

# Sentiments Detection and Amazon Product Review

Y.BHAGYA LAKSHMI<sup>1</sup>, Dr.V.Bhaskara Murthy<sup>2</sup>

<sup>1</sup>MCA Student, B V Raju College, Kovvada, Andhra Pradesh, India.

<sup>2</sup>HOD & Professor, B V Raju College, Kovvada, Andhra Pradesh, India.

## ABSTRACT:

Today, digital reviews play a pivotal role in enhancing global communications among consumers and influencing consumer buying patterns. E-commerce giants like Amazon, Flipkart, etc. provide a platform to consumers to share their experience and provide real insights about the performance of the product to future buyers. In order to extract valuable insights from a large set of reviews, classification of reviews into positive and negative sentiment is required. Sentiment Analysis is a computational study to extract subjective information from the text. In the proposed work, over 4,000,00 reviews have been classified into positive and negative sentiments using Sentiment Analysis. Out of the various classification models, Naïve Bayes, Support Vector Machine (SVM) and Decision Tree have been employed for classification of reviews. The evaluation of models is done using 10 Fold Cross Validation.

**Keywords:** *ML, Random forest, SVM, Naive bayes, high accuracy.*

## 1. INTRODUCTION

With an ever increasing demand of smart phones, the mobile phone market is expanding at an exponential pace. With such a boom in the smart-phone industry, there is a need to realize the holistic review of the brand and the model of phone. There are numerous brands present in the market, out of which some are dominant and occupy quite a big part of the industry. For instance, Samsung, Apple, etc. are names associated with brands which are famous throughout the world. Electronic commerce plays a vital role in increasing the sales of the mobile phones and

influencing consumer buying patterns. Reviews available on such e-commerce platforms act as a guiding tool for the consumers to make informed decisions. Retail websites like Amazon.com offer different options to the reviewers for writing their reviews. For instance, the consumer can provide numerical rating from 1 to 5 or write comments about the product.

As there are innumerable products manufactured by many different brands, so providing relevant reviews to the consumers is the need of hour. Number of reviews associated with a product or a

brand is increasing at an alarming rate, which is no less than handling the big data. Classifying the reviews on the basis of sentiment of customers into positive and negative sentiment provides sentiment orientation of the review, hence results in better judgement. Segregation of reviews on the basis of their sentiment can help future buyers to evaluate positive and negative feedback constructively and reach at better decisions as per their requirements. This evaluation acts as a testimony to the users who are looking to know the details and specifications of the smartphones; thereby increasing user credibility. In this research, unstructured data of Mobile Phone Reviews have been extracted from Amazon.com. It has been filtered to remove noisy data and has been pre-processed to evaluate sentiment of the reviews using supervised learning. The reviews have been classified using machine learning classification models like Naïve Bayes, Support Vector Machine (SVM) and Decision Tree and have been cross validated to find the best classifier for this purpose.

## 2. RELATED STUDY

Data analytics has enabled to unravel the hidden patterns in data. Volume, Velocity and Variety define Big Data [1]. Veracity and Value are two more Vs that play an important role in Big Data. The volume

and the relentless rapidity at which data are being generated every day are exceeding the computing capacity of many IT departments. E-commerce websites are loaded with a large set of diverse reviews for various products. These reviews can be used to determine consumer behavior and make informed decisions. Reviews can be both structured and unstructured. Valuable business insights can be fetched by filtering the irrelevant data. Big Data has enabled businesses to flourish and improvise on the basis of evidence rather than intuition. It aids in gaining insights on better targeted social influencer marketing, segmentation of customer base, recognition of sales and marketing opportunities, detection of fraud, quantification of risks, better planning and forecasting, understanding consumer behavior, etc. [2]. Sentiment analysis implies identifying sentiment of reviews on the basis of positive, negative and neutral connotations. Sentiment analysis can be performed at three levels, viz.

