# Sentiment Analysis of Online Product Reviews by Hybridized Random Forest with Naïve Bays Algorithm

**Dr. (Mrs) Radha Pimpale[1], Dr. Rahul Khokale[2]**
Asst. Prof, Department of Information Technology,
Priyadarshini Bhagwati College of Engineering, Nagpur
Professor and Head of Department, Department of Computer Science,
Napur Institute of Engineering, Nagpur

## ABSTRACT

*This paper deals with the analysis of online product reviews.. Data used in this study are the online product reviews collected from different online web sites. These online reviews are extremely useful to new customer for buying product, to make a decision whether the product to be purchased is worth or not, to develop market strategies etc. In this paper, we have extracted positive, negative and neutral sentiments about the online product from the sentences. The classification techniques used for categorization are Support Vector Machine (SVM), Naïve Bayesian, Random Forest (RF) and Hybrid approach is considered for sentiment analysis.*

## Keywords

Sentiment analysis, online product review, Support Vector Machine, Naïve Bayesian, Random Forest

## 1. Introduction

Sentiment analysis is the computational task of automatically determining what feelings, a writer is expressing in text about the product. Sentiment is an feeling and thought. Sentiment analysis involves in mining the naturally expressed text to understand the feeling of people towards the interested online product. Sentiment mining and analysis has found many application in areas of Marketing, crime, healthcare, tourism, fraud detection, finance etc.

Public opinion about online product most important for new customer to buy product online. It studies customers sentiments towards certain product. If customers purchase product through online, customers are able to post their own thought about product through online shopping sites so that it might be useful to the other customer.

Customers can freely post their own opinions about the quality of product. Product review indicating whether the view is positive, negative, or neutral. Some message is tagged based on the emoticons (☺as positive, ☻as negative).

Data used in this paper is a set of product reviews collected from online shopping sites. Some product review receives inspections before it can be posted and each review must have a rating of the product. The rating is based on a star-scaled system, where the highest rating has 5 stars and the lowest rating has only 1 star. This paper we categories sentiment analysis based on

different classification techniques of big data Analysis. We categories paper based on review level and sentence level.

Review is most important part to purchase online product, based on review, people can buy that product. We can also categories based on sentence level. Sentences such as nice product, worth product make importance to buy product.

An algorithm is proposed and implemented to analyzing the online product, it evaluating survey responses and determining whether product reviews are positive or negative., Performance of three classification models such as Naïve Bayesian, Random Forest, and Support Vector Machine are evaluated and compared with newly hybrid Random forest with Naïve based Algorithm results.

## 2. Related Work

Hu and Liu [4] summarized a list of positive words and a list of negative words, respectively, based on customer reviews. The positive list contains 2006 words and the negative list has 4783 words. Both lists also include some misspelled words that are frequently present in social media content.

Pang and Lee [5] suggested to remove objective sentences by extracting subjective ones. They proposed a text-categorization technique that is able to identify subjective content using minimum cut.

Gann et al.[6] selected 6,799 tokens based on Twitter data, where each token is assigned a sentiment score, namely TSI(Total Sentiment Index), featuring itself as a positive token or a negative token.

Devika M D, Sunitha C, Amal Ganesha [7] compare different classifier and compare them as n-gram evaluation is useful than word by word. . It was shown that using unigrams (a bag of words) as features in classification performed quite well with either naïve Bayes or SVM.

## 3. Data collection

Data used in this paper is a set of product reviews collected from online shopping sites such as Amazon, Flipcart. Those online reviews were posted by over millions of reviewers (customers) towards different products.

Each review includes the following information: 1) reviewer Name 2) product ID; 3) rating;

4) Date of review 6) review text. 7).Every rating is based on a 5-star scale 8) certified buyer resulting all the ratings to be ranged from 1-star to 5-star .

Data collection (a) Data based on product categories (b) Data based on review categories.

