# Word Alignment Model Based Optimized Opinion Words and Targets Co-Extraction from Online Reviews

**Manjushri Gagare[1], S . R Todmal[2]**
Department of Computer Engineering
[1, 2]Imperial College Of Engineering And Research

**Abstract-**Data mining is the field of pattern discovering process from large volume data sets. In this field, opinion mining is the subfield that analyzes the opinions of people, their sentiments, appraisals and emotions regarding available products and services. In today's internet world, rating systems becomes more popular for getting reviews of first time launched products. These online rating systems make use of opinion mining that is co-extractions of opinion targets and words from online reviews. But such systems are affected by the problem of accuracy. Such problem is solved here with the Word Alignment Process. This process finds out the relation between opinion words and targets, which is further used to find out the each candidates confidence. In this paper to balance the load and make faster execution, partial supervised learning technique is used along with syntactic patterns. This paper makes use of IBM4 word alignment model and hill climbing algorithm. Further this system used for web based application to recommend the merchandise based on review analysis. The system is tested with Sentiment Net Corpora and Google N-gram corpora, and proved that the IBM4 is more time efficient and accurate than IBM3 word alignment model.

*Keywords*-Opinion mining, opinion targets extraction, opinion words extraction.

## I. INTRODUCTION

Data mining is a process of searching, gathering and analyzing a large set of data in the database and discovering the relationships. There are many challenges that have given arise in data mining and one among them is opinion mining. Sentimental analysis is also known as Opinion mining which involves analyzing the emotions of the people towards building a system to categorize and collect opinions about various products and services. Interest has been growing rapidly in opinion mining in recent years because it mainly has many numbers of applications.

The opinion mining tasks at hand can be broadly classified based on the level at which it is done with the various levels being namely,

- The document level,
- The sentence level and
- The feature level.

At the document level, sentiment classification of documents into positive, negative, and neutral polarities is done with the assumption made that each document focuses on a single object O(although this is notnecessarily the case in many realistic situations such as discussion forum posts) and contains opinion from a single opinion holder. At the sentence level, identification of subjective or opinionated sentences amongst the corpus is done by classifying data into objective (Lack of opinion) and subjective or opinionated text. Subsequently, sentiment classification of the aforementioned sentences is done moving each sentence into positive, negative and neutral classes. This system makes the assumption that a sentence contains only one opinion which as in our previous levels is not true in many cases. An optional task is to consider clauses.

At the feature level, the various tasks that are looked at are:

Task1: Identifying and extracting object features that have been commented on in each review/text.
Task 2: Determining whether the opinions on the features are positive, negative or neutral.
Task 3: Grouping feature synonyms and producing a feature-based opinion summary of multiple reviews/text.

The main objective in co-extracting is collecting the opinions about the online reviews of the product. From opinion reviews, customers can obtain the information and the durability of the product, which direct their purchase actions. For the movement, manufactures can obtain the immediate response from the reviews and get an opportunity to improve the quality of their product. In opinion mining, extracting opinion targets and opinion words are the two fundamental sub tasks. Opinion targets are objects about which users opinions are expressed, and opinion words are which indicates opinion polarities. Extracting them can provide the essential information for obtaining fine-grained analysis on customer

opinions. Thus, it has attracted a lot of attention. If a word is an opinion word other words with which that word is having opinion relations will have high probability to be opinion target, and vice versa.

In this way, extraction is alternatively performed and mutually reinforced between opinion targets and opinion words. Although this strategy has been widely employed by previous approaches, it still has several limitations. They also propose a method which will formulate opinion targets or words extraction as a co-ranking task. All nouns/noun phrases are regarded as opinion target candidates and all adjectives or verbs are regarded as opinion word candidates, who are widely adopted by pervious methods. Then each candidate will be assigned a confidence and ranked, and the candidates with higher confidence than a threshold will be extracted as the results.