document level, sentence level and phrase level [3]. A lot of prior research has been done in this field where words and phrases have been classified with prior positive or negative polarity [4]. This prior classification is helpful in many cases but when contextual polarity comes into the picture, the meaning derived from positive

or negative polarity can be entirely different. The contextual polarity of the phrases was taken into consideration and ambiguity was removed [5]. Also, a refined method has been devised to establish contextual polarity of phrases by using subjective detection that compressed reviews while still maintaining the intended polarity [6]. Delineated study has been conducted on tweets available on Twitter, movie reviews to build the grounds on sentiment analysis and opinion mining. A sentiment classifier has been built to categorize positive, negative and neutral sentiments not only in English but also for other languages using corpus from Twitter [7]. The polarity of smartphone product reviews has been found only on the basis of positive and negative orientation of the review [8]. A system has been built using support vector machine where sentiment analysis is carried out by taking into consideration sarcasm, grammatical errors and spam detection [9]. An enhanced Naïve Bayes model by combining methods like effective negation handling, word ngrams and feature selection has been utilized to conduct sentiment analysis [10]. Sentiment analysis is not only confined to the English language but has been implemented for various languages. Sentiment analysis of Chinese text by implementing four feature

selection methods and five classifiers viz. Centroid classifier, K-nearest neighbor, Window classifier, Naïve Bayes and SVM has been done [11]. Through this learning paradigm it was concluded that SVM outperforms all the other learning methods in terms of sentiment classification. Sentiment analysis on travel reviews using three machine learning models namely, Naïve Bayes, SVM and character based N-gram model has been performed in which SVM and N-gram approaches have better performance than Naïve Bayes [12]. It has been observed that in maximum number of cases SVM showcases best performance in comparison to other classification models.

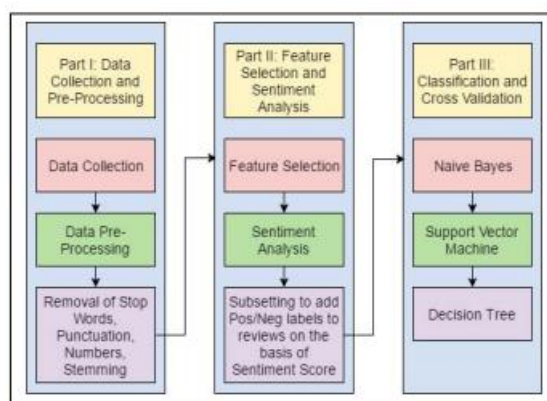
### EXISTING SYSTEM

Number of reviews associated with a product or a brand is increasing at an alarming rate, which is no less than handling the big data. Classifying the reviews on the basis of sentiment of customers into positive and negative sentiment provides sentiment orientation of the review, hence results in better judgment.

### 3 PROPOSED SYSTEM

Sentiment analysis is not only confined to the English language but has been implemented for various languages. Sentiment analysis of Chinese text by implementing four feature selection methods and five classifiers viz. Centroid

classifier, K-nearest neighbor, Window classifier, Naïve Bayes and SVM has been done [11]. Through this learning paradigm it was concluded that SVM outperforms all the other learning methods in terms of sentiment classification. Sentiment analysis on travel reviews using three machine learning models namely, Naïve Bayes, SVM and character based N-gram model has been performed in which SVM and N-gram approaches have better performance than Naïve Bayes [12]. It has been observed that in maximum number of cases SVM showcases best performance in comparison to other classification models



The proposed framework of the research work is conducted in three different modules as shown in Fig. 1.

### A. Dataset and its Features

The first module includes data collection and pre-processing of data. A large sample of online reviews is collected from the e-commerce giant Amazon.com. The data set consists of over 400,000 reviews for approximately 4500 mobile

phones. It includes six features as explained in table 1.

### B. Approach

The approach followed by the proposed framework is described in Fig. 1. Initially, the experimental data is collected from an e-commerce website Amazon.com. Each data set is in the Comma Separated Values (CSV) file format and available as supplement. In the second step, data are pre-processed to remove stop words, punctuation marks, whitespaces, digits and special symbols. 'tm' package [19] is employed for text mining. In the third step, feature selection is performed to extract relevant features from the data set. In the given data set out of the six features, only three features, i.e., Product Name, Brand Name and Reviews have been considered. In the fourth step, sentiment orientation of the reviews is determined. In the fifth step, 'Pos/Neg' tags are appended to the dataset to corresponding to each review to conduct supervised learning. The sixth step involves training and testing the classified data using Naïve Bayes, SVM and Decision Tree models. The accuracy so obtained is validated using 10-fold cross validation.

### ANALYSIS:

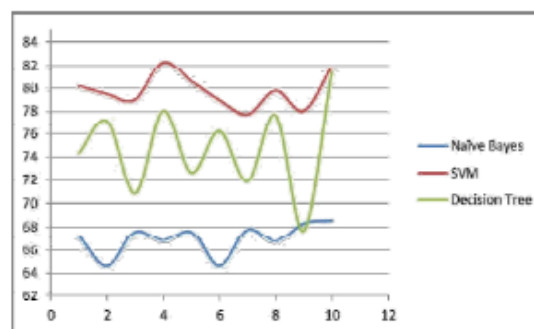
An inbuilt package named 'Syuzhet' [20] has been used to conduct

Sentiment Analysis. This package encompasses three sentiment dictionaries. NRC sentiment dictionary is used to extract eight different emotions and their corresponding valence in the text including all the reviews. Ten different emotions represented are: anger, anticipation, disgust, fear, joy, sadness, surprise, trust, positive and negative.

