

# An efficient word alignment model for co-extracting opinion targets and opinion words from online reviews

J. Yesudoss<sup>1</sup>, T. Banusankari<sup>2</sup>

<sup>1</sup> Assistant Professor, <sup>2</sup> Research scholar, Master of Philosophy, Department of Computer Science, Sri Ramakrishna Mission Vidyalaya College of Arts and Science, Tamil Nadu, India.  
jydoss@gmail.com, banusankarimphil@gmail.com

## Abstract

**Objectives:** The main objective of this research is to improve the topical relations by extracting the opinion targets as well as opinion words, and achieve the higher performance using word alignment model concept.

**Methods:** Partially Supervised Word Alignment Model (PSWAM) is used for word alignment in existing system. The Latent Dirichlet Allocation (LDA) model is used for discovering opinion word relation extraction in proposed system.

**Findings:** The proposed method achieves high performance in terms of sensitivity and specificity.

**Application/Improvements:** The proposed system is done by using Latent Dirichlet Allocation (LDA) which is used to increase the performance for number of dataset more efficiently.

**Keyword:** Opinion mining, word alignment model, opinion targets extraction, opinion words extraction.

## 1. Introduction

Generally, data mining is the search for hidden patterns present in huge databases. Data mining scans via a huge volume of data to find out the patterns and correlations between patterns. Data mining requires the use of data analysis tool to determine previously unknown, valid patterns and relationships from the data. Such kind of tool can enclose statistical model, mathematical algorithms and machine learning methods. Thus, data mining technique is the way of getting analysis and prediction results more than gathering and running data. Data mining can be executed on data signified in quantitative, textual or multimedia forms. Data mining application could use several parameters to inspect the data. They contain the concepts such as association, sequence analysis, classification, clustering and forecasting.

Opinion mining is an important factor in the domain of data mining and it is also called as Sentiment analysis. The opinion mining is used to analyze the people's opinions, emotions, assessments and attitudes. Along with the explosive growth of user created messages, web sites and social networks has become a significant media for where millions of users can communicate their opinions [1]. This is typically hard to discover an accurate reason of opinion variations because they might involve complicated factors. It is examined that the promising topics suggested in variation period can be highly connected to authentic reasons behind the opinion variations. While people communicate their opinions, they frequently state reasons for some specific events or topics to support their present views and ideas [2] [3].

In opinion mining, the important issue is to mine opinion targets, which is described as the objects or classes also on customers have articulated their opinions, classically as nouns, adjectives or phrases. To mine and examine opinions from online reviews, it is unacceptable to simply attain the overall sentiment about a product. In many scenarios, users suppose to discover fine grained opinions<sup>1</sup> about a characteristic or feature of manufactured goods that is examined. In such scenario, the word alignment model is improved to investigate the number of document reviews more significantly. The scenario used the method word alignment model along with partially supervised approach for evaluating the reviews. It is used to estimate the opinion targets and opinion words [4]. This research work is introduced the approach named as constrained hill climbing algorithm [5] which is used to analyze the review sentences from the specified documents.

In [6] discussed analysis of opinions using double propagation approach. This research scenario is focused on the identification of relations by using a parser and it is used to enlarge the initial opinion lexicon and to mine targets [7]. The proposed double propagation method is used to produce the important syntactic relations. In [8] suggested support vector machine based technique for classifying the opinion targets from given document. In [9] discussed the

opinion mining concepts using LDA based hybrid approaches. This model is focused on the creation of model to mutually determine the aspects as well as aspect specific opinion words [10].

## 2. Materials and Methods

### 2.1. Word Alignment Model

This model develops opinion relation recognition as a word alignment process. It uses the word-based alignment model to execute monolingual word alignment that is extensively utilized in several tasks such as collocation extraction. In performance, each sentence is simulated to produce an equivalent corpus. A bilingual word alignment algorithm is used to the monolingual scenario to arrange a noun or noun phrase (possible opinion targets) along with its modifiers (possible opinion words) in sentences. In specific, if it to directly concerns the typical alignment model to this task, an opinion target candidate (noun or noun phrase) will arrange with the inappropriate words rather than potential opinion words (adjectives or verbs), such as prepositions and conjunctions [11]. Therefore, the scenario introduced some restrictions in the alignment model as follows:

- 1) Nouns or noun phrases (adjectives/verbs) should be arranged with adjectives/verbs or a null word. Arranging to a null word implies that this word either has no modifier or modifies nothing;
- 2) Other dissimilar words, such as prepositions, conjunctions and adverbs, could only support with themselves.