For example: "This phone is amazing, but the resolution of the display is bad" Here, the customer will be keen on knowing the reviews which gives the good or positive opinion on the phone and the bad or negative opinion on the display resolution, not just the reviews overall sentiments. After this extraction, the former step is to provide the relation among these words.For this, the graph on co ranking algorithm is used and the relation graph is used to provide the relation among them

In this paper study about the related work done, in section II, the proposed approach modules description, mathematical modeling, algorithm and expected result in section III. And at final we provide a conclusion in section IV.

## II. LITERATURE REVIEW

Kang Liu, LihengXu, and Jun Zhao [1], in this paper, authors propose the most complex word alignment model called the "IBM-3 model". It is also called as the fertility based model. "Word Alignment Model" is based on the syntactic patterns and nearest neighbour rule. IBM-3model has the capacity of capturing opinion relations which is more effective opinion word and opinion target extraction. This paper has mainly focused on opinion words and opinion targets and detecting the relations among them.

Minqing Hu and Bing Lu [2], in this paper, authors aim for mining and summarizing all the reviews of the customer based on that product. Here authors only mine the reviews and the features of the product based on the reviews of the user as negative or positive review opinion. Here, the work is mainly concerned with the Positive and the negative review orientation which is based on the adjective word or seed. The main objective here is to provide the huge number of customer reviews a feature based summary for the merchandise sold online, and the evaluation metric is based on precision and recall.

L.Zang, B.Liu, S.H.Lim, and E.O'Brien-Strain [3], in this paper, authors propose a ranking algorithm which is based on the web page called HITS. It is for relevance for applying the compute feature. In this proposed algorithm state-of-art problems which are used for the double propagation feature extraction. In this paper the feature ranking and the feature extraction are the two approaches that are proposed to deal with the problems of coextracting the opinion reviews. Here in this feature each candidate is ranked with the importance. The HIT algorithm is used for web page and relevance ranking.

Kang Liu, LihengXu, Jun Zhao [4], in this paper, authors propose the word-based translation model (WTA) which is used for the extraction of opinions. Here the association between the opinion words and opinion targets are mined together. In WTM, the word positions the frequencies and other attributes are compared with the adjacent method which can be considered globally. This will give the ranking frame work for the opinion targets. The main objective is to formulate the opinion words and opinion targets as the word alignment task. The mining association between the opinion target and the opinion word is the two major components for extracting the opinion targets.

Fantago Li, SinnoJailin Pan, Ou Jin, Qiang Yang and Xiaoyan Zhu [5], in this Paper, authors propose the framework that is based on domain adaption method. This is the domain for co-extracting the sentiment-and –topic – lexicon based interests. The algorithm such as Relational Adaptive bootstrapping (RAP) is used to expand the seeds in target domain. The topic seeds and high confidence sentiment is generated and expanded by the target domain. The topic-lexicon co-extraction and sentimental analysis is a twofold framework.

G.Qiu, L.Bing, J.Bu and C.Chen [6], in this paper, authors propose the novel propagation based method as the solution for the target extraction and the opinion lexicon expansion .They are also better in performance compared to state-of-art method. Here the additional requirements of resources are not required. The initial steps of the opinion lexicon are used for the extraction of the opinion relation. Here the system extracts the opinion words from the previous iteration seeds of the opinion words and later uses these words to target it through the identification process of syntactic

relations. Here the relation between the opinion words and target words are used for the relation identification.

Robert C. Moore [7], in this paper, author has described the descriptive approach for training of simple word alignment model which has more accuracy than the complex generative method. The IBM, HMM and LogLikelihood-Based Model is used for the measurement of associations, the LLR score for pair of words is high when there is a strong positive association.

Fangato Li, Chao Han, Minlie Huang, Xiaoyan Zhu, YinhJu Xia, Shu Zhang and Hao Yu [8], in this paper, authors propose a framework known as the w machine learning which is based on the conditional random fields (CRF). CRF has the rich features for extracting positive and negative opinions. X.Ding,

B.Liu, and P.S.Yu [9], in this paper, authors propose the semantic orientation opinion methods, here both implicit and explicit methods of opinion are considered. Here the summarization of review is based on the object feature. Object feature, opinion extraction and opinion polarity detection are the purpose of the new machine learning framework which is based on Conditional Random Fields (CRFs). CRF can integrate many features than the Lexicalized HMM model.

