



National
College *of*
Ireland

Sentiment Analysis of Spam Reviews Using Bert- Large with SoftMax Classifier

MSc Research Project

Data Analytics

Student Name: Atul Vasant Lambhate

Student Number: x20203624

Supervisor: Dr. Abubakr Siddig

National College of Ireland Project

Submission Sheet – 2021/2022

Student Name: Atul Vasant Lambhate

Student ID: x20203624

Programme: Msc Data Analytics **Year:** 2022

Module: Research Project

Lecturer: Dr. Abubakr Siddig

Submission Due Date: September 19th 2022

Project Title: Sentiment Analysis of Spam Reviews Using Bert-Large with SoftMax Classifier

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the references section. Students are encouraged to use the Harvard Referencing Standard supplied by the Library. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action. Students may be required to undergo a viva (oral examination) if there is suspicion about the validity of their submitted work.

Signature: Atul Vasant Lambhate

Date: September 19th 2022

Sentiment Analysis of Spam Reviews Using BERT-Large with SoftMax Classifier

Research Project

Atul Vasant Lambhate

x20203624

MSc in Data Analytics

19th September 2022

Abstract

The advent of the internet has made customers turn to e-commerce platforms for their shopping avenues. These e-commerce websites host a variety of products on their platform. To improve the user's shopping experience, these platforms created a digital word-of-mouth phenomenon in the form of reviews. Customers can help make use of these reviews to facilitate their purchases whereas they can also give vital information about the product to the platform. Identifying sentiment associated with a product helps the platform maintain its image in the market as bad product reviews and customer experiences can hamper its reputation. While manual sentiment analysis is a long and tedious process for humans to perform, machines however can be made to do it without the help of human interference. This study involved implementing machine learning models to analyze the sentiments embedded in the product reviews about the musical instruments listed on the Amazon website. Three classifiers mainly Naïve Bayes, Support Vector Machines, and Long Short Term Memory Neural Networks (LSTM) are evaluated in the study using the accuracy metric. LSTM is found to be the most accurate of them all with an accuracy of 67 percentage.

Keywords: Sentiment Analysis, Spam Reviews and Natural Language Processing, Deep Learning.

Contents

1 Introduction	3
1.1 Overview of sentiment analysis of Spam review using the BERT model	3
1.2 Motivation of the research	3
1.3 Problem Statement	4
1.4 Research Aim	4
1.5 Research Objective	4
1.6 Research Question	4
1.7 Outline of Research	4

1.8	Paper Plan	5
1.9	Summary	5
2	Literature Review	5
2.1	Introduction	5
2.2	CNN-LSTM Vs RNN-LSTM for Sentiment Analysis	6
2.3	Research Overview on BERT for Sentiment Analysis on Spam Reviews	6
2.4	Fine-Tuning BERT with Different Deep Learning Approaches	7
2.4.1	BERT with Bi-LSTM, Text-CNN and Fast-Text Model	7
2.4.2	BERT for Spam Entity Recognition	8
2.4.3	Multi-Modal BERT Approach for Sentiment Analysis	9
2.5	Factors affecting Performance of BERT	9
2.6	BERT with SoftMax classifier facilitating the sentiment analysis on spam reviews	10
2.7	Pre-training of Deep Bidirectional Transformers for Language Understanding	11
2.8	Gaps in Literature	11
2.9	Conceptual framework	12
2.10	Summary	12
3	Methodology	13
3.1	Dataset	13
3.2	Exploratory Data Analysis (EDA)	13
3.3	Natural Language Processing	14
3.4	Dimension Reduction	14
3.5	Bert Tokenization	15
3.6	Evaluation	15
4	Implementation, Evaluation, and Results	15
4.1	Implementation, Evaluation and Results of Naïve Bayes Classifier	17
4.2	Implementation, Evaluation and Results of Support Vector Machine Classifier	17
4.3	Implementation, Evaluation and Results of Long Short Term Memory Model	17
5	Conclusions and Future Work	20

