

# Extracting the opinion targets and opinion words by using modified PSWAM with sentimental analysis

Harshala R Patil<sup>1</sup>, Prof. R. B. Wagh<sup>2</sup>

<sup>1</sup>RCPIT, Shirpur

<sup>2</sup>RCPIT, Shirpur

(E-mail: harshalapatil158@gmail.com, rajnikantw@gmail.com)

*Abstract*— As the trend of online shopping is growing rapidly, the peoples leaned towards online shopping. Many application asks for a shopping experience that is review. As a result, a large number of reviews are available for a single product. The Internet provides the biggest forum to express an opinion about any product whether it is good or bad. By overlooking these reviews the manufacturer gets the idea about the product improvement in a timely fashion. The retailer can get an idea about which products are trending. And end user can read review status before purchasing so that he or she can understand the product is good or bad. The important task is to find opinion relation between words in the sentence. To extract opinion targets, opinion words and identifying the relations between them as an alignment process the modified partially-supervised word alignment model (PSWAM) is used. The TF-IDF score of possible opinion target and words are calculated. Then, the confidence of each candidate calculated and the candidates with higher confidence will be extracted as the opinion targets or opinion words. This model captures opinion relations more precisely, especially for long span relations as compared to previous methods based on the nearest-neighbor rules. Sentimental analysis is used to find sentiment about the product they purchased. The word of the bag is used to find out the negative, positive orientation of the sentence. Because of the usage of partial supervision, the proposed model obtained a better result as compared to the unsupervised alignment model.

**Keywords**— *Opinion Mining, Opinion Words, Opinion Targets, TF-IDF.*

## I. INTRODUCTION

Now-a-days e-commerce technology becomes more popular, the convenience of online shopping has attracted more and more people. In order to get product feedback timely and to update the future customer with other's shopping experiences of the same product, it is common for merchants to allow their customers to leave product reviews. As the number of customers increases, the number of reviews about the product grows as well. So mining of reviews become an important process. Using opinion mining customer, who wants to buy the product, can have an idea about the product's quality and the manufacturer can improve the product on time. With the rapid expansion of e-commerce, more and more products are sold on the Web and more and more people are

buying products on the Web. In order to enhance customer satisfaction and their shopping experiences, it has become a common practice for online merchants to enable their customers to review or to express opinions on the products that they buy. With more and more common users becoming comfortable with the Internet, an increasing number of people are writing reviews. As a consequence, the number of reviews that a product receives grows rapidly. Some popular products can get hundreds of reviews at some large merchant sites. It makes it very hard for a potential customer to read them to help him or her to make a decision on whether to buy the product.

Recently, the number of online shopping customers have dramatically increased due to the rapid growth of e-commerce, and the increase of online merchants. With the rapid expansion of e-commerce, more and more products are sold on the Web and more and more people are buying products on the Web. To enhance customer satisfaction, it has become a common practice for online merchants and product manufacturers to allow customers to review or express their opinions on the products or services they using. The customers can now post a review of products at merchant sites, e.g., amazon.com, cnet.com, and epinions.com. With more and more common users becoming comfortable with the Internet, an increasing number of people are writing reviews. As a result, more and more reviews are increasing on the web about products. It makes difficult for the potential to read all reviews. These online customer reviews, thereafter, become a cognitive source of information which is very useful for both potential customers and product manufacturers. Customers have utilized this piece of this information to support their decision on whether to purchase the product. For product manufacturer perspective, understanding the preferences of customers is highly valuable for product development, marketing and consumer relationship management. Since customer feedbacks influence other customer's decision, the review documents have become an important source of information for business organizations to take its development plans.

There are two main types of textual information facts and opinions, a major portion of current information processes methods such as web search and text mining work with the former. Opinion Mining refers to the broad area of natural language processing, computational

linguistics and text mining involving the computational study of opinions, sentiments and emotions expressed in the text. A thought, view, or attitude based on emotion instead of the reason is often referred to as a sentiment. Hence, an alternate term for Opinion Mining is Sentiment Analysis. This field finds critical use in areas where organizations or individuals wish to know the general sentiment associated with a particular entity - be it a product, person, public policy, movie or even an institution. The opinion mining has many application domains that include science and technology, entertainment, education, politics, marketing, accounting, law, research and development. With the tremendous growth of the World Wide Web and e-commerce, huge volumes of opinionated texts in the form of blogs, reviews, discussion groups and forums are available for analysis making the Web the fastest, most comprehensive and easily accessible medium for sentiment analysis for researchers. However, finding opinion sources and monitoring them over the Web can be a tough task because a large number of diverse sources exist on the Web and each source also contains a huge volume of information. From a human's perspective, it is both difficult and tiresome to find relevant sources, extract pertinent sentences, read them, summarize them and organize them into a usable form. An automated and faster opinion mining and summarizing system are thus needed.

