

Sentiment Analysis of Amazon Product Reviews Using Machine Learning Algorithms

Dr. U. D. Prasan¹, M. Ruchita², T. Dinesh³, S. Manohar⁴, S. Prasanth⁵, Saurab⁶

¹Professor & HOD of Computer Science Engineering, Aditya Institute of Technology and Management (AITAM), Tekkali, Andhra Pradesh, India

^{2,3,4,5,6} Students, IV B. Tech CSE, Aditya Institute of Technology and Management (AITAM), Tekkali, Andhra Pradesh, India

ABSTRACT

Sentiment analysis or opinion mining is one of the major tasks of NLP (Natural Language Processing). Sentiment analysis has gained much attention in recent years. Nowadays customers wanted to purchase anything just at one click of a mouse button. Online shopping is becoming even more popular due to its high level of convenience. Online sellers and merchants ask their purchasers to share their opinions about the products they have bought. As a result, millions of reviews are generated daily, which makes it difficult for a customer to make a good decision or whether to buy the product or not. Analyzing large number of reviews is also hard and time consuming for product manufacturers. The goal of our project is to understand and analyze the Amazon user review data set to identify whether most of the reviews are positive or negative with the help of logistic regression and support vector machine classification algorithms.

Keywords – Classification, Natural Language Processing, Logistic Regression, Sentiment analysis, Support Vector Machine

Date of Submission: 01-04-2023

Date of acceptance: 11-04-2023

I. INTRODUCTION

Sentiment is an emotion or attitude prompted by the feelings of the customer. Sentiment analysis is often referred to as opinion mining, because the opinion given by the customer will be mined to reveal the polarity of the opinion. It comes under machine learning. Since the online data is tremendously growing day-by-day, it is considered to be very important. As customers express their reviews and thoughts about the brand more openly than ever before, sentiment analysis has become a powerful tool to monitor and understand online conversations. Sentiment analysis is a process where the dataset consists of emotions, attitudes or assessment which takes into account the way a human think. In a sense, trying to understand the positive and the negative aspect is a very difficult task. People are able to post their own content through various social media, such as forums, micro-blogs, or online social networking sites. Online data has several flaws. The first flaw is that since people can freely post their own content, the quality of their opinions cannot be guaranteed. The second flaw is that ground truth of such online data is not always available. A ground truth is more like a tag of a certain opinion

about indicating whether the opinion is positive or negative. The purpose of our project is to analyze and understand the overall sentiment expressed by customers about a particular product on the Amazon platform. Sentiment analysis is a technique that involves using natural language processing (NLP) and machine learning algorithms to automatically identify and extract the subjective information from textual data such as product reviews.

II. LITERATURE SURVEY

There are many works done in the area of sentiment analysis.

R. Singh and A. Singh [1] paper discusses the use of a decision tree algorithm for sentiment analysis of Amazon product reviews. The authors first discuss the importance of sentiment analysis in today's world and how it can be used to make informed decisions. They then explain the dataset used in the study, which consisted of 5000 Amazon product reviews. The dataset was preprocessed to remove stop words and perform stemming. The results of the study showed that the decision tree algorithm achieved an accuracy of 77.6%, with precision, recall, and F1 score of 0.78, 0.77, and 0.77 respectively.

tively. The authors finally conclude that decision tree algorithm is a useful tool for sentiment analysis of Amazon product reviews and can be used to make informed decisions about products.

S. S. Shrivastava and S. K. Shukla [2] paper focuses on the use of the random forest algorithm for sentiment analysis of Amazon product reviews. The authors begin by introducing the importance of sentiment analysis in today's world and how it can be used to make informed decisions. They then discuss the dataset used in the study, which consisted of 1000 Amazon product reviews. The dataset was preprocessed to remove stop words, perform stemming, and convert the text into a numerical representation using the bag-of-words model. The results of the study showed that the random forest algorithm achieved an accuracy of 84.8%, with precision, recall, and F1 score of 0.87, 0.85, and 0.84 respectively.

S.S.Reddy and V.Krishnaiah [3] paper discusses the use of the K-nearest neighbor (KNN) algorithm for sentiment analysis of Amazon product reviews. The results of the study showed that the KNN algorithm achieved an accuracy of 82.3%, with precision, recall, and F1 score of 0.83, 0.82, and 0.82 respectively. The authors finally conclude that the KNN algorithm is a useful tool for sentiment analysis of Amazon product reviews and can be used to make informed decisions about products.

