# Sentiment Analysis Based on Multiple Reviews by using Machine learning approaches

Stephina Rodney D'souza

M.E. Student
Department of Computer Engineering
St. Francis Institute of Technology
Mumbai, India
step.6dsouza@gmail.com

Kavita Sonawane

Professor
Department of Computer Engineering
St. Francis Institute of Technology
Mumbai, India
kavitasonawane@sfitengg.org

*Abstract*—Sentiment is an attitude, thought, or judgment prompted by feeling. Sentiment Analysis can be defined as the process of analyzing online pieces of writing to determine the emotional tone they carry. With the vast growth of the social media content on the Internet in the past few years, people now express their opinion on almost anything in discussion. With respect to this, Bag–of–Words (BoW) is the most popular way to model text in statistical machine learning (ML) approaches. However, the performance of BoW sometimes remains unlimited due to some fundamental deficiencies in handling the polarity shift problem and other few challenges like quality of the opinions, hidden state representations, polarity categorization etc. To come across these challenges our focus will be on Dual Sentiment Analysis which processes the Sentiment with all the perspectives (positive, negative or neutral). This may lead towards the accurate prediction for final decision making based on the reviews given by the customers. The proposed work is being experimented on the Amazon Product reviews specifically the Mobile device reviews. This work aims at overcoming the limitation of existing system and improving the accuracy.

*Keywords— Amazon customer reviews, classification, Dual Sentiment Analysis (DSA), training and prediction.*

## I. INTRODUCTION

Sentiment analysis can be defined as the process in which it identifies and extracts the subjective information of a particular source material which helps a business to understand the social sentiment of their brand, product or service while monitoring online conversations. It aims to understand the attitude of a speaker, or writer related to a particular topic, document or an event.

In today's era, buyers intend to do online shopping as it not only saves time but also helps to view the different products available on various shopping websites owing to its variety of options, low cost value and quick supply systems [1]. People often gaze over the products and reviews before buying the product on amazon itself and other shopping website [2]. Due to the rapid increase of World Wide Web (WWW), people often express their sentiments over internet through social media, blogs, ratings and reviews. Hence, there is a need to analyze the concept of expressing sentiments and calculate the insights for exploring business [3].

After a user purchases a particular product, he/she provides reviews based on their experience which are useful for future customers that seek opinions to support them in order to take up a particular decision regarding the product [4]. Product reviews available online are significantly utilized in mining opinions as customers rely heavily on learning the sentiments indicated in the text. The concept of sentiment analysis or sentiment classification has been on the rise since 2000 (Liu, B. 2012), where its goal is to evaluate the text in accordance to the sentimental polarities of the users' thoughts or opinions, e.g., positive or negative which is generally present in the form of unstructured data [5].

Data mining tools and algorithms are utilized to discover and analyze the sentiments and attitudes of consumer behavior on products they have purchased or want to buy (Jack, L., & Tsai, 2015) [7]. The BoW model impacts the performance of machine learning based systems under the polarity shift circumstances. Therefore, it is important to improvise the performance and execution of machine learning models to give more accurate results.

## II. RELATED WORK

A lot of research has been carried out for performing the sentiment analysis on opinion mining for online reviews, Mohan Kamal Hassan, Sana Prasanth Shakthi and R Sasikala (2017) [1], and Abinash Tripathya, Ankit Agrawalb (2015) [2], the classification of the product review is performed by tagging the keyword as well as analyzing the fundamentals of determining, positive and negative approach towards the product. The results show the categorization of positive and negative comments, as well as focus on the irrelevancy and repetition of characters that are eliminated.

Both the Basari, A. S. H., Hussin, B (2013) [9], uses the concept of opinion mining and classification of text using feature extraction and hybrid approach of Support Vector Machine (SVM) and Particle Swarm Optimization (PSO) for sentiment analysis, where it increases the data as well as the size and complexity [3]. The result obtained using the dual optimization problem is that by using the Particle Swarm Optimization (PSO) the accuracy is 77%, whereas the actual rate for SVM classifier should be 97%. [5].