The work is partly based on and closely related to opinion mining efficiently and sentence sentiment classification. Much research has been done on sentiment analysis of review text and subjectivity analysis i.e. determining whether a sentence is subjective or objective. Another related field of research is feature/topic-based sentiment analysis, in which opinions on particular attributes of a product are determined. Most of this work concentrates on finding the sentiment associated with a sentence (and in some cases, the entire review). Though there has been some work in review summarization, and assigning summary scores to products based on customer reviews, there has been relatively little work on ranking products using customer reviews.

To make application easy to use for the end user, many techniques were used and introduced by researchers. In previous methods, the most adapted technique was a nearest-neighbor rule and syntactic patterns. The nearest neighbor rules align the nearest adjective/verb to a noun/noun phrase in a limited window as its modifier. Clearly, this strategy can not obtain effective and precise results because there exist long and different opinion expressions. Then several heuristic syntactic patterns were designed. However, online reviews generally have informal writing styles that include grammatical errors, typographical errors, and punctuation errors. This makes the existing parsing tools prone to generating errors, which are usually trained on formal formats of texts such as news reports. The standard WAM is trained in a completely unsupervised manner which leads to producing unsatisfactory alignment of words in a sentence. It is not able to give the result precisely. And completely supervised WAM is impossible to implement practically.

An opinion target is word or object occurred in a sentence about which customer expresses their opinion in their reviews, it can be the noun or noun phrase. The opinion words are the words that is used to modify an opinion target in a sentence, it can be an adjective or a verb. In opinion mining extracting opinion words and targets are two fundamental tasks these subtasks are also known as product feature extraction. Product feature extraction can provide the essential information for obtaining fine-grained analysis on customer review. Thus, it has obtained lots of attention in marketing, selling sector. For Example,

“This Mobile has a good and clear screen.”

In the above example, “good” and “clear” are usually used to describe “screen”, so that there are opinion relations between them. If we know that “clear” is opinion word then “screen” is supposed to be an opinion target in this domain. Further, Opinion Target “screen” is used to find out that “big” is most likely an opinion word. The extraction is performed alternatively between opinion words and targets until there is no item left to extract. Then, a constrained POS Tagger is used to find possible opinion targets and possible opinion words in the provided sentence. The opinion relation graph is drawn. A random walk based co-ranking algorithm is performed to calculate the candidate confidence. The TF-IDF score of opinion target and opinion word is calculated. Then, Candidates with higher confidence are extracted as opinion words for opinion target, while calculating confidence TF-IDF score is considered. Stanford POS tagger is used in natural language processing. And at the same time, we performed a sentimental analysis. Reviews in the selected category are divided into Positive, Negative.

## II. RELATED WORK

The process of extracting opinion target and opinion word is not new tasks in opinion mining but the user wants to know about the opinion about the product so that the user can decide it is feasible or not buy it. There are significant efforts focused on all these tasks.

Kang Liu, Liheng Xu, and Jun Zhao [1] have proposed the complex partially supervised word alignment model called the “IBM-3 model”. To obtain the optimal alignments in sentences, an EM-based algorithm is adopted to train the model as it is partially supervised. In this proposed system, to calculate confidence possible opinion targets and opinion words a random walk based algorithm was used. This model has a good ability to detect opinion relations between words, which leads to more effective opinion word and opinion target extraction than previous methods. The focus is mainly on finding opinion words and opinion targets and detecting the relations among them.

L.Zang, B.Liu, S.H.Lim, and E.O'Brien-Strain [2] have proposed the method that uses a ranking algorithm which is based on the web page called HITS. The experiments on diverse real-time datasets were performed. In this method, the feature ranking and

feature extraction are the two fundamental tasks that are introduced to deal with the problems of extracting the opinion reviews. In this case feature ranking is applied to each extracted feature candidate. The feature importance is determined by two factors – feature relevance and feature frequency. The HIT algorithm is specially used for finding feature importance and rank them high.

Minqing Hu and Bing Lu [3], aim for mining and summarizing all the reviews given by the customer. The customer reviews are collected, mined and feature based summary is provided. The main focus is on mining the large dataset of customer reviews and collecting the features of the products. This mining and summarizing the review is based on the reviews of the user as a negative review opinion or positive review opinion. The main concern is with the Positive and the negative review orientation of the review written by the customer, which is based on the adjective word or seed used by the customer to define that product. Here the part-of-speech Tagging technique is used to align the words. The huge number of customer reviews dataset provided.

