# Words Segmentation-based Scheme for Implicit Aspect Identification for Sentiments Analysis in English Text

Dhani Bux Talpur[1]

School of Information and Communication Engineering
Guilin University of Electronic Technology
Guilin, China

Guimin Huang[2,*]

Guangxi Key Laboratory of Trusted Software
School of Computer Science and Information Security
Guilin University of Electronic Technology, Guilin, China

*Abstract*—**Implicit and Explicit aspects extraction is the amassed research area of natural language processing (NLP) and opinion mining. This method has become the essential part of a large collection of applications which includes e-commerce, social media, and marketing. These application aid customers to buy online products and collect feedbacks based on product and aspects. As these feedbacks are qualitative feedback (comments) that help to enhance the product quality and delivery service. Whereas, the main problem is to analyze the qualitative feedback based on comments, while performing these analysis manually need a lot of effort and time. In this research paper, we developed and suggest an automatic solution for extracting implicit aspects and comments analyzing. The problem of implicit aspect extraction and sentiments analysis is solved by splitting the sentence through defined boundaries and extracting each sentence into a form of isolated list. Moreover, these isolated list elements are also known as complete sentence. As sentences are further separated into words, these words are filtered to remove anonymous words in which words are saved in words list for the aspects matching; this technique is used to measure polarity and sentiments analysis. We evaluate the solution by using the dataset of online comments.**

*Keywords*—*Implicit aspect; explicit aspects; polarity; sentiments analysis*

## I. INTRODUCTION

With the advancements in the field of technology more and more peoples are in touch with online shopping websites and these numbers are increasing day by day. This innovative move transfer street shopping into online shopping. The most popular trending websites includes Taobao, JD, Alibaba and Amazon etc. These e-commerce websites generally provide an easy and accessible platform for customers, where consumers can share the experience with feedback regarding products. With the help of these feedbacks it is easy to extract opinions and aspects of entities from various online comments of consumers as these reviews can help to provide opinion which can further use for prediction.

The feedback help consumers know about popular trends and aspects of these products to buy. Recently, different approaches are introduced on this field as few models were also proposed to process specific task. These specific tasks are the basic part of the NLP application is words segmentation, the procedure of separating and dividing the sentence into a

single token of words is called Word Tokenization [1]. In Natural Language Processing (NLP) the term tokenization or word segmentation is thought as the most important task [2]. Mostly each application of NLP needs at a certain level the process of breaking its text into distinct tokens for processing. The tokenization and extraction method is done by identifying word borders in languages like English where punctuation marks or white spaces are used to isolate words [3].

Many sentiment analysis tools and applications have been developed to mine the opinions in user-generated content on the Web. However, the performances are very poor due to the complexity of natural language [4,5,6]. In essence, sentiment analysis is still a problem of natural language processing (NLP), which deals with the natural language documents, which are also called unstructured data [7]. Prior researches show that sentiment analysis is more difficult than the traditional topic-based text classification [8]. Although various methods have been projected to conduct sentiment analysis, it is still difficult to deal with some linguistic phenomena, such as negation and mix-opinion text. This indicates to low accuracy of sentiment classification [9,10]. Besides, it is insufficient to only determine the polarity of the opinions, since an opinion without a target is of limited use. The task of extracting the opinions and their targets simultaneously is also called aspect-level sentiment analysis in the research literature and is more difficult to achieve [11].

Furthermore, in order to achieve the finest information that is required for such analysis, the different aspects and features of a product or service must be identified in the comments section. There are different examples of such features include size, price, service and parts of product aspect which are mentioned in this text. Some of the examples are illustrated below.

"The mobile size is very large but picture quality is awesome and price is cheap". In this sentence 'size','price' and 'quality' are all aspects on which sentiments is expressed.

In the proposed work, it is consider how to extract the implicit aspects and sentiments analysis on an aspect level.