The authors Pang, B., & Lee, L. (2008) [10] have explained that sentiment analysis is deeply involved in research. It signifies the techniques that can be used to

understand the new difficulties and circumstances because of opinion mining when contrasted with conventional fact-based analysis. This paper also talks about natural language processing and its various functionalities like text processing, information extraction, document classification and many more. It uses the different concepts of text processing such as lemmatization, word segmentation, tokenization, normalization, and stemming. Various constraints that have come accrossed by the authors is that it leads to revealing of more information to the user, which thus leads to users interference leading to plagiarism.

### A. Threats and vulnerabilities in existing systems

During the process of opinion mining and classification of the product review by tagging the keyword, the major concern is regarding the quality of the product, and coming across the irrelevancy of the product [1].

In order to carry out the working by reviewing the dataset where the pre-processing is performed, after which the vectorization which includes the TD-IDF and classification is carried out, the reputation of characters is eliminated thus leading to loss of meaningful information [2].

To carry out the process of feature extraction, the authors need to have an understanding about polarity categorization and opinion mining. The challenge faced is lacking a collection of Opinion data viz, its objective or subjective. Also the major concern is that if we increase the data it leads to increase in size and complexity [3].

### B. Challenges to be addressed

Based on the review of literature for sentiment analysis, we have addressed the following challenges by the system.

- Accuracy of a particular data [1][2].

- Capturing the high level semantic meanings behind each text data.

- The quality of the opinions cannot be guaranteed: such as meaningless, fake, irrelevant opinions etc [8][9].

- Polarity Shifting and Categorization [4].

- Polarity Shift elimination- which affects the classification performance of machine learning based on sentiment analysis systems [2][3].

## III. PROPOSED SOLUTION

With the intention of achieving the desired accuracy along with addressing the existing challenges mentioned above (B), the proposed solution will be focusing on generation of reviews which are trustworthy which will guide the customer in accurate decision making.

The following Fig. represents on various phases (outline) of the proposed system [2][3].
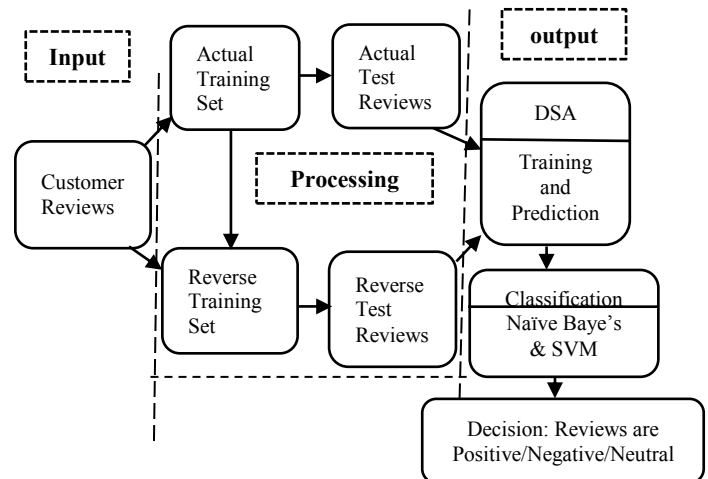


Fig. Proposed system: Dual Sentiment Analysis (DSA) Based on Multiple Reviews

The system will consider the set of reviews from any shopping website such as Amazon, Flipkart etc. for classifying the results as positive and negative using the classification algorithms such as Support Vector Machine (SVM), Naïve Baye's etc. Using the property of sentiment classification, the system will perform the data expansion model which will create reversed reviews opposite to that of the actual reviews for every training and test review. By considering the actual and the reversed reviews, we can apply the Dual Sentiment Analysis (DSA) algorithm [6][12] for training and prediction. We can also use the classifier by expanding all the combined probabilities of the actual and the reversed training dataset. The Prediction algorithm can make predictions by assessing both the sides of each review, i.e., how positive or negative the actual review is, and vice versa in the reversed review will be assessed [16].

### A. Methodology: Phases with Algorithms and Implementation details

**1. Data Preprocessing**

**a. Sentence Tokenization:** We will split a text into tokens which includes words, numbers and symbols, and eliminate the whitespaces, punctuation marks from the sentences and individual words.

**b. Stop Word Removal:** The words that will be frequently appearing in a sentence/review will be removed eg. "too", "more", "very", etc.

**c. Stemming:** The Porter Stemmer algorithm will be applied to perform the text normalization to convert a text to its standard form, expand abbreviations, and so on.

## 2. DSA with Data Expansion techniques

### a. Algorithmic Steps for Dual Sentiment Analysis

i. Download the Amazon Product Dataset and consider the reviews of a particular data.