Fangato Li, Chao Han, et al. [4] have proposed the method that is based on feature-based summarization of reviews. They introduced a new machine learning framework which is based on conditional random fields. This is the new method for co-extracting the sentiments and also topic lexicons. The algorithm such as Relational Adaptive bootstrapping (RAP) is used to expand the seeds in the target domain in the corpus. The twofold effective framework was seen that is topic-lexicon co-extraction and sentimental analysis. The framework can employ an effective rich feature and also extract object feature, Positive opinion and Negative opinion.

Ana-Maria Popescu and O. Etzioni [5] has developed the model that identifies the corresponding customer opinion to determine their sentiment polarity. The relaxation labeling technique is proposed, it mainly focuses on the extraction of features and identifying the customer opinions about the extracted feature and then it is used for deciding the sentiment polarity. Here, OPINE is introduced which is an unsupervised information extraction system. The purpose of OPINE is to mine and build a model of important features of products, evaluation by reviewers and relative quality across the product. The explicit features are required to parse the customer reviews information.

Shahzad Qaiser, Ramsha Ali [6] designed the system, examining the relevance of key-words to documents in corpus. The study is focused on how the algorithm can be applied on number of documents. First, the working principle and steps which should be followed for implementation of TF-IDF are elaborated. Secondly, in order to verify the findings from executing the algorithm, results are presented, then strengths and weaknesses of TD-IDF algorithm are compared. Finally, the work is summarized and the future research directions are discussed. Finally the TF-IDF on which this study is focused on. One can sort the output of algorithm either in

ascending order or descending order based on their occurrences or their TF-IDF score so that the keywords having greater occurrences or greater TF-IDF score would come on the top in decreasing order or the keywords having lower occurrences or lower TF-IDF score would come on the top in increasing order. That can really help in analyzing or slicing the data to generate reports or visualizations. The program can be executed with minimum, a few microseconds time to a few seconds or a minute, depending on the size of the provided dataset.

Zheng Lu, Weiyao Lin[7] has proposed to infer user search goals for a query by clustering its feedback sessions represented by pseudo-documents. First, introduce feedback sessions to be analyzed to infer user search goals rather than search results or clicked URLs. Both the clicked URLs and the unclicked ones before the last click were considered as user implicit feedbacks and taken into account to construct feedback sessions. Therefore, feedback sessions can reflect user information needs more efficiently. Second, they map feedback sessions to pseudo documents to approximate goal texts in user minds. The pseudo-documents can enrich the URLs with additional textual contents including the titles and snippets. Based on these pseudo-documents, user search goals can then be discovered and depicted with some keywords. Finally, a new criterion CAP (Classified Average Precision) is formulated to evaluate the performance of user search goal inference. Experimental results on user click-through logs from a commercial search engine demonstrate the effectiveness of our proposed methods. The complexity of our approach is low and our approach can be used in reality easily. For each query, the running time depends on the number of feedback sessions.

### III. METHODOLOGY

The “Fig. 1” defines the process flow of the system. It extracts the possible opinion targets and possible opinion words from the input sequence. The system requires the input as any review statement. It also identifies the relation between them. After that estimate the confidences of each candidate. Extract the candidate among them which have higher confidence.

#### A. The POS Tagger

The Stanford POS tagger is applied to find out the types of words in the sentence. The classifier classifies the words in the provided sentence. The abbreviations are used to define a word, for example, for “Noun” it used NN. The accuracy of the POS tagger is 90%. If a sentence contains the ambiguity in any form, the POS tagger is not able to identify that so that it can't resolve.

#### B. The TF-IDF Algorithm

Tf-idf stands for *Term* frequency-inverse document frequency.

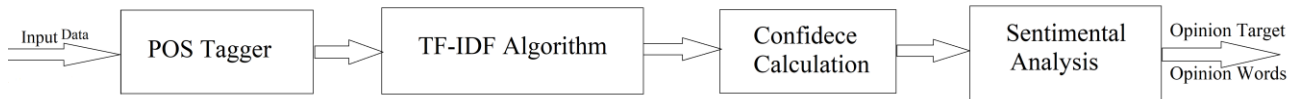


Fig : Modified PSWAM's Process Flow

The importance increases proportionally to the number of times a word appears in the document but is offset by the frequency of the word in the corpus (dataset).

Tf-Idf is a weighting scheme that assigns each term in a document a weight based on its term frequency (tf) and inverse document frequency (Idf). The terms with higher weight scores are considered to be more important. Here, we are calculating only possible opinion target and possible opinion words.