1 Introduction

BERT (Bidirectional Encoder Representations from Transformer) is an eminent model that helps in establishing a network with language processing to understand context and problems. The bidirectional approach is optimised with the ideas to predict different contexts. Carrying different complex tasks including text classification, question answering, semantic similarities, summarization etc, and the BERT model performs the evaluation of different aspects. Nowadays, social media platforms are playing a significant role in gathering population sentiment in determining public opinion across the globe. Social media platforms like Facebook, Instagram, Twitter and others are there for classifying emotions that help in conducting sentiment analysis. Reflecting upon the attributes of Natural Language Processing (NLP) both positive and negative viewpoints of the global population can be evaluated. This particular research work will focus on the understanding of sentiment analysis of spam reviews with the help of the BERT model. Converting public thoughts into sentiment, research work will help in determining the utility of the BERT model with the help of the SoftMax classifier. Providing an overview of sentiment analysis using the neural-based model of BERT and classifier, the research aim, objectives and questions will be determined. In addition to that, a problem statement and motivation of the research work will be determined. Further, a brief glance at the dissertation structure will be provided outlining the overall research work.

1.1 Overview of sentiment analysis of Spam review using the BERT model

Business success is entirely dependent on the consumers in today's modern e-commerce market. The buyer's appreciation for a product and please recognise that as its success. Before shopping all the consumers review the opinions of other people who used that product before (Su et al.; 2020) The impact of these reviews according to different theories has been analysed by different researchers for several decades. According to the result of this analysis, the researchers concluded that spam reviews have several negative impacts on consumers. By creating fine-tuned BERT which utilises the SoftMax classifier as an extra layer to determine the spam reviews quickly. It is more relevant and less time-consuming than the previous methods. As a result, a new sub-specialty in sentiment analysis research has been enhanced that may be used by a wide range of businesses. Sentiment analysis will be developed based on the reviews of positive and negative feelings and information, which can be classified. This review of analysis from several models was covered by using Support Vector Machine (SVM), LSTM and Naive Bayes (Tang et al.; 2020). For translating the message text into spam and not spam feelings transformation related to a free train created a language model which is known as Bidirectional Encoder Representative (BERT) is constructed. In these proposals, the efficacy of different models and classifiers for sentiment analysis in BERT models.

1.2 Motivation of the research

The researcher has developed a significant contribution of the BERT model, which is a SoftMax classifier for sentiment analysis of spam reviews in this research. The research also covers the reliable approach of sentiment analysis with the improvement of forward and reverse feet of the BERT model. Through the research work, the readers can un-

derstand a fine-tuned BERT, which will be used by the SoftMax classifier. This research worked to determine the significance of the soft match classifier as an extra layer that can detect spam reviews effectively. The topic that is covered in this research is very reliable which helps the consumers in today's world to identify these approaches of sentiment analysis, which can be improved with forward, and reverse feeds of BERT.

1.3 Problem Statement

The enhancement of spam reviews on social media platforms reduces the trustworthiness of the users. The use of the BERT model has several disadvantages too which can impact the sentiment analysis of spam reviews. Reflecting upon the attributes of natural language processing has negative viewpoints on the global population that can be enhanced (Zhong et al.; 2019). In today's world, most consumers increasingly believe in the review of product information before purchasing those products. The usefulness of online reviews is increasingly becoming an effective tool for businesses to enhance their sales. Sometimes the customers are impeded by fake reviews, which give an untruthful product quality. Therefore, the detection of fake reviews is required.

1.4 Research Aim

The aim of the research work is to understand the significant contribution of the BERT model with SoftMax classifier for sentiment analysis of spam review. The research work also evaluates the reliable approach of sentiment analysis with improved BERT's forward and reverse feed.

1.5 Research Objective

1. To create a fine-tuned BERT that will use the SoftMax classifier
2. To determine the significance of the SoftMax classifier as an extra layer to detect spam reviews efficiently.
3. To identify the reliable approach of sentiment analysis that can be improved with BERT's forward and reverse feeds.
4. To understand the complex tasks of the BERT model including next sentence prediction, Named-entity-recognition etc.