A. Mukherjee and S. K. Das [4] paper focuses on the use of the Naive Bayes algorithm for sentiment analysis of Amazon product reviews. The authors begin by discussing the importance of sentiment analysis and its applications in various fields. They then describe the dataset used in the study, which consisted of 1000 Amazon product reviews. The dataset was preprocessed to remove stop words and perform stemming. The results of the study showed that the Naive Bayes algorithm achieved an accuracy of 81.7%, with precision, recall, and F1 score of 0.82, 0.82, and 0.81 respectively. The authors conclude that the Naive Bayes algorithm is a useful tool for sentiment analysis of Amazon product reviews and can be used to make informed decisions about products.

M.H.Rahman, M.M.Rahman, and M.H.Kabir [5] paper presents a comparative study of three popular machine learning algorithms, namely K-Nearest Neighbor (KNN), Decision Tree, and Random Forest, for sentiment analysis of Amazon product reviews. The authors start by introducing the importance of sentiment analysis and its applications in

various fields. They then describe the dataset used in the study, which consisted of 2000 Amazon product reviews. The dataset was preprocessed to remove stop words and to perform stemming. The

results of the study showed that the Random Forest performed the best with an accuracy of 87.35%, followed by Decision Tree with an accuracy of 82.5%, and KNN with an accuracy of 80.5%. The authors conclude that Random Forest is the most suitable algorithm for sentiment analysis of Amazon product reviews.

S. S. Patil, K. R. K. Singh, S. S. Deokate [6] paper proposes a machine learning-based approach for sentiment analysis of Amazon product reviews. The authors start by introducing the importance of sentiment analysis and its applications in various fields. They then describe the dataset used in the study, which consisted of 3000 Amazon product reviews. The dataset was preprocessed to remove stop words, punctuation, and perform stemming. Naive Bayes and Logistic Regression achieved accuracies of 78.4% and 80.2%, respectively. The authors conclude that SVM is the most suitable algorithm for sentiment analysis of Amazon product reviews.

B.C.Patil and V.R.Reddy [7] paper proposes a K-Nearest Neighbor (K-NN) classification algorithm for sentiment analysis of Amazon reviews. The authors start by introducing the importance of sentiment analysis and its applications in various fields. They then describe the dataset used in the study, which consisted of 1500 Amazon reviews. The dataset was preprocessed to remove stop words, punctuation, and perform stemming. The results of the study showed that the K-NN algorithm achieved the highest accuracy of 76.8% when $K=7$. The precision, recall, and F1 score for the positive, negative, and neutral sentiment categories were also reported for different values of K .

A.Bhattacharya and S.Basak [8] paper proposes a comparison of Naive Bayes algorithm and Decision Tree algorithm for sentiment analysis of Amazon product reviews. The authors begin by discussing the importance of sentiment analysis and the challenges associated with the sentiment analysis. They then describe the dataset used in the study, which consisted of 500 Amazon product reviews. The dataset was preprocessed by removing stop words, punctuation, white spaces and performing stemming. The accuracy of Naive Bayes algorithm is greater than the accuracy of the Decision Tree

algorithm. The result of the study showed that the Naive Bayes algorithm achieved an accuracy of 74.4%, while the Decision Tree algorithm achieved an accuracy of 72.8%. The precision, recall, and F1 score for the positive, negative, and neutral sentiment categories were also reported for both Naive Bayes algorithm and Decision Tree algorithm.

N.Naem, N.Naureen, and M.Shahbaz[9] paper proposed a sentiment analysis approach using the random forest algorithm to classify Amazon product reviews as positive, negative, or neutral. The authors collected a dataset of Amazon product reviews, preprocessed the data, extracted features, and trained the random forest model. The experimental results showed that the proposed approach achieved good accuracy in sentiment analysis of Amazon product reviews. The paper concluded that the random forest algorithm can be used effectively for sentiment analysis of Amazon product reviews.

III. IMPLEMENTATION

Amazon is one of the largest E-commerce sites as for that there are innumerable amount of reviews that can be seen. We used data named Amazon product data which was provided by Kaggle. Our dataset comes from Consumer Reviews of Amazon Products. This dataset has 30,847 rows and 12 attributes. Each example includes the type, name of the product as well as the text review and the rating of the product etc. For preparing the desired data a simple code was written in python to remove these useless features. We used various Natural Language Processing (NLP) techniques to prepare the data.