Typically, the tf-idf weight is composed by two terms-Normalized Term Frequency (tf) and Computing the Term Frequency (Tf).

#### 1) Normalized Term Frequency (tf)

Frequency indicates the number of occurrences of a particular term  $t$  in document  $d$ . Therefore,

$$TF(t, s) = N(t, s). \quad (1)$$

Since we are dealing with the term frequency which rely on the occurrence counts, thus, longer documents will be favored more. To avoid this, normalize the term frequency.

$$TF(t, d) = \frac{N(t, d)}{\|D\|}. \quad (2)$$

Where,  $\|D\|$  is number of term in the document. Let's suppose, we have a document "T1" containing 10 words and the word "Alpha" is present in the document exactly 10 times.  $TF = 1/10 = 0.1$

#### 2) Inverse Document Frequency (IDF)

When the term frequency of a document is calculated, it can be observed that the algorithm treats all keywords equally. The inverse document frequency assigns lower weight to frequent words and assigns greater weight for the words that are infrequent. For example, we have 10 sentences and the term "technology" is present in 5 of those sentences, so the inverse document frequency can be calculated as,

$$IDF = \log e (10/5) = 0.3010 \quad (3)$$

### C. PSWAM

Opinion relation identification is defined as a word alignment process. To carry out monolingual word alignment, the word-based alignment model is used. Replicated every sentence to achieve a parallel corpus.

Replicated every sentence to achieve a parallel corpus. The constrained Hill-Climbing algorithm used to

find the alignments in the sentence. The constraint applied are as follows:

- Noun/ Noun phrases (adjectives/ verbs) must aligned with Adjective/ Verbs (Noun/ Noun phrases) or NULL words. Alignment with the null word specifies that it has no modifier or it modifies nothing.
- Other unrelated words such as preposition, adverbs and conjunction, symbols, must be aligned to themselves.

#### 1) Calculating the Opinion Association among the words

From the alignment results, we obtain a set of word pairs, each of which is composed of a noun/noun phrase (opinion target candidate) and its corresponding modified word (opinion word candidate). Next, the alignment probabilities between a potential opinion target  $wt$  and a potential opinion word  $wo$  are estimated using,

$$P\left(\frac{wt}{wo}\right) = \frac{Count(wo/wt)}{count(wo)} \quad (4)$$

Where,  $P\left(\frac{wt}{wo}\right)$  means the alignment probability between these two words. Similarly, we obtain the alignment probability  $P\left(\frac{wt}{wo}\right)$  by changing the alignment direction in the alignment process. Next, we use the score function in and to calculate the opinion association  $OA(wt, wo)$  between  $(wt)$  and  $(wo)$  is,

$$OA(wt, wo) = (\alpha * P\left(\frac{wt}{wo}\right) + (1-\alpha) * p(wo/wt))^{-1} \quad (5)$$

#### 2) Estimating candidate confidence with TF-IDF algorithm

We then calculate the confidence of each opinion target/word candidate, and the candidates with higher confidence are extracted as opinion targets or opinion words. We assume that two candidates are likely to belong to a similar category if they are modified by,

$$c_t^{k+1} = (1 - \mu) * M_{to} * c_o^k + \mu * (TF - score)_t \quad (6)$$

$$c_o^{k+1} = (1 - \mu) * M_{to} * c_t^k + \mu * (TF - score)_o \quad (7)$$

Where,  $c_t^{k+1}$  and  $c_o^{k+1}$  are the confidence of an opinion target candidate and opinion word candidate, respectively, in the  $k+1$  iteration.

$c_o^k$  and  $c_t^k$  are the confidence of an opinion target candidate and opinion word candidate, respectively, in



the  $k$  iteration.  $M_{to}$ , is the opinion associations among candidates.  $(TF - score)_t$  and  $(TF - score)_o$  are the TF-IDF score of possible opinion target and possible word.  $\mu$  is set to 0.3. The word alignment model modified by adding or considering TF-IDF score while calculating the Confidences of the candidates.

The candidate with higher confidence is collected as the opinion word for opinion target. The higher value of opinion word associated with opinion target as opinion relation is formed.

#### D. Sentimental Analysis

To perform sentimental analysis of the selected dataset, we have used the “words of bag” method. The adjectives/verbs present in the input are compared with the “words of bag”. If a word is positive oriented means found in the positive list of words then it will be considered as the positive opinion word and opinion target. If a word is negatively oriented means found in the negative list of words then it will be considered as the negative opinion word and opinion target. Our minor contribution is to generate positive, negative feedback.

#### IV. EXPERIMENTAL RESULTS

We have selected the four datasets of customer review with different size and different products as shown in below table. And positive negative segregation is done by using word of bag. Also Opinion Target and Opinion words are extracted efficiently by using TF-IDF score.