Feature	Description
Product Name	Model name of mobile phone
Brand Name	Manufacturing brand
Price	Price of the mobile in dollars
Rating	User rating between 1 to 5
Reviews	User reviews provided for every mobile phone
Review Votes	Number of people who found the review helpful

Represents that the number of positive reviews is more than double the negative reviews. This implies that the data are imbalanced as the target variable has imbalanced proportion of classes. Therefore, running machine learning models for classification would yield biased predictions and misleading accuracies. To avoid such scenarios, data balancing is employed. There are different methods that can be used to transform imbalanced data into balanced data like under sampling and oversampling. For treating imbalanced data, under sampling technique has been used. Under sampling means to reduce the number of observations from majority class to balance the data set. The balanced data so

obtained has almost equal number of positive and negative reviews.



**Fig.1. Scatter Plot.**

#### 4. CONCLUSION

An evolutionary shift from offline markets to digital markets has increased the dependency of customers on online reviews to a great extent. Online reviews have become a platform for building trust and influencing consumer buying patterns. With such dependency there is a need to handle such large volume of reviews and present credible reviews before the consumer. Our research is aiming to achieve this by conducting sentiment analysis of mobile phone reviews and classifying the reviews into positive and negative sentiment. After balancing the data with almost equal ratio of positive and negative reviews, three classification models have been used to classify reviews. Out of the three classifiers, i.e., Naïve Bayes, SVM and Decision Tree, predictive accuracy of SVM is found to be the best. The accuracy results have been cross validated and the highest value of accuracy

achieved was 81.75% for SVM among the three models. In future, the work can be extended to perform multiclass classification of reviews which will provide delineated nature of review to the consumer, hence better judgement of the product. It can also be used to predict rating of a product from the review. This will provide users with reliable rating because sometimes the rating received by the product and the sentiment of the review do not provide justice to each other. The proposed extension of work will be very beneficial for the e-commerce industry as it will augment user satisfaction and trust.

## REFERENCES

- [1] S. Erevelles, N. Fukawa, and L. Swayne, "Big data consumer analytics and the transformation of marketing," *Journal of Business Research*, vol. 69, no. 2, pp. 897–904, 2016.
- [2] P. Russom et al., "Big data analytics," TDWI best practices report, fourth quarter, pp. 1–35, 2011.
- [3] S. Erevelles, N. Fukawa, and L. Swayne, "Big data consumer analytics and the transformation of marketing," *Journal of Business Research*, vol. 69, no. 2, pp. 897–904, 2016.
- [4] V. Hatzivassiloglou and K. R. McKeown, "Predicting the semantic orientation of adjectives," in *Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 1997, pp. 174–181.
- [5] T. Wilson, J. Wiebe, and P. Hoffmann, "Recognizing contextual polarity in phrase-level sentiment analysis," in *Proceedings of the conference on human language technology and empirical methods in natural language processing*. Association for Computational Linguistics, 2005, pp. 347–354.
- [6] B. Pang and L. Lee, "A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts," in *Proceedings of the 42nd annual meeting on Association for Computational Linguistics*. Association for Computational Linguistics, 2004, p. 271.
- [7] A. Pak and P. Paroubek, "Twitter as a corpus for sentiment analysis and opinion mining," in *LREc*, vol. 10, no. 2010, 2010.
- [8] M. WAHYUDI and D. A. KRISTIYANTI, "Sentiment analysis of smartphone product review using support vector machine algorithm-based particle swarm optimization." *Journal of Theoretical & Applied Information Technology*, vol. 91, no. 1, 2016.
- [9] D. N. Devi, C. K. Kumar, and S. Prasad, "A feature based approach for sentiment analysis by using support vector machine," in *Advanced Computing (IACC), 2016 IEEE 6th International Conference on*. IEEE, 2016, pp. 3–8.
- [10] V. Narayanan, I. Arora, and A. Bhatia, "Fast and accurate sentiment classification using an enhanced naive bayes model," in *International Conference on Intelligent Data Engineering and Automated Learning*. Springer, 2013, pp. 194–201.