*for example: Mobile phone data/reviews.*

ii. Check for all the reviews that are available for a particular product.

iii. Using the Dual Sentiment Analysis (DSA) perform the reverse on the original reviews.

iv. Compare the original and reverse reviews and apply DSA in which the training and prediction will be carried out.

v. Apply the Classification algorithms for comparing the results generated by the algorithms to improve the accuracy.

By utilizing all the actual reviews in the data expansion technique in the training model, we may observe that all the actual reviews may not have a definite polarity. Hence, to handle this inconsistency, a new data expansion technique should be proposed that will select only a certain segment of actual reviews that generate the desired output.

i. *"The phone's features are very interesting and the cost is low. I love it."*

ii. *"The phone's features are somewhat interesting though the cost is high. I don't hate it."*

The first example has a very strong sentiment with a low polarity shift rate. Hence this review has a definite sentiment, and its reversed review will also have a definite one as it helps in removing the ambiguity from the sentence by using the new data expansion technique i.e. Dual Sentiment Analysis (DSA).

### b. Data Expansion Technique

Based on the following rules for each actual review and an antonym dictionary, the reversed review is generated.

**Sentence Word Reversion:** Based on the antonym and sentiment classification, we will perform the reversing of the text review.

**Reversing Label:** This will help to identify the positive/negative/neutral.

*For example,*

*Sentence: "The phone's features are very interesting and the cost is low. I love it."*

*Sentence Label: Positive*

*Applying Rule 1 i.e. antonym to this example "interesting" becomes "boring", "low" becomes "high" and "love" becomes "hate".*

### c. Dual Sentiment Analysis

**Training Model:** The actual training reviews in the training stage will be reversed to their opposites by generating a newly reversed review dataset which can be referred as "actual training set" and "reversed training set". The actual and the reversed reviews will be in correspondence with each other in the data expansion technique. We will prepare a training classifier by collecting all possible combinations of the probabilities in the actual training and reversed reviews datasets.

**Prediction Model:** Both the actual and the reversed datasets will be used for predicting the final result. The prediction model can also predict how positive and negative are the actual and reversed reviews.

*For positive review, the syntax is as follows:*

```
posdata = []
with open('data/original_file.txt', 'r',encoding="utf-8") as myfile:
```

*For negative review,*

```
negdata = []
with open('data/reverse_file.txt', 'r',encoding="utf-8") as myfile:
```

### d. Classification Algorithms

**SVM (Support Vector Machine):** It is a supervised machine learning algorithm which can be used to perform both the classification and regression techniques. It works by mapping the data, and finds an optimal boundary between the possible outputs.

**Naïve Baye's:** The categorization of the text judges the documents or reviews in a system. It is a family of probabilistic algorithms that take the advantage of probability theory and Baye' Theorem to predict the tag of a text.

```
for cl in classifier_list:
    if cl == 'svm':
        classifierName = 'SVM'
        classifier = SklearnClassifier(LinearSVC(), sparse=False)
        classifier.train(trainfeats)

    else:
        classifierName = 'Naive Bayes'
        classifier = NaiveBayesClassifier.train(trainfeats)
```

## IV. RESULTS AND DISCUSSION

### A. Experimentation details:

The database used is the Amazon Product data which consists of different datasets such as Clothes (Indian or Western), Books, Furniture's etc. which provides details regarding the customers, their reviews, and customer id based on a particular product.

The Training set is used for generation of reviews/text and testing is used to check the overall performance of the system. In machine learning algorithms, these can be used for decision making or predictions of a particular review.

### B. Performance Evaluation Parameters:

To evaluate the performance of proposed system, various standard classifier evaluation metrics such as precision, recall, f-measure can be used. These are explained as follows [2]:

Accuracy can be calculated as the ratio of number of correctly predicted reviews to the number of total number of reviews present in the corpus. It can be also defined as Precision.

Precision: It gives the exactness of the classifier. The equation is as follows: ....(1)

$$\frac{True\ Positive\ (TP)}{True\ Positive(TP) + False\ Positive(FP)}$$

The Ideal case value for Precision is 0.87 and Worst case value is 0.80.