### 2.2. Partially-Supervised Word Alignment Model

The typical word alignment model is classically trained in a completely unsupervised manner, which may not obtain precise alignment results. Thus, to progress alignment process, the algorithm execute a partial supervision on the statistic model and utilize a partially- supervised alignment model to integrate partial alignment links into the alignment process. In this research, the partial alignment links are considered as conditions for the trained alignment model.

#### 2.2.1 Parameter Estimation for the PSWAM

Unlike the unsupervised word alignment model, the arrangements created via the PSWAM should be as reliable as probable along with the labeled partial alignments. To accomplish this objective, the model improves an EM-based algorithm. For training an easier arrangement model, such as the IBM-1 and IBM-2 models, the users imply achieve every probable alignment from the experiential information data. Those incompatible alignments along with pre-provided partial alignment links is clean out; consequently, they will not be counted for parameter assessment in succeeding iterations. However, in this scenario, we select a more complex alignment model, the IBM-3 model, which is a fertility- based model. For training IBM-3 model, it is NP-complete and unfeasible to specify every possible arrangements. It specifies that the typical EM training approach is time consumption and not practical. To solve the above mentioned issue, GIZA++ produces a hill-climbing technique, which is a local optimal solution to speed up the training process.

#### 2.2.2 Obtaining Partial Alignment Links by Using High-Precision Syntactic Patterns

In nature, the model can alternate to manual labeling. However, this approach is both time consideration and unfeasible for numerous domains. The scenario requires an automatic process for partial alignment creation. To perform this goal, it is transformed to syntactic parsing. As stated in the initial segment, though present syntactic parsing tools cannot acquire the whole correct syntactic tree of familiar sentences, straight syntactic dealings is still achieved exactly. Hence, some higher accuracy lower syntactic models are considered to confine the opinion relations amongst words for initially producing the partial alignment links. It is then sent to further alignment process.

### A. LDA with GIZA++ tool for word alignment model

Latent Dirichlet allocation (LDA) is proposed method which is used for improving the topical relations in given documents. This method is an efficient model which permits sets of annotations to be described via unobserved groups that clarify why some parts of the information are alike. For example, if observations are words collected into documents, it posits that each document is a mixture of a small number of topics and that each word's creation is attributable to one of the document's topics. LDA is an instance of a topic model and it is originally accessible as a graphical model for topic discovery.

**B. Calculating the Opinion Associations among Words**

From the alignment results, the scenario attains a group of word pairs, each of which is collected of a noun/noun phrase and its equivalent customized word [12]. The estimation of the alignment probability between two words is as follows.

$$P(w_t|w_o) = \frac{Count(w_t, w_o)}{Count(w_o)} \quad (1)$$

Where  $P(w_t|w_o)$  means the alignment possibility among these two words. Likewise, the scenario acquire the alignment likelihood  $P(w_o|w_t)$  through varying the alignment track in the alignment procedure. Subsequently, we employ the score value and to determine the opinion relationship among  $w_t$  and  $w_o$  using the given below formula.

$$OA(w_t w_o) = (\alpha * P(w_t|w_o) + (1 - \alpha)P(w_t|w_o))^{-1} \quad (2)$$

Where  $\alpha$  is the harmonic factor which is utilized to unite these two alignment possibilities and OA is opinion association. In this scenario, we set  $\alpha = 0.5$ .

**C. Estimating candidate confidence with graph co-ranking**

After extracting the opinion relations among opinion target candidates and opinion word candidates, the scenario finish the creation of the opinion association grid. Then the scenario compute the assurance of every opinion target or word entrant on this graph, and the candidates along with superior confidence than a threshold are mined as opinion targets or opinion words. In this research work, the method considers that two candidates are possibly to belong to a related group if they are adapted through identical opinion words or change parallel opinion targets. Thus, it can forward the confidences amongst diverse candidates, which specify that the graph-based approaches are appropriate.

**3. Results and Discussion**

In this section, the performance metrics are evaluated by using existing and proposed methodologies. The performance metrics are such sensitivity and specificity metrics. The existing PSWAM model is shown the lower sensitivity and specificity values for number of dataset files. The proposed LDA model has shown the higher sensitivity and specificity values for number of dataset files. From the experimental result, we can conclude that the proposed LDA model is better than the existing method in terms of higher performance. An experimental result shows that the proposed method achieves high performance in terms of sensitivity and specificity.