Yuanbin Wu, Qi Zhang, XuanjingHaung, Lide Wu [10], in this paper, authors propose the opinion mining for the unstructured documents. Dependency tree is constructed for the extraction of relation between opinion expression phrase and product features. Here, opinion expression, emotional attitude and product feature is all combined to form the opinion phrase unit which are useful for opinion mining tasks. The phrase dependency tree, SVM-WTree and SVM-PTree are used for the extraction of features

### III. PROPOSED APPROACH

#### A. Problem Statement

For given online reviews, develop a system for co-extraction optimized mining of opinion words and opinion targets, by using constrained hill climbing algorithm and word alignment process with IBM4 algorithm.

#### B. Proposed System Overview

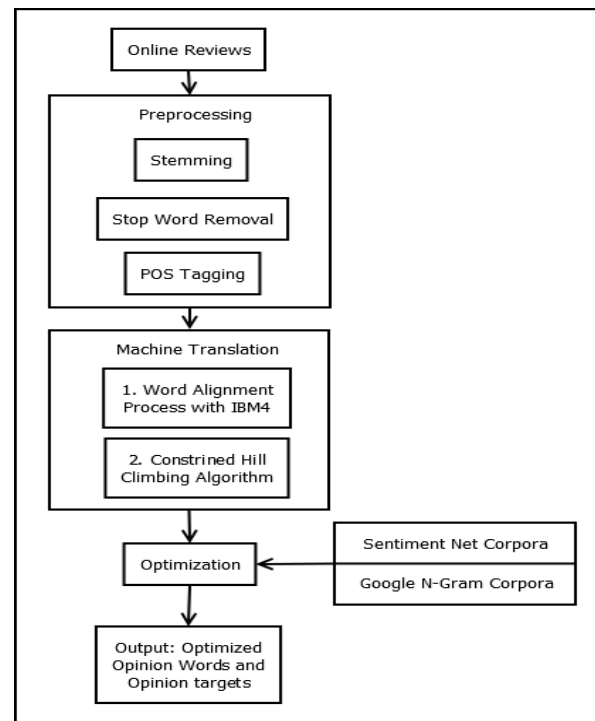Figure 1 depicts the architectural view of proposed system



Figure 1.Proposed System Architecture

Initially this system takes online reviews aninput.This system performs extraction of opinion words and opinion targets as a process of co ranking. This system considers, all nouns or noun phrases as a opinion words and all adjectives or verbs as a opinion targets in given review sentences. Initially all online reviews are goes through preprocessing step. In this, stemming is performed, stop words are removed and POS tagging is also performed. These preprocessed data is pass to the machine translation module. In this module, initially word alignment process is conducted by using IBM4 which formulates opinion relation identification. Then alignment of sentences is calculated by suing contained hill climbing algorithm. These extracted opinion words and targets are optimized against two corpora named as: Sentiment Net Corpora and Google N-Gram Corpora.

#### C. Mathematical Model

Let S is system, such that,
S = {Input, Process, Output}
Input: Online Reviews extracted from web applications
R = {r1, r2, ,r3, …, rn}
Where, R is the set of online reviews.
For these reviews, system needs to extract opinion words and targets.

Output:

Optimized Opinion Words

OW = {ow1, ow2, …, own}
Where, OW is the set of optimized opinion words extracted by system.

Optimized Opinion Targets
OT = {ot1, ot2, …,otn}
Where, OT is the set of optimized opinion targets extracted by system.

Process:

1. Preprocessing
   Preprocessing of input online reviews is performed with three techniques such as:
   P = {p1, p2, p3}
   Where,
   P1 = Stemming
   is the process of reducing inflected (or sometimes derived) words to their word stem or root formgenerally a written word form.