## 4.Classification Approach for Sentiment Analysis

SVM and Naïve Bayes and Random forest classifiers are used to classify. SVM is used to classify text as either positive or negative. Naïve bayes algorithm is used to classify sentiment[7]. Random Forest, is multiple decision Tree are ensemble to improve performance

## 4.1 SVM

Support Vector Machine" (SVM) is a supervised machine learning algorithm which can be used for both classification or regression analysis. SVM is mostly used in classification problems. In this algorithm, each data item is plot as a point in n-dimensional with the value of each star considered. Classification can be done by the hyper-plane that differentiate the two classes positive and negative and SVM works well with clear margin of separation. SVM is effective in high dimensional spaces.

## 4.2 Naïve Bayes

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods. Naive Bayes model is easy and fast to predict class of test data set. It also performs well in multi class prediction.

For sentiment analysis we applied naïve bayes supervised text classification algorithm. Firstly we removed punctuations, numbers, web link and whitespaces and then come up with naïve bayes gives result in table that evaluate or count the positive, negative or neutral comments. We also calculate data frequency and find which word has higher frequency.

Bayes theorem provides a way of calculating posterior probability P(c|x) from P(c), P(x) and

P(x|c). Look at the equation below:

$$P(c \mid x) = \frac{P(x \mid c) P(c)}{P(x)}$$

$$P(c \mid X) = P(x_1 \mid c) \times P(x_2 \mid c) \times \cdots \times P(x_n \mid c) \times P(c)$$

## 4.3 Random Forest Algorithm

Random forest algorithm is a supervised classification algorithm. It creates the forest with a number of trees. Random Forest algorithm (RF) or the Random Forest classifier can be used for both classification and the regression task.

In general, the more trees in the forest the more robust the forest looks like. In the same way in the random forest classifier, the higher the number of trees in the forest gives the high accuracy results. Random forest classifier will handle the missing values.

Random Forest is a type of ensemble machine learning algorithm called Bootstrap Aggregation or bagging. the Bagging ensemble algorithm and the Random Forest algorithm for predictive

modelling. We are applying the bootstrap method for estimating statistical quantities from samples and the Bootstrap Aggregation algorithm for creating multiple different models from a single training dataset. Random forest algorithm has accuracy, it runs on large dataset , it generate accurate when large proportion data is missing , generated forest is used for future used so provide accurate result. The Random Forest algorithm, the algorithm works on Bagging and to use it for predictive modelling and produce results in a very powerful classifier

## 5. Sentiment Classification Using Supervised Learning
Sentiment classification is usually formulated as a two-class classification problem, *positive* and *negative and Neutral*. Training and testing data used are normally product reviews. Since online reviews have rating scores assigned by their reviewers, e.g., 1-5 stars, the positive and negative classes are determined using the ratings. For example, a review with 4 or 5 stars is considered a positive review, and a review with 1 to 2 stars is considered a negative review. assigning all 3-star reviews the neutral class. Sentiment classification is essentially a text classification problem.

## 6. Implementation
The comments are collected from the customer in the form rating to review the online product. The implementation steps may include,

1. Retrieving the feedback (i.e.) data Collection or review data
2. Pre-processing
3. Applying different classifying techniques

Dataset The dataset contains online product reviews along with their associated binary sentiment polarity labels. The dataset is obtained online shopping website such as Flipcart, Amazon. The user can view their interested product

## 6.1 Hybrid Random Forest Algorithms

In this Hybrid algorithm , we considered feature of Random forest classifer has capability that it can handle large amount of data with missing values , it also has capability to save random forest tree for further prediction. After implementations we considered the test feature and calculated their votes in the form of two classes for performing operation of Naïve Based algorithm

## 6.2 Hybrid Algorithm

1. Randomly select "K" features from total "m" features where k << m
2. Among the "K" features, calculate the node "d" using the best split point
3. Split the node into daughter nodes using the best split
4. Repeat the 1 to 3 steps until "l" number of nodes has been reached
5. Build forest by repeating steps 1 to 4 for "n" number times to create "n" number of trees
6. Takes the test features and use the rules of each randomly created decision tree to predict the outcome and save the predicted outcome (target)
7. Calculate the values for each predicted target
8. Convert the Values into data set , put dataset into a frequency table

9. Create Likelihood table by finding the probabilities like positive probability and negative probability customer review
10. Naive Bayesian equation to calculate the posterior probability for each class. The class with the highest posterior probability is the outcome of prediction.