The recent research has concluded that there was no parallel development available in which aspect extraction and sentiments analysis work together. The research work is

*Corresponding Author

divided into different portions. Section 2 describes an overview of related work with detailed methods of implicit sentiment analysis. However in Sect. 3 defines the research methodology and the results. Moreover in section 4 a detail overview of experiments and discussion analysis has been illustrated. Finally, the Section 5 gives brief outline of conclusion and future study.

## II. RELATED WORK

The sections of literature review discuss the literature of relevant fields which are concerned with implicit aspects extraction and sentiments analysis of customer reviews. Based on review analyzing the background and leverage some previous methodology of opinion mining and sentiment analysis to design solution by using distant methodologies. Furthermore sentiment analysis also deals with the computational treatment of sentiment, views and subjectivity in English text [12]. This approach is also known as opinion mining. Earlier developed schemes of natural language processing deals with sentiments analysis includes diverse aspects concerning how information about emotions, attitudes, perspectives and social identities is conveyed in the language.

There are two main approaches used for sentiments analysis which are further separated into lexicon-based and machine learning methodologies. The following approaches are combined to resolve issues related to sentiment analysis and aspects extraction Fig. 1 show different approaches:

The most useful approach in text classification as it is used for sentiment analysis with the help of machine learning. There are two different types of machine learning datasets required to perform classification one is known as training set and other is test set. Machine Learning uses text-classification approach for sentiment analysis. Two datasets are required to perform machine learning. The first one is called as training set and the other is test set. The training set is trained by an automatic classification to learn characteristics of documents while the test set is used to validate the performance of automatic classifier. It is also use to document level or sentence level sentiment analysis to predict whether the reviews are positive, neutral or negative. This segmentation can further be classified into supervised learning and unsupervised learning.

This can also be further classified into supervised learning and unsupervised learning.

*a) Supervised Learning:* Supervised learning is also known as label learning. It has various classifier algorithms to execute sentiment analysis of an opinion. It requires a labeled training data set to determine the invisible instances. Support Vector Machine (SVM), Naïve Bayes (NB), NBtree, J48 and LDA are some of the techniques used for this classification.

*b) Unsupervised Learning:* Unsupervised learning is also known as unlabeled learning. It uses set of inputs for clustering where labels are not known during the training of the data sets. Classification is performed based on opinion words and phrases.
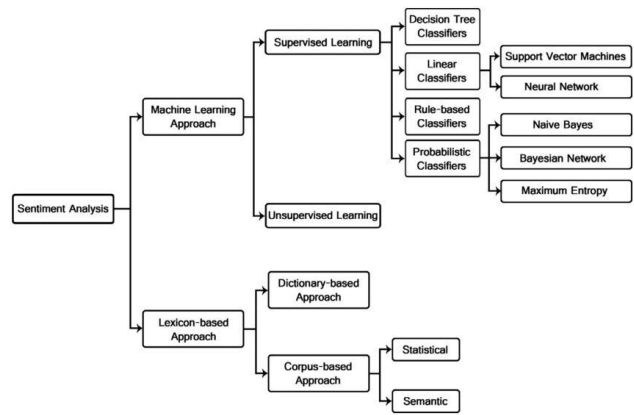


Fig. 1. Different Approaches for Sentiments Analysis.

The lexicon based approach depends on a collection of known and pre-determined domain sentiment terms. It is further classified into dictionary-based and corpus-based semantic methods used to find the polarity and score from large corpora. The dictionary-based approach searches the opinion word and finds its synonyms and antonyms in the seed list to measure the polarity from a dictionary. The corpus-based approach is an independent domain approach which is developed by the author according domain used. So, the sentiment words in the corpus are context specific.