Recall: It can be defined as the process in which it measures the completeness of the classifier.
The equation is as follows: …(2)

$$\frac{True\ Positive(TP)}{True\ Positive(TP)\ +\ False\ Negative(FN)}$$

The Ideal case value for Recall is 1 and Worst case value is 0.

F-measure: It can be defined as the harmonic mean of precision and recall.
The equation is as follows: ….(3)

$$\frac{2\ x\ Precision * Recall}{Precision + Recall}$$

The Ideal case value for F-measure is 1 and Worst case value is 0.

### C. Result obtained: Analysis

1. Once the reviews are obtained from the Amazon Product dataset, we then perform the splitting of a particular text/review by using;

$splitter = r"!|\backslash?|\backslash.\{3\}|\backslash.\backslash D|\backslash.\backslash s|,|and\backslash s|but"$

This is carried out so that we remove all the unwanted data like the blank-spaces, special characters etc present in the review/statement.

2. Next, the original reviews are considered and the DSA is performed to reverse the original set of reviews regarding a particular product.

For example,
**Original Review:**
a) They look good and stick good! I just don't like the rounded shape because I was always bumping it and Siri kept popping up and it was irritating.
b) They stick on great and they stay on the phone. They are super stylish and I can share them with my sister. :)
c) These are awesome and make my phone look so stylish! I have only used one so far and have had it on for almost a year! Great quality!

**On Reversing the Original Review:**

a) They look evil stick evil I just don't dislike the rounded shape because I was always bumping it was soothe.
b) They are super styleless.
c) Make my phone look so styleless was in imperfect condition.

3. On generating the reverse from the original ones, we will now determine the equations and threshold for the reviews.

$negfeats = [(featx(f), 'neg')\ for\ f\ in\ word\_split(negdata)]$
$posfeats = [(featx(f), 'pos')\ for\ f\ in\ word\_split(posdata)]$

$negcutoff = int(len(negfeats)*3/4)$
$poscutoff = int(len(posfeats)*3/4)$

The overall strength of the positive and negative review is computed using the following eaqns,
trainfeats=negfeats[:negcutoff]+posfeats[:poscutoff]
testfeats=negfeats[negcutoff:] + posfeats[poscutoff:]

**Determination of Threshold:**
Threshold helps us to calculate the positive_precision and negative_precision whose value ranges from 0 to 1 in a given data.

In this work to segregate the reviews based on the sentiments we apply the following threshold. In order to determine the class labels for positive and negative reviews

we have set the threshold to 50%, and it is applied as follows:

*Positive sentiments >= 0.5;* then it falls under the Positive Class.

*Negative sentiments <= -0.5;* it falls under the Negative Class.

*Neutral > -0.5 < 0.5 ;* if it lies in between then it falls under the neutral class.

Using the equations (1), (2), (3) we get the following performance evaluation parameters. They are as follows:

Table No. 1

| Para-meters / Algorithms | Accuracy | Positive precision | Negative precision | Precision | Recall | F-measure |
|---|---|---|---|---|---|---|
| Naïve Baye's | 0.668 | 0.60 | 1.0 | 0.800 | 0.66 | 0.62 |
| SVM | 0.917 | 0.95 | 0.88 | 0.92 | 0.91 | 0.91 |

Based on the experimentation of the proposed method and analysis of results we can say that the proposed work has a good impact on accurate prediction for final decision making based on the reviews given by the customers. It also helps us in understanding the following:

- Using the preprocessing technique such as Sentence Tokenization, Stemming, Stop-word removal etc. will help to refine the reviews in the desired format that leads to reduction in ambiguity and computing complexities for further processing.
- The Dual Sentiment Analysis (DSA) provides a clear idea of whether the sentiment is positive or negative to the customer as it removes the ambiguity, solves the polarity shift problem and provides clarity of the sentiment to the user. It also helps to overcome the loss of meaningful information which happens during the Stop-word removal process. Although the reversing increases the complexity but it provides an accurate desired result for decision making purpose.
- Depending upon the positive_precision, negative_precision and threshold values obtained, Naïve Baye's and SVM (Support Vector Machine) algorithms are used to classify the reviews into positive/negative/neutral as it helps in training and predicting the reviews that are generated by DSA. On comparing the two algorithms in the above Table No.1, we conclude that the accuracy of SVM is comparatively better (0.91%) than as compared to Naïve Baye's (0.66%).