   P2 = Stop Word Removal
   **stop words** are words which are filtered out before or after processing of natural language data.
   P3 = POS Tagging
   Part-of-speech tagging (POS tagging or POST), also called grammatical tagging or word
   category disambiguation, is the process of marking up a word in a text (corpus) as corresponding to a particular part of speech, based on both its definition and its context—i.e., its relationship with adjacent and related words in a phrase, sentence, or paragraph

2. Machine Translation
   In this step, opinion relations between opinion targets and opinion words are captured by using two methods.
   M = {m1, m2}
   Where,
   m1 = Word Alignment Process with IBM4
   This step formulates opinion relation identification.
   m2 = Constrained Hill Climbing Algorithm
   This step calculates alignment for sentences

   O = {ow, ot}
   Where, ow is the opinion words and ot is the opinion targets extracted.

3. Optimization
   For optimization of O, two different corpora are used.
   Corpora 1: Google N-Gram Corpora

Corpora 2: Sentiment Net Corpora
OO = {oow, oot}
Where, OO is the set of optimized opinion words and optimized opinion targets.

**D. Algorithm**

**Input:** Review sentences $S_i = \{w_1, w_2, \ldots, w_n\}$
**Output:** The calculated alignment $\hat{a}$ for sentences
1 **Initialization:** Calculate the seed alignment $a_0$ orderly using simple model (IBM-1, IBM-2, HMM)
2 **Step 1**: Optimize toward the constraints
3 **while** $N_{ill}(\hat{a}) > 0$ **do**
4    **if** $\{a: N_{ill}(a) < N_{ill}(\hat{a})\} = \emptyset$ **then**
5      break
6    $\hat{a} = argmax_{a \in nb(\hat{a})} Pro(f|e, a)$
7 **end**
8 **Step 2**: Toward the optimal alignment under the constraint
9  **for** $i < N$ and $j < N$ **do**
10   $M_{i,j} = -1$, if $(i, j) \notin \hat{A}$;
11 **end**
12 **while** $M_{i_1, j_1} > 1$ or $S_{j_1, j_2} > 1$ **do**
13   **If** $(j_1, a_{j_2}) \notin \hat{A}$ or $(j_2, a_{j_1}) \notin \hat{A}$ **then**
14    $S_{j_1, j_2} = -1$
15   **end**
16   $M_{i_1, j_1} = \arg\max M_{i,j}$
17   $S_{j_1, j_2} = \arg\max S_{i,j}$
18   **If** $M_{i_1, j_1} > S_{j_1, j_2}$ **then**
19    Update $M_{i_1, *}, M_{j_1, *}, M_{*, i_1}, M_{*, j_1}$
20    Update $S_{i_1, *}, S_{j_1, *}, S_{*, i_1}, S_{*, j_1}$
21    set $\hat{a} := M_{i_1, j_1}(a)$
22   **end**
23   **else**
24    Update $M_{i_1, *}, M_{j_2, *}, M_{*, i_1}, M_{*, j_2}$
25    Update $S_{j_2, *}, S_{j_1, *}, S_{*, j_2}, S_{*, j_1}$
26    set $\hat{a} := S_{j_1, j_2}(a)$
27   **end**
28 **end**
29 **return** $\hat{a}$;

## IV. RESULTS AND DISCUSSION

**A. Expected Result**

In this section discussed the experimental result of the proposed system.

Optimal Graph:

In the following graph 13.1 shows the optimal graph for the different reviews, x-axis shows the different reviews and y-axis the shows the rating of that reviews.
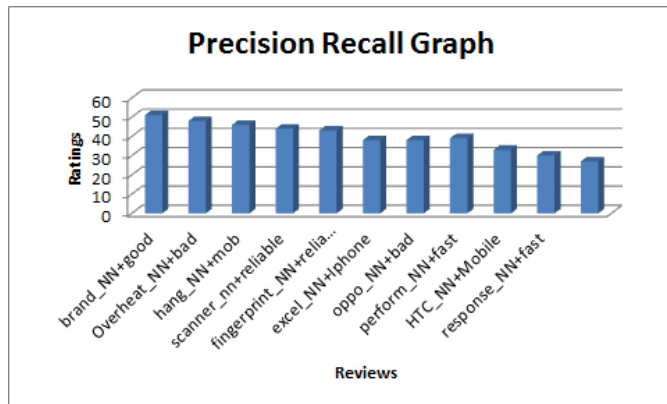