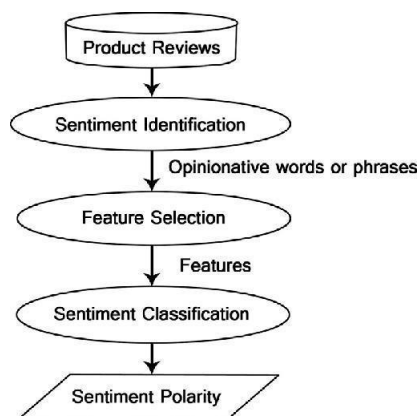


Fig: workflow

Tokenization: It is the process of separating a sequence of strings into individuals such as words, keywords, phrases, symbols and other elements known as tokens. Tokens can be

individual words, phrases or even whole sentences. In the process of tokenization, some characters like punctuation marks are discarded. The tokens work as the input for different processes like parsing and text mining.

Stop word removal: Stop words are those objects in a sentence which are not necessary in any sector of text mining. So we generally ignore these words to enhance the accuracy of the analysis. In different formats there are different stop words depending on the country language etc. **POS tagging:** The process of assigning one of the parts of speech to the given word is called Parts of Speech tagging. It is generally referred to as POS tagging. Parts of speech generally contain nouns, verbs, adverbs, adjectives, pronouns, conjunctions and their subcategories. Parts of speech tagger or POS tagger is a program that does this job. **TF-IDF:** TF-IDF is an information retrieval technique which weighs a term's frequency (TF) and also inverse document frequency (IDF). Each word or term has its own TF and IDF score. The TF and IDF product scores of a term is referred to as the TF*IDF weight of that term. Simply we can state that the higher the TF*IDF score (weight) the rarer the term and vice versa. TF of a word is the frequency of a word. IDF of a word is the measure of how significant that term is throughout the corpus. Converting all the capital letters to lowercase. **Stemming and reducing inflectional forms to a stemma form.** Lemmatizing to group together the different inflected forms of a word so they can be analyzed as a single item. The input data can be transformed into a reduced set of features (feature vectors). This process is called feature extraction. For feature labeling, we used two files containing positive and negative words collected from the dictionary. The resultant set of pre-processing is compared with these files. The positive words are labeled as '1' whereas negative words are labeled as '0'.

IV. ALGORITHMS

Logistic Regression and Support Vector Machine are the best supervised machine learning algorithms to classify the data.

4.1 Logistic Regression

Logistic Regression is one of the most popular machine learning algorithms. Which comes under the supervised Learning technique. It is used for prediction of a categorical dependent variable using a given set of independent variables. Logistic Regression predicts the output of a categorical dependent variable. Therefore,

the outcome must be a categorical or discrete value. It can be either yes or no, 0 or 1, true or false etc. but instead of giving the exact values as 0 and 1 it gives the probabilistic values which lie between 0 and 1. Logistic Regression is much similar to the Linear Regression except that how they are used. Linear Regression is used to solve Regression problems. Whereas Logistic regression is used for solving the classification problems. In Logistic Regression, instead of fitting a regression line, we fit an "S" shaped logistic function which predicts two maximum values (0 or 1). Logistic Regression is a significant machine learning algorithm because it has the ability to provide the probabilities and classify new data using continuous and discrete datasets.

4.2 Support Vector Machine

Support Vector Machine is a universal learner. Support Vector Machine has defined both input and output format. The output is either positive or negative and input is vector space. The text document is not suitable for learning. Those texts are transformed into a structured format. The texts transformed into a format which matches into the input of machine learning algorithm. The score of the texts is calculated and then the score is given as input to Support Vector Machine. Support Vector Machine (SVM) has been proved one of the most powerful learning algorithms for text categorization but text categorization sometimes may produce errors. To decide which one is better between texts a comparison of text classifier is required. SVMs can efficiently perform nonlinear classification. SVM simultaneously minimize the empirical classification error and maximize the geometric margin. So SVM is called Maximum Margin Classifier.

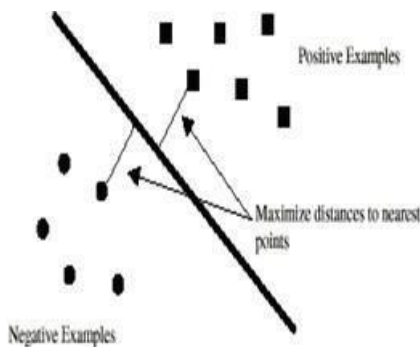


Fig: support vector classifier

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme

cases are called as support vectors, and hence the algorithm is termed as Support Vector Machine. The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n -dimensional space into classes so that we can easily put the new data point into the correct category in the future. The performance measure is used in this case.

