Linguistics.
- [15] Turney, P. D. and Littman, M. L. (2003). Measuring praise and criticism: Inference of semantic orientation from association. ACM Transactions on Information Systems, 21(4):315–346.
- [16] Liu, Sentiment Analysis: A Multi-Faceted Problem IEEE Intelligent Systems, 2010.
- [17] Jack G. Conrad, Frank Schilder, "Opinion mining in legal blogs", In Proceedings of the 11th International Conference on Artificial Intelligence and Law (ICAIL07).
- [18] Kushal Dave, Steve Lawrence, David M. Pennock, "Mining the Peanut Gallery: Opinion Extraction and Semantic Classification of Product Reviews.
- [19] Gong Tianxia, "Processing Sentiments and Opinions in Text: A Survey".
- [20] Zhang Wei, "Opinion Mining and Sentiment Analysis: A Survey".

Authors

Dr. Suresh Jain received the B.E. Degree in Civil Engineering, from Maulana Azad National Institute of Technology, Bhopal, India in May 1986, M.E. in Computer Engineering from Shri Govindram Sakseria Institute of Technology & Science, Indore, India in April. 1988 and the Ph.D. Degree in Computer Science from Devi Ahilya University Indore, India in March 2007. He is presently the Professor & Dean Academics, K.C.B. Technical Academy, Indore India. since April 2008, he was a Professor, Institute of Engineering & Technology, Devi Ahilya University, Indore (August 2007- March 2008), Reader Institute of Engineering & Technology, Devi Ahilya University, Indore (May. 2000- August 2007). His research interests are in Machine Learning, Grammatical Inference, Artificial Intelligence, Graphics and Multimedia, Database Applications. He has published a book: Introduction to Programming in Pascal - II, Nakoda Publisher, 2000. He is a Senior Member of the Computer Society of India and ISTE societies and Life Member of the Indian Society for Technical Education.



Ms. Pragya Shukla has received her B.E. In Computer Engineering From Government Engineering College, Bhopal in the year 1996. At present she is working as Associate Professor in department of Computer Engineering, Institute of Engineering and Technology DAVV Indore, India since 1998. Pursuing her Ph.D. under the faculty of Engineering, Devi Ahilya University, Indore, India. She has published nearly 9 research papers in National/International Journals. Her Areas of Artificial Intelligence, Discrete Structures, DBMS, Computer Architecture. Her Areas of Artificial Intelligence, Discrete Structures, DBMS, Computer Architecture.



Anshu Jain received her B.E degree in Computer Science and Engineering from Maharana Pratap college of Technology, Gwalior India in 2007 and Pursuing M.E (Computer Engineering with sp. Software Engineering) from Institute of Engineering & Technology, Indore, India She has published nearly 4 research papers in National/International Journals. Her Primary interests include Opinion Mining And Information Retrieval.



Hitesh Bandiya received his B.E degree in Computer Science and Engineering from Mandsaur Institute of Technology, Mandsaur India in 2005 and Pursuing M.E (Computer Engineering with sp. Software Engineering) from Institute of Engineering & Technology, DAVV, Indore, India He has published nearly 4 research papers in National/International Journals. His Primary interests include Information Retrieval and Knowledge Discovery.