Fig. 2 Optimal Graph

Precision, Recall, F-measure Graph:

Here, X-axis represents precision, recall and F-measure, Y-axis represents precision, recall and F-measure in percentage.
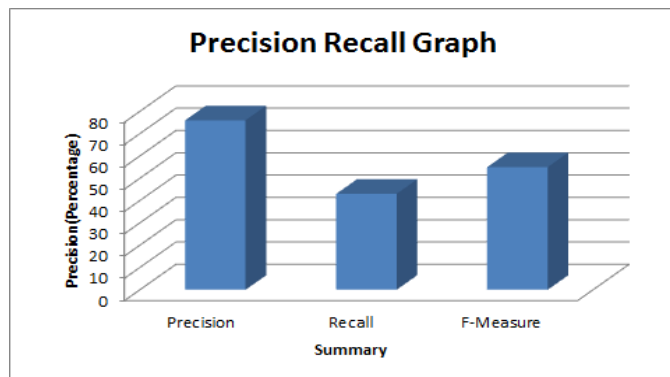


Fig. 2: Precision, Recall and F-measure Graph

## V.CONCLUSION AND FUTURE SCOPE

This paper concludes that the popularity of internet arises the need of online reviews study of any product to clarify the user thoughts about newly launched product. For this purpose, this paper develops the co-extraction of opinion words and targets by using IBM4 word alignment model. To improve the performance, system is combined with hill climbing algorithm. This system overcomes the accuracy and time efficiency problems of IBM3 word alignment model. The performance of system is tested on Sentiment Net Corpora and Google N-Gram Corpora. From experimental results, we conclude that the proposed system is very effective for extraction of opinion words and opinion targets.

## ACKNOWLEDGMENT

## REFERENCES

[1] Kang Liu, LihengXu, and Jun Zhao, "Co-Extracting Opinion Targets and Opinion Words from Online Reviews Based on the Word Alignment Model". IEEE Transactions on Knowledge and data engineering, march 2015.

[2] M. Hu and B. Liu, "Mining and summarizing customer reviews," in Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, Seattle, WA, USA, 2004, pp. 168–177.

[3] L. Zhang, B. Liu, S. H. Lim, and E. O'Brien-Strain, "Extracting and ranking product features in opinion documents," in Proc. 23th Int. Conf. Comput. Linguistics, Beijing, China, 2010, pp. 1462–1470.

[4] K. Liu, L. Xu, and J. Zhao, "Opinion target extraction using word based translation model," in Proc. Joint Conf. Empirical Methods Natural Lang. Process. Comput. Natural Lang. Learn., Jeju, Korea, Jul. 2012, pp. 1346–1356.

[5] F. Li, S. J. Pan, O. Jin, Q. Yang, and X. Zhu, "Cross-domain co extraction of sentiment and topic lexicons," in Proc. 50th Annu. Meeting Assoc. Comput.Linguistics, Jeju, Korea, 2012, pp. 410– 419.

[6] G. Qiu, L. Bing, J. Bu, and C. Chen, "Opinion word expansion and target extraction through double propagation," Comput. Linguistics, vol. 37, no. 1, pp. 9–27, 2011.

[7] R. C. Moore, "A discriminative framework for bilingual word alignment," in Proc. Conf. Human Lang. Technol. Empirical Methods Natural Lang. Process., Vancouver, BC, Canada, 2005, pp. 81–88.

[8] F. Li, C. Han, M. Huang, X. Zhu, Y. Xia, S. Zhang, and H. Yu, "Structure-aware review mining and summarization." in Proc. 23th Int. Conf. Comput. Linguistics, Beijing, China, 2010, pp. 653–661.

[9] X. Ding, B. Liu, and P. S. Yu, "A holistic lexicon-based approach to opinion mining," in Proc. Conf. Web Search Web Data Mining, 2008, pp. 231–240.

[10] Y. Wu, Q. Zhang, X. Huang, and L. Wu, "Phrase dependency parsing for opinion mining," in Proc. Conf.

Empirical Methods Natural Lang. Process., Singapore, 2009, pp. 1533–1541.