## V. SIGNIFICANCE OF DUAL SENTIMENT ANALYSIS

DSA has a wide range of applications in Online shopping as well as social media sites as it helps a customer to take up a particular decision by providing the accurate solution to the user's query.

It can be applied in Online Banking systems, Schools and Colleges Admission process as its main focus is to compare all the feedbacks/reviews using DSA and generate an accurate solution for the same.

## VI. CONCLUSION

The research study of implementing the Dual Sentiment Analysis (DSA) seems to be very attractive and interesting in the field of machine learning, as it helps to overcome the challenges such as the BoW model, Polarity shift categorization problems etc.

Experimentation of the proposed work with Amazon product data specifically the Mobile device reviews with the intention of overcoming the limitations of existing system and improving the accuracy. It has proved the same in the results obtained. Based on the result and analysis we can conclude that DSA contributed positively in final decision making process that would help the customer to gain the correct insights about the reviews; in turn to choose or recommend the correct product.

## REFERENCES

[1] M K Hassan, S P Shakthi and R Sasikala, **"**Sentimental analysis of Amazon reviews using naïve bayes on laptop products with MongoDB and R", *School of Computer Science and Engineering, VIT University*, Vellore-632014, India, 2017.

[2] A Tripathya, A Agrawalb, SK Rath, "Classification of Sentimental Reviews Using Machine Learning Techniques", *3rd International Conference on Recent Trends in Computing,* pp 821-829, 2015.

[3] A Bhatt, A Patel, Harsh CK Gawande, "Amazon Review Classification and Sentiment Analysis", *(IJCSIT) International Journal of Computer Science and Information Technologies*, Vol. 6, no. 6, pp 5107-5110, 2015.

[4] Efstratios "Sentiment Polarity Detection From Amazon Reviews: An Experimental Study", *Sygkounas, Giuseppe Rizzo, Raphael Troncy*, 2016.

[5] Liu B Sentiment Analysis and Opinion Mining, "Synthesis Lectures on Human Language Technologies", *Morgan & Claypool Publishers,* 2012.

[6] El-Din, D. M., "Enhancement bag-of-words model for solving the challenges of sentiment analysis", *International Journal of Advanced Computer Science and Applications*, Vol 7, pp 244-247, 2016.

[7] Jack. L. & Tsai. Y. D., "Using Text Mining of Amazon Reviews to explore User-Defined Product Highlights and Issues", *In Proceedings of International Conference on Data Mining*, January 2015.

[8] Dragoni, M., Tettamanzi, A., Pereira, C, "DRANZIERA: An Evaluation Protocol For Multi-Domain Opinion Mining", *10th International Conference on Language Resources and Evaluation (LREC),* 2016.

[9] Basari, A. S. H., Hussain, B., Ananta, I. G. P., & Zeniarja J., "Opinion Mining of Movie Review using hybrid method of SVM-PSO", *Cioroianu. L., NLP*, pp 453-462, 2013.

[10] Ikeda, D., Takamura, H., Ratinov, L. A., & Okumura, M., "Learning to shift the polarity of words for sentiment classification", In Proceedings of the Third International Joint Conference on Natural Language Processing, Vol 1, 2008.

[11] Lincy, W., & Kumar, N., "A Survey on Challenges in Sentiment Analysis", *International Journal of Emerging Technology in Compute. Journal of Big Data, Vol 2, no.5, 2015.*

[12] Fang, X., Zhan, J. "Sentiment analysis using product review data*",* Science and Electronics, Vol 21, no.3, pp 409-412, 2016.*

[13] Rui Zhao and Kezhi Mao, "Fuzzy Bag-of-Words Model for Document Representation*", IEEE Transcations on Fuzzy Systems*, Vol. 26, no. 2, April 2018.

[14] R. Zhao and K. Mao, "Topic-aware deep compositional models for sentence classification", *IEEE/ACM Trans. Audio, Speech, Language Process.*, Vol. 25, no. 2, pp 248–260, Feb. 2017.

[15] G Angulakshmi, Dr.R Manicka Chezian, "An Analysis on Opinion Mining: Techniques and Tools", Vol 3, no.7, 2014.

[16] S. ChandraKala1 and C. Sindhu, "OPINION MINING AND SENTIMENT CLASSIFICATION: A SURVEY," *ICTACT journal,* Vol. 3, no.1, pp 420-427, Oct 2012.