V. CONCLUSION

Sentiment analysis deals with the classification of texts based on the sentiments they contain. This focuses on a typical sentiment analysis model consisting of three core steps, namely data preparation, review analysis and sentiment classification and describes representative techniques involved in those steps. Sentiment analysis is an emerging research area in text mining and computational linguistics, and has attracted considerable research attention in the past few years. Future research shall identify neutral reviews also and explore sophisticated methods for opinion and product feature extraction, as well as new classification models that can address the ordered labels property in rating inference. Applications that utilize results from sentiment analysis is also expected to emerge in the near future.

Algorithm	Accuracies
Logistic Regression	91%
Random Forest	87%
KNN	84%
Decision tree	85%
Support Vector Machine	91%

Fig: table of accuracy scores of various models

Thus, we get maximum accuracy in Logistic Regression and Linear SVC. So, we choose Logistic Regression and Linear SVC as our final algorithms for sentiment analysis of Amazon product reviews.

REFERENCES

- [1]. R. Singh and A. Singh, "Sentiment analysis of Amazon product reviews using decision tree algorithm," 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC), Erode, India, 2019, pp. 1064-1067, doi:10.1109/ICCMC.2019.8725549.
- [2]. S. S. Shrivastava and S. K. Shukla, "Sentiment analysis of Amazon product reviews using random forest algorithm," 2017 5th Internat

- ionalConferenceonAdvanced Computing & Communication Technologies(ACCT),Rohtak,India,2017,pp.95-98,doi:10.1109/ACCT.2017.16.
- [3]. S.S.ReddyandV.Krishnaiah,"Sentimentanalysis ofAmazonproductreviewsusingKNNalgorithm,"2016IEEEInternationalConferenceonRecentTrends inElectronics,Information&CommunicationTechnology(RTEICT),Bangalore,India,2016,pp.804807.Poi.10.1109/RTEICT.2016.7807884.
- [4]. I.MukherjeeandS.KDas.,“Sentimentanalysisof AmazonproductreviewsusingNaiveBayesalgorithm,” 2016 IEEE 7th Annual InformationTechnology, Electronics and Mobile Communication Conference (IEMCON), Vancouver, BC, 2016, pp. 1-6,doi:10.1109/IEMCON.2016.7746337.
- [5]. M. H. Rahman, M. M. Rahman, and M. H.Kabir, "A comparative study of KNN, Decision Treeand Random Forest algorithms for sentiment analysisof Amazon product reviews," 2020 IEEE InternationalConference on Computer, Communication, Chemical,MaterialsandElectronicEngineering (IC4ME2),Dhaka,Bangladesh,2020,pp.1-4,doi:10.1109/IC4ME2.2020.9334648.
- [6]. S.S.Patil,K.R.K.Singh,S.S.Deokate,S. S. Ingle, "Sentiment Analysis of Amazon ProductReviews using Machine Learning," 2021 InternationalConferenceonComputingandInformationTechnologies (ICCIT), Wollongong, Australia, 2021,pp.1-6,doi:10.1109/ICCIT52534.2021.9422872.
- [7]. B.C.PatilandV.R.Reddy,"SentimentAnalysisofAmazonReviewsusingK-NNClassification,"2021InternationalConferenceonEmergingTrends inInformationTechnologyandEngineering(icETITE),Pune,India,2021, pp.1-5,doi:10.1109/icETITE51432.2021.9451919.
- [8]. A.BhattacharyaandS.Basak,"SentimentAnalysis ofAmazonProductReviewsusingNaiveBayes andDecisionTree,"2021InternationalConferenceonInventiveResearchinComputingApplications (ICIRCA), Coimbatore, India, 2021, pp.239-243,doi:10.1109/ICIRCA52104.2021.9603899.
- [9]. N.Naeem,N.Naureen,andM.Shahbaz,"SentimentAnalysisofAmazonProductReviewsUsing RandomForest,"2020InternationalConference onAdvancesinComputing, CommunicationControl and Networking(ICACCCN),Bangalore,India,2020,pp170174,doi:10.1109/ICACCCN49741.2020.9289365.