The lexicon-based scheme is use for sentiments or subjectivity lexicons [13, 14, 15] to develop a corpus-based dictionary by remodeling the insights of the co-occurrence and conjunction method. These lexicons contain sentiment words that are also called opinion words listed with their polarity and strength. They tried the dictionary with the *Cornell Movie Review Data* and claim the accuracy of their measure up to 72.5%, which match the accuracy of machine learning classifiers [16]. In sentiment classification lexicon- based approach is applied by author, part of peech tagging, SentiWordNet and WordNet combined with a weighted model provided by Natural Language Processing (NLP)[17]. The dissimilarity in sentiments classification, opinion extraction goals are to produce wealthier information and need in-depth analysis of reviews [18].

The majority of approaches in the literature worked on extraction of the explicit features in sentences system are require to automatically find out and analyze the online opinionated texts (texts with opinions or sentiments), sentiment analysis grow out of its need. Some of the approaches are already mentioned in the reviews. However, implicit features have equally balance as explicit feature in the review. [19] It presents a fine-grained method for the labels of polar sentiment in text and for explicit sentiments on the other hand as well as implicit expressions of sentiment of polar facts.[20] The proposed scheme deals with a double-implicit problem in opinion mining and sentiment analysis. This methodology deals with recognizing features and polarities of opinionated reviews not consisting of opinion words and aspect terms. [21] It also Propose fuzzy-based information in engineering approach that has been developed for sentiment classification of a distinct group of such sentences that include the change or deviation from the desired range or value.

In the previous research works, most of the researches is done on implicit aspects extraction and sentiments analysis but a limited work has been done on the segmentation problem. The proposed classification approach uses words segmentation and a rule-based method to extract implicit features. However, aspect-based opinion mining gives attention on the task of extracting opinion about aspects. This research paper develops the implicit aspect extraction with a segmentation scheme and sentiment analysis approach for English text.

## III. RESEARCH METHODOLOGY

### A. Framework for Sentiment Analysis

In the real-world the sentiment analysis based on prior research identified by research gaps and requirements. In order to fulfill these research gaps is an effective approach for sentiment analysis; it is an innovative framework which is proposed in this research. By implementing proposed framework of sentiments analysis on customer reviews can predict the aspect level. It can not only analyze single-opinionated text but also mixed-opinion text. The framework provides an innovative way to detect sentimental words in the text through sentiment lexicon to gain more information from text. In addition to this the framework offers an effective approach to conduct aspect- level or feature-based sentiment analysis. Fig. 2 illustrates the component which includes segmentation, words matching, aspect-level analysis and sentiment classification.

### B. Dataset

In this research paper can aid in observing a detail overview of data sets which are based on the proposed method and functionality. The data of dataset is the collection of reviews that are available online on websites includes implicit and explicit features. Whereas this dataset includes an electronics item which includes mobiles, laptops and cameras etc. these customers reviews are randomly collected from different online shopping websites and web portals. There are 4600 collected reviews from different domains like DVD player, mobile and camera in which analysis is made on 2024 reviews from these reviews based on this identification of the polarity of text through lexicon and sentiments analysis which are given in Table I.

### C. Aspect-Level Analysis

The methodology proposed is a novel approach that concludes aspect level analysis which helps to explore potential aspects and feature which are clearly mentioned in the sentence. The parallel sentiment expression expressed regarding aspects need the extraction at the same time of the entity. Furthermore, many organization and companies are using aspect-based analysis.

There are some of advantages of aspect-level sentiment analysis which is opinion and its targets which can be captured and provide clear sentiment information of components of a product and services or its attributes. It is similar to evaluate sentiment classification in order to estimate the performance of aspect-level sentiment analysis were each aspect is needed to be examined on behalf of human verdict.
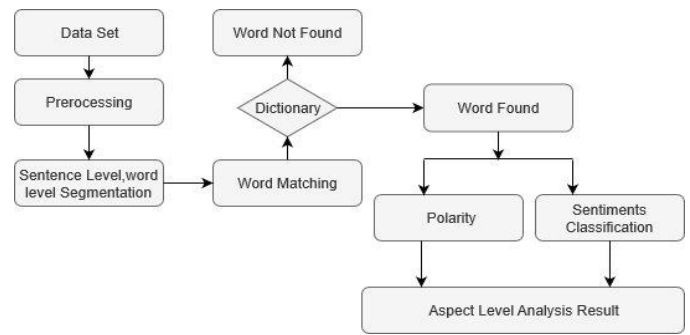


Fig. 2. Framework of System.