## 7. Software Used

Software used for this study is R, an open source machine learning software. The classification models selected for categorization are: Naïve Bayesian, Random Forest, and Support Vector Machine.

## 8. Results

Data collected from reviews that have rating 5 and 4-star ratings are labelled as positive, ratings 1-star and 2-star reviews considered as negative and 3 star considered as neutral

The evaluation measures used for result evaluation are mean , standard deviation and accuracy, precision, recall,
For the collected dataset of 10,000 data, for classification

Table I : Performance Evaluation

| Techniques Used | Precision | Recall | Accuracy |
|---|---|---|---|
| SVM | 60.21 | 65.23 | 65.78% |
| Naïve Bayes | 69.23 | 71.53% | 70% |
| Random Forest | 74.56% | 75.32% | 75% |
| Hybrid Random Forest  with Naïve Bayes | 80.53% | 83.54% | 85% |

## 9.Conclusion

In this paper, we have evaluated through naïve bayes , support vector machine and random forest algorithm as well as Hybridised naïve bays with random forest. Random forest with Naïve bayes is gives almost accurate results gives 85% accurate result, so we found that Random forest with naïve bayes supervised learning algorithm works better. A sentiment categorization process has been proposed along with detailed descriptions of each step. Experiments for both sentence-level categorization and review-level categorization have been performed.

## REFERENCES

1. Bing Liu," Sentiment Analysis and Opinion Mining", April 22, 2012, Book

2. Duc Tam Hoang, Sentiment Analysis: Polarity Dataset", March 4, 2014,Book

3. Xing Fang* and Justin Zhan, Sentiment analysis using product review data", Fang and

Zhan Journal of Big Data (2015) 2:5 DOI 10.1186/s40537-015-0015-2

4. Hu M, Liu B (2004) Mining and summarizing customer reviews In: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, 168–177.. ACM, New York, NY, USA.Google Scholar

5. Pak A, Paroubek P (2010) Twitter as a corpus for sentiment analysis and opinion mining In: Proceedings of the Seventh conference on International Language Resources and Evaluation.. European Languages Resources Association, Valletta, Malta.Google Scholar

6. Liu B (2010) Sentiment analysis and subjectivity In: Handbook of Natural Language Processing, Second Edition.. Taylor and Francis Group, Boca.Google Scholar

7. Liu B (2014) The science of detecting fake reviews. http://content26.com/blog/bing-liu-the-science-of-detecting-fake-reviews/.

8. Lin Y, Zhang J, Wang X, Zhou A (2012) An information theoretic approach to sentiment polarity classification In: Proceedings of the 2Nd Joint WICOW/AIRWeb

   Workshop on Web Quality, WebQuality '12, 35–40.. ACM, New York, NY, USA.View ArticleGoogle Scholar

9. Sarvabhotla K, Pingali P, Varma V (2011) Sentiment classification: a lexical similarity based approach for extracting subjectivity in documents. Inf Retrieval14(3): 337–353.View ArticleGoogle Scholar

10. Devika M Dª*, Sunitha Cª, Amal Ganesha (2016) , "Sentiment Analysis:A

    Comparative Study On Different Approaches", 1877-0509 © 2016 The Authors. Published by Elsevier B.V.,, ICRTCSE 2016, doi: 10.1016/j.procs.2016.05.124

11. Chetashri Bhadanea,Hardi Dalalb, Heenal Doshic," Sentiment analysis: Measuring opinions", International Conference on Advanced Computing Technologies and

    Applications (ICACTA-2015),, 1877-0509 © 2015 Published by Elsevier B.V., doi: 10.1016/j.procs.2015.03.159