**3.1. Sensitivity**

Sensitivity is an absolute quantity, the smallest absolute amount of change that can be detected by a measurement. Sensitivity refers to the test's ability to correctly detect records from the given documents.

$$Sensitivity = \frac{\text{number of true positives}}{\text{number of true positives} + \text{number of false negatives}}$$

In this graph, x axis is taken for two methods of and y axis is taken for sensitivity. From the Figure.1 the proposed scenario shows the highest sensitivity rather than existing method. LDA provides higher performance for number of dataset file in proposed method.

Figure 1. Sensitivity comparison

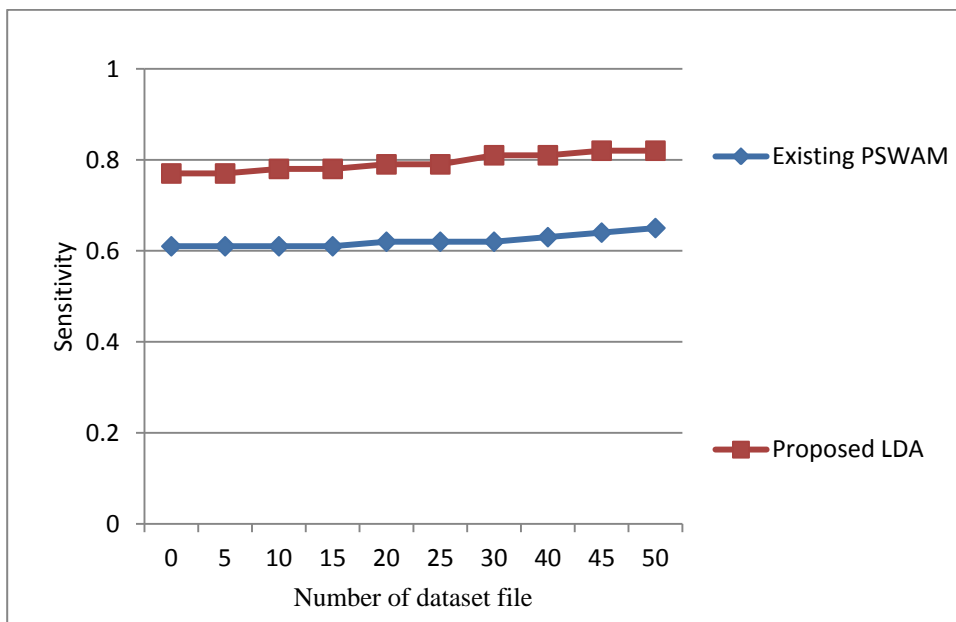
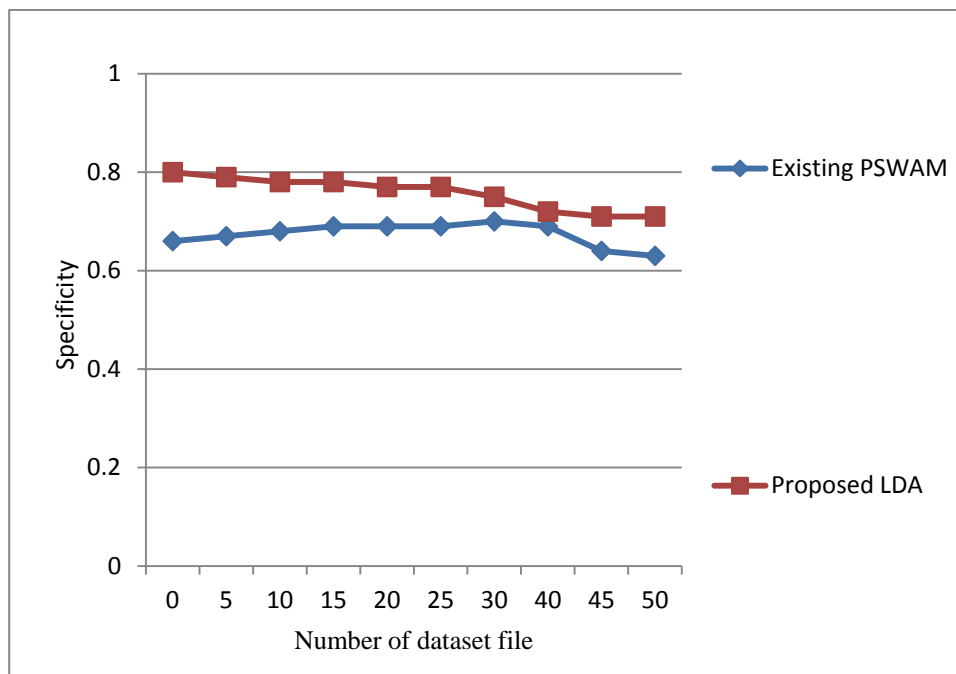