TABLE. I. NUMBER OF REVIEWS

| Domain | Reviews | Words | Source |
|---|---|---|---|
| DVD Player | 850 | 12980 | Amazon |
| Mobile | 600 | 9682 | Amazon |
| Camera | 654 | 11893 | Amazon |

### D. Sentence Segmentation into Words

In this process sentence segmentation is converted into words in which English word are segmented and created from sentences and word boundaries are identified and marked. The boundaries are denoted by space character before and after of each word. The process begins with identifying word boundaries after that the isolated words are listed in the list which is clearly displayed in Fig. 3. The words are retrieved and analyzed to validate words to identify from lexicon. Token is created to identify the correct spelled word this can be accomplished by removing any unnecessary characters which are part of the original words. The process of filtration involves traversing with other words and to remove any or all special characters such as @, $, %, &, #. All these characters part of string and use previously. By using these hard space characters, new paragraph and newline symbols are cropped. The filtration process which is used to eliminate occurrence of English alphabet from documents and articles as it is most useful in some situation.



{The}
{mobile}
{size}
{is}
{very}
{large}
{but}
{picture}
{quality}
{is}
{awesome}
{and}
{price}
{is}
{cheap}

Fig. 3. Sentence Tokenization in Words.

## E. Score Calculation and Word Matching

After the process of sentence segmentation and words matching the score calculation of English text the tokens are searched and matched with lexicons. In reality English text, contains multiple sentiment elements. Tokens are searched and matched with the help of lexicon after each token from English text has created. English text contains multiple sentiment elements. After the processes of segmentation and removal of stop words the remaining text is matched with lexicon and the values of each word is fetched at later stage. The value of sentiment English text can be calculated based on the corresponding calculation formulas. The text contains sentiments words and final sentiment value of English text which is calculated as follows:

Fig. 4 shows the score calculation algorithm in this algorithm the FinalValue represents the final sentiment value of the text. Sc represents the sentiment value of each aspect. The sentiment value will appear in following three states:

If FinalValue > 0, this output directs the sentiment of the text is positive.

If FinalValue = 0, it displays the sentiment of the English text as neutral.

If FinalValue < 0, it indicates the negativity of the text.

```
Input: English Text

Sentiment Polarity of Text

1.FinalValue=0;

2.A text divide into sentence n ,sentence divide into words w,

3. for(i=1;i++;i<=n)

4. { Sc=0;

5.for(j=1;j++;j<=w)

6.  {Sc=sc+0;}

7.FinalValue=FinalValue+SC;

8. If (FinalValue>0)

9. Sentiment of text is positive;

10. else If (FinalValue<0)

11.  Sentiment of text is negative;

12.  else

13.  Sentiment of text is neutral

14. end if

15. end if
```

Fig. 4.    Score Calculation Algoritm.

## IV. EXPERIMENTAL RESULTS AND ANALYSIS

In the section of experiments result and analysis the classification algorithm involves sentiment analysis in sentence tokenization which is performed by input of comments. The processing of tokenization split text into words as it is a critical stage in algorithm. However, the Sentences are tokenized into single words. The text classification method is used to extract implicit aspects with the help of word segmentation approach. The experiment in Table II illustrates the five different classification algorithms which are used to predict accuracy in the text.