1.6 Research Question

1. How to create a fine-tune BERT along with the SoftMax classifier?
2. How well BERT large will uncased encoder tune with SoftMax classifier to perform sentiment analysis on spam review?
3. How to analyse the complex methods of sentiment analysis using BERT analysis?

1.7 Outline of Research

With the adoption of digitalization, the market of the e-commerce industry is continuously growing with the dependency of consumers. People can get a wide number of options by the use of social networking platforms like Amazon, Flipkart, Twitter, Facebook and many more. In their busy life people are continuously purchasing various types

of production services from E-Commerce brands and offering their feedback (Su et al.; 2020). Feedback comes in both positive and negative ways in which they have expressed their feelings regarding the use of the production services. It is important to examine the spam reviews to offer better services and avoid the reviews on the website to benefit the brand image of the company. The research work will focus on outlining the various methods to improve sentiment analysis by the application of fine-tuned BERT to examine the emotions of the consumers. It will outline the ways of creating a fine-tuned BERT that will use the SoftMax classifier to perform sentiment analysis on spam reviews of the consumers. It will help in improving the purchasing decisions.

1.8 Paper Plan

The plan of the paper is focused on dividing the research work into several chapters. Based on the plan of the paper further research work will be developed (Zhong et al.; 2019). The first chapter's introduction refers to the research work that comes with defining the research aim, objectives, motivation, plan for the paper and questions. The literature review refers to the second chapter of the research that reflects on defining the existing literature on the research topic aligning with the developed research questions. Chapter 3 reflects on identifying the methods of continuing the research work called following the KDD process referred to as Methodology. Chapter four reflects on the design specification and the last chapter reflects on developing the scopes of Future Work is the fifth Chapter.

1.9 Summary

From the above discussion, it can be concluded that the introduction chapter provides a glance at the overall research aim, objectives and questions. Aligning with the research topic of sentiment analysis of spam review, the utility of the BERT large model with Softmax classifier is determined. Discussing the motivation of the research and drafting the problem statement, the chapter has provided an overview of the importance of the BERT model. For sentiment analysis of Spam review, the BERT model along with the Softmax classifier plays a significant role. The above introductory section has also provided a glance at the dissertation structure, outlining the overall research work.

2 Literature Review

2.1 Introduction

The literature review is a significant part of the whole report that offers an evaluative report on the various literature and research articles on the same topic. The literature review takes the assistance of various scholarly articles that discuss the significance of sentiment analysis in the domain of spam review. The literature review also analyses the various techniques like deep learning, combined technical approaches and machine learning capabilities. The significant application of Deep and machine learning methodology addresses real-world concerns like language processing, image classification and text classification. The literature review part discusses how the Deep and machine learning approaches to aid in detecting spam reviews while leveraging trained frameworks. The literature review part also enlightens on the significant benefits of deep learning approaches in natural language processing. The various other approaches like long/short term memory,

TextBlob, TF-IDF and other approaches help in improving the outcome with the application of BERT approaches. The previous literature would also shed light on the fact of fine-tuning the BERT layer that addresses better intervention on the SoftMax classified and the large encased encoder of BERT. With the use of various scholarly literature, the researcher targets to analyse how sentiment analysis has been widely used in order to take preventive measures for avoiding spam reviews, especially on various e-commerce websites.

2.2 CNN-LSTM Vs RNN-LSTM for Sentiment Analysis

The E-Commerce sector witnesses the emergence of spam reviews, which are degrading the effective communication and interaction between the customers and the business. The need of avoiding and addressing spam reviews on the commerce website with the use of sentiment analysis with various classifiers and the BERT approach is a significant effort toward minimizing the ill effects of spam reviews. As stated by (Asonam; 2022), a chunk of text undergoes the analysis of sentiment schools that suggest entities categories and sentiments course with the use of machine learning, natural language processing and text analytics. LSTM networks are interconnected with the RNN extensions in learning the various temple sequential data precisely designed for standardised RNNs. As per (Picone; n.d.), this has been widely applied in various deep learning practical applications like natural language processing, speech recognition and stock forecasting. The LSTM model has been widely used for sentiment classification that distinguishes between the negative and positive reviews of the customers on an e-commerce website. The integration of CNN layers in the LSTM model targets to derive and extract the significant features within the input data targets to offer a sequential prediction. This also derives better outcomes and analyse of the activity in video labelling activity recognition and image labelling. This facilitates the process of predicting a sentiment across a text while making better use of spam reviews in any E-Commerce sector. Both the model in the LSTM sentiment analysis approach would help in classifying a text in text analytics while analysing and determining the significant meaning of the same.