Figure 2. Specificity comparison



**3.2. Specificity**

Specificity (also called the true negative rate) measures the proportion of negatives that are correctly identified as such. Specificity relates to the test's ability to correctly detect records from the specified documents.

$$\text{Specificity} = \frac{\text{number of true negatives}}{\text{numebr of true negatives} + \text{number of false positives}}$$

In this graph, x axis is taken for two methods of and y axis is taken for specificity. From the Figure.2 the proposed scenario shows the highest specificity rather than existing method. LDA provides higher performance for number of dataset file in proposed method.

## 4. Conclusion

The existing scenario method is used for co-extracting opinion targets and opinion words by using a word alignment model. The existing research is focused on discovering opinion relations among opinion targets and opinion words. The proposed scenario is focused on the discovering the topical relations using LDA method. From the result we can conclude that the proposed method is better than the existing method. The items along with higher probabilities are mined out. In the proposed system the LDA model is used to provide effective and efficient topical relations among sentences. From the experimental result we can conclude that, proposed scenario yields higher performance rather than existing scenario. The performance is superior in terms of sensitivity and specificity values. Hence the proposed method is higher accuracy by using LDA method rather than existing scenario.

## 5. Acknowledgement

We the authors assure you that, this is our own work and also assure you there is no conflict of interest.

## 6. References

1. Manju, S. Revathi, E. V. R. M. Kalaimani, R. Bhavani. Product Aspect Ranking Using Semantic Oriented Sentiment Classifier. *International Journal of Scientific Engineering and Research (IJSER)*. 2014; 2(10), 25-28.
2. S. Akilandeswari, A.V.Senthil Kumar. A novel approach for mine infrequent weighted itemset using coherent rule mining algorithm. *Indian Journal of Innovations and Developments*. 2015; 4 (3), 1-6.
3. Lee, JiHye, SeungYeob Yu. Cognition difference on online public opinion dissonance between Korean and Chinese Netizens: Its causes, functions and solutions. *Indian Journal of Science and Technology*. 2015; 8(26), 1-12.
4. Liu, Kang, LihengXu, Jun Zhao. Extracting opinion targets and opinion words from online reviews with graph co-ranking, Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. 2014; 1.
5. Liu, Kang, et al, Opinion target extraction using partially-supervised word alignment model, Proceedings of the Twenty-Third international joint conference on Artificial Intelligence, AAAI Press, 2013.
6. Qiu, Guang, et al, Opinion word expansion and target extraction through double propagation, *Computational linguistics*. 2011; 37(1), 9-27.
7. R. Suganthi, P. Kamalakannan, Exceptional patterns in multi database mining, *Indian journal of Innovations and Developments*. 2015; 4(4),1-4.
8. Ma, Tengfei, Xiaojun Wan. Opinion target extraction in Chinese news comments, Proceedings of the 23rd International Conference on Computational Linguistics: Posters, Association for Computational Linguistics, 2010.
9. Zhao, Wayne Xin, et al, Jointly modeling aspects and opinions with a MaxEnt-LDA hybrid, Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2010.
10. Zol, Sapna, Preeti Mulay. Analyzing sentiments for generating opinions (ASGO)-A New Approach. *Indian Journal of Science and Technology*. 2015; 8.S4, 206-211.
11. Liu, Zhiyuan, Xinxiong Chen, Maosong Sun, A simple word triggers method for social tag suggestion, Proceedings of the Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2011.
12. M. Hu, B. Liu. Mining opinion features in customer reviews, in Proceedings 19th National Conference Artificial Intelligence., San Jose, CA, USA, 2004, pp. 755–760.

*The Publication fee is defrayed by Indian Society for Education and Environment (iSee). [www.iseeadyar.org](http://www.iseeadyar.org)*

### Citation:

J. Yesudoss, T. Banusankari. An efficient word alignment model for co-extracting opinion targets and opinion words from online reviews. *Indian Journal of Innovations and Developments*. 2015; 4 (7), November.