TABLE. II.    EVALUATION TABLE WITHOUT USING WORDS MATCHING ALGORITHMS CLASSIFICATION METHODS

| Classifier | Precision | Recall | F-Measure |
|---|---|---|---|
| Support Vector Machine | 0.685 | 0.660 | 0.657 |
| Naivebayes | 0.744 | 0.429 | 0.735 |
| NBTree | 0.782 | 0.472 | 0.697 |
| J48 | 0.801 | 0.576 | 0.781 |
| LDA | 0.843 | 0.595 | 0.801 |

Fig. 5 and Fig. 6 illustrate the results of proposed experiment model. In Fig. 5, the results are based on random algorithm by using machine learning classifiers like SVM, Naivebayes, Nbtree, J48 and LDA to prediction the performance of text classification. It clearly illustrate that the accuracy rate and F-measure is very high in LDA as compared to others. But the performances of Recall have decrease in all classifiers measure rate. The experiment analysis shows in Table III that not every classification is good in accuracy measure.
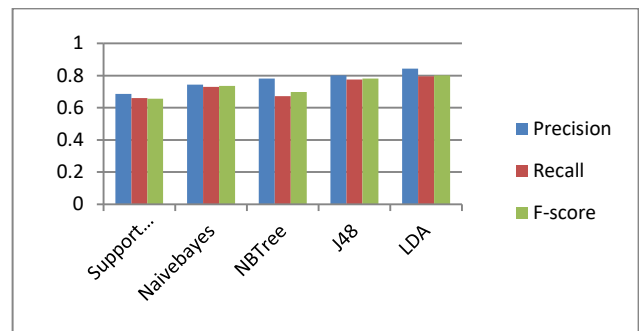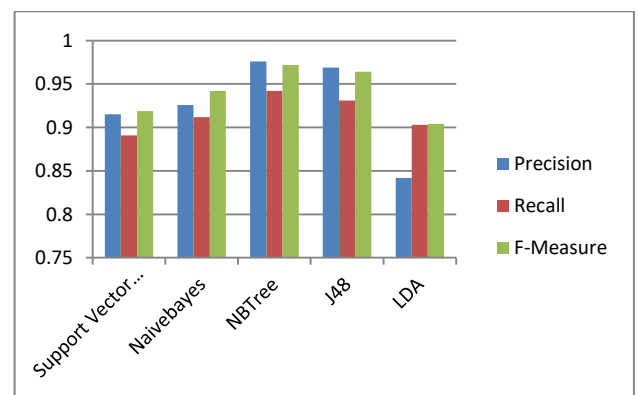


Fig. 5.    Analysis of different Classifiers.



Fig. 6.    Classifier with Framework.

TABLE. III.    EVALUATION TABLE USING WORDS MATCHING ALGORITHMS WITH CLASSIFIERS

| Classifier | Precision | Recall | F-Measure |
|---|---|---|---|
| Support Vector Machine | 0.915 | 0.891 | 0.919 |
| Naivebayes | 0.926 | 0.912 | 0.942 |
| NBTree | 0.976 | 0.942 | 0.972 |
| J48 | 0.969 | 0.931 | 0.964 |
| LDA | 0.842 | 0.903 | 0.904 |

In Fig. 6, the classification framework is analyzed by using proposed word segmentation and matching algorithm. In the proposed algorithm data is purified and send for classification. Furthermore, it clearly displays that NBtree classification have the highest accuracy in all three measures precision, recall and F-Measure. The proposed method helps text classification to increase its prediction and performance.

## V. CONCLUSION

To conclude, the analysis and experiments based on customer reviews help the different companies to improve the product quality and services. The research paper illustrates a new task of extracting implicit aspects from qualitative feedback comment from customers. The techniques like word segmentation and opinion mining are used to compare the effectiveness of different classifiers performance as well as text sentiment analysis. In the experiments section the observation has been made that NB tree classification provides better performance with F-Score. Furthermore, in future work the research will focus on recognizing the semantic polarity based on score of sentences and by applying distinct machine learning algorithm to identify positive and negative effects of the News Articles, Blogs and reviews etc.

## ACKNOWLEDGEMENT

### REFERENCES

[1] Mahar, J. A., Shaikh, H., Memon, G. Q., "A Model for Sindhi Text Segmentation into Word Tokens", Sindh University Research Journal (Science Series), Vol.44 (1) pp.43-48 (2012).