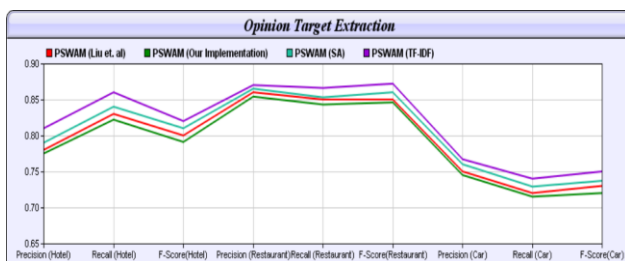


Chart-1: Experimental comparison among different Opinion Target Extraction technique.

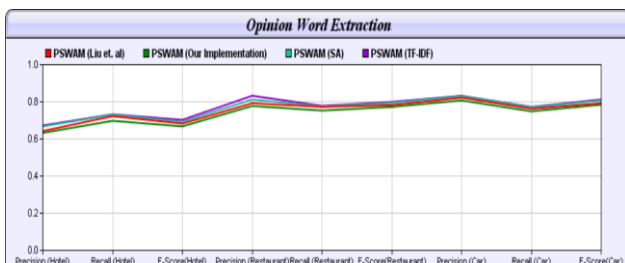


Chart-2: Experimental comparison among different Opinion word Extraction technique

In “Chart-1” and “Chart-2”, the Red Line shows the results of the Existing System implemented by Liu et. al [1]. The green line shows the results of the implementation of the existing system and the line with blue color shows the improvement in the results with minor contribution by changing the value of harmonic factor 0.4. The violate line shows the improvement in results by using TF-IDF algorithm The performance of different techniques are compared. Table 2 and Table 3

show the resultant values of opinion target and opinion of word extraction. The greater value of Precision, Recall and F-score indicate that the proposed system extracts the opinion target and opinion words efficiently. Where P denotes the Precision, R denotes Recall and F denotes F-score. The datasets of the Hotel, Car and Restaurant are used to compare the performance of the existing and proposed system. This proves the effectiveness of the proposed system.

#### V. CONCLUSION

In this paper, we described work on the mining opinion words by using a Modified partially supervised word alignment model. The purpose is extracting opinion words and opinion target and detecting opinion relations between them by using a partially supervised word alignment model that uses the TF-IDF score. The dynamic contribution is focused on calculating confidence by using TF-IDF score with improvement and sentiment analysis on customer review and categorize them as Positive reviews, Negative reviews and Detecting association between opinion targets and opinion words. Here, the model gives the positive, negative opinion about the product so that customer can decide whether to purchase a product or not and the manufacturer gets idea bout to increase the quality of the product in a timely manner. The TF-IDF score is calculated for opinion target and opinion words so that system able to extract opinion target and opinion words more efficiently. The experimental results show that our approach improved the performances of the mining task.

#### REFERENCES

- [1] Kang Liu, Liheng Xu, and Jun Zhao (2015), “Co-Extracting Opinion Targets and Opinion Words from Online Reviews Based on the Word Alignment Model.” IEEE transactions on knowledge and data engineering, vol. 27, no. 3, March 2015.
- [2] L. Zang, B. Liu, S.H. Lim and E. O’Brien-Strain (2010), “Extracting and ranking product features in opinion documents”, in Proc. 23th Int. Conf. Comput. Linguistics, Beijing, China, 2010, pp. 1462–1470.
- [3] M. Hu and B. Liu (2004), “Mining and summarizing customer reviews”, in Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, Seattle, WA, USA, 2004, pp. 168–177.
- [4] F. Li, S. J. Pan, O. Jin, Q. Yang, and X. Zhu (2012), “Cross-domain co-extraction of sentiment and topic lexicons”, in Proc. 50th Annu. Meeting Assoc. Comput. Linguistics, Jeju, Korea, 2012, pp. 410–419.
- [5] A.-M. Popescu and O. Etzioni(2005), “Extracting product features and opinions from reviews”, in Proc. Conf. Human Lang. Technol. Empirical Methods Natural Lang. Process., Vancouver, BC, Canada, 2005, pp. 339–346.
- [6] Shahzad Qaiser, Ramsha Ali (2018), “Text Mining : Use of TF-IDF to Examine the Relevance of Words to Documents.”, in International Journal of Computer Applications (0975 – 8887) Volume 181 – No.1, July 2018
- [7] Zheng Lu, Weiyao Lin (2013), “A New Algorithm for Inferring User Search Goals with Feedback Sessions”, in IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 25, NO. 3, MARCH 2013