2.3 Research Overview on BERT for Sentiment Analysis on Spam Reviews

The application of the BERT model in sentiment analysis signifies the unsupervised spam reviews based on the algorithmic use of the fine-tuning process. The development of the LSTM model in sentiment analysis helps in classifying both the negative and positive categories of reviews that further undergo fine-tuning of the hyperparameters in order to ensure the improvement of the model specific. As stated by (Tang et al.; 2020), he fine-tuning process of the BERT Model generates more algorithmic intervention and prediction on the reviews while delivering better outcomes in order to avoid spam reviews on the e-commerce website to a huge extent. This helps in offering a better interaction and connection between the customers and the business while making better effects in helping the researchers with the uses of significant sentiment analysis for eliminating the same. The spam reviews involve the classification of traditional neural networks and various convolutional network approaches. The pre-trained models are useful in offering a perspective on the real-time classification of texts, which might be driven by the use of CNN as well as the RNN model. Neural networks that help in improving the pre-

training methods and various sets of executable algorithms drive the text classification in sentiment analysis. This sentiment analysis uses the BERT module that helps in extracting the semantic representation while the CNN module enables the extraction of various underline birds. As per (Wu et al.; 2019), the contextual interaction in CNN extracted information by the use of BiLSTM. The pre-trained language model of BERT helps in outperforming the significance of spam reviews that majorly standardise the success and customer satisfaction in the e-commerce website. The use of logistic regression as a standardised model in predicting and avoiding the family views would help in generating more predictions and forecast reports with the use of deep learning and machine learning practices (Zhong et al.; 2019), supports the need of analysing how the SoftMax classifier with the use of the fine-tuning model of BERT ensures detection of spam reviews while adapting a resilient method of sentiment analysis.



Figure 1: Social Sentiment Analysis
Source: Zhong et al., 2019

2.4 Fine-Tuning BERT with Different Deep Learning Approaches

Different deep learning approaches are determining the significance of the fine-tuning process of BERT that helps in generating more significant approaches while managing the core approach of sentiment analysis (Su et al.; 2020). The Fast-Text Model, Text-CNN, and Bi-LSTM are significant factors that help in streamlining the deep learning approaches while managing the learning algorithms in the spam reviews in various e-commerce websites.

2.4.1 BERT with Bi-LSTM, Text-CNN and Fast-Text Model

The conceptual framework of the bidirectional BERT model that has been widely used in sentiment analysis for spam reviews put better emphasis on the various deep learning approaches. The significant integration of the Fast-Text Model, Text-CNN and Bi-LSTM determine how the performance management of the sentiment analysis in text analytics

is managed and involved to a huge extent. (Selvapandian et al.; 2019), support the conceptual framework of the BERT model that helps in determining the sarcastic attitude and customer demand in order to appreciate a certain service or product across the converging platform of social media. The timely gained accurate information at the most significant approach that drives the investigation of spam reviews across the converging social media platforms. The technological intervention in various locations in order to determine the sentiment recommendations of the public union of first a rich semantic text representation. The masked language model and next sentence prediction are the two pre-training task which has been assigned by BERT Model for analysing radius spam reviews and comments for undergoing sentiment analysis (Mochihashi; 2020). This helps in analysing the attitude and opinion of the local uses on various social events and news, which is useful in guiding the sentiment analysis process across a diverse place. The text CNN model leveraging and extracting the lexical context by the Bi-LSTM offering a significant analysis of spam reviews extracts the local characteristics. The change performance and various other practices on sentiment analysis in order to address this time review are highly driven by the fine-tuned BERT model as a whole.