[2] Haruechaiyasak, C.; Kongyoung, S.; Dailey, "A comparative study on Thai word segmentation approaches", Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology, 2008. ECTI-CON 2008. 5th International Conference on, vol.1, no., pp.125-128, 14-17 May 2008.

[3] Nadir D. And Sarmad H, "Urdu word segmentation. In Human Language Technologies", The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics (HLT '10). Association for Computational Linguistics, Stroudsburg, PA, USA, 528-536.

[4] Sobkowicz, P., Kaschesky, M. and Bouchard, "Opinion mining in social media: Modeling, simulating, and forecasting political opinions ", Government Information Quarterly, 29(4), pp.470-479.

[5] Mohammad, S. M., Kiritchenko, S. and Zhu, X "NRC-Canada: Building the state-of-the-art in sentiment analysis of tweets", arXiv preprint arXiv:1308.6242.

[6] Maynard, D. G. and Bontcheva, "Challenges of Evaluating Sentiment Analysis Tools on Social Media", Proceedings of LREC.

[7] Liu, B. (2012), "Sentiment analysis and opinion mining', Synthesis lectures on human language technologies" 5(1), pp.1-167.

[8] Pang, B. and Lee, L. "Opinion mining and sentiment analysis", Foundations and trends in information retrieval, 2(1-2), pp.1-135.

[9] Vinodhini, G. and Chandrasekaran, R. M, "Sentiment analysis and opinion mining: a survey", International Journal, 2(6).

[10] Park, S., Lee, W. and Moon, "Efficient extraction of domain specific sentiment lexicon with active learning", Pattern Recognition Letters, 56, pp.38-44.

[11] Lin, Y., Zhang, J., Wang, X., and Zhou, "An information theoretic approach to sentiment polarity classification", In Proceedings of the 2nd Joint WICOW/AIRWeb Workshop on Web Quality, pp. 35-40. ACM.

[12] Pang, B. and Lee, "Opinion mining and sentiment analysis", Foundations and trends in information retrieval, 2(1-2), pp.1-135.

[13] Turney, "Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews", In Proceedings of the Association for Computational Linguistics (ACL), pp. 417–424.

[14] Hu, M., and Liu, "Mining and summarizing customer reviews", In R. Kohavi, J. Gehrke, W. DuMouchel, & J. Ghosh (Eds.), KDD-2004 - Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 168-177.

[15] Hu, M., and Liu, "Mining and summarizing customer reviews", In Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 168-177.

[16] Ding, X., Liu, B., and Yu, P. S, "A holistic lexicon-based approach to opinion mining", In Proceedings of the 2008 International Conference on Web Search and Data Mining, pp. 231-240. ACM.

[17] Marjan Van de Kauter ." The good, the bad and the implicit: a comprehensive approach to annotating explicit and implicit sentiment", Lang Resources & Evaluation (2015) 49:685–720.

[18] Alexandra Balahur , Jesús M. Hermida, Andrés Montoyo," Implicit Polarity and Implicit Aspect Recognition in Opinion Mining", Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, pages 20–25, Berlin, Germany, August 7-12, 2016. c 2016 Association for Computational Linguistics.

[19] C acilia Zirn, Mathias ,Nieper," Fine-Grained Sentiment Analysis with Structural Features", Proceedings of the 5th International Joint Conference on Natural Language Processing, pages 336–344, Chiang Mai, Thailand, November 8 – 13, 2011.

[20] Alexandra Balahur, Jesús M. Hermida, Andrés Montoyo, "Detecting implicit expressions of emotion in text: A comparative analysis".

[21] Amir Hossein Yazdavar a, Monireh Ebrahimi a, Naomie Salim ," Fuzzy Based Implicit Sentiment Analysis on Quantitative Sentences", Journal of Soft Computing and Decision Support Systems 3:4 (2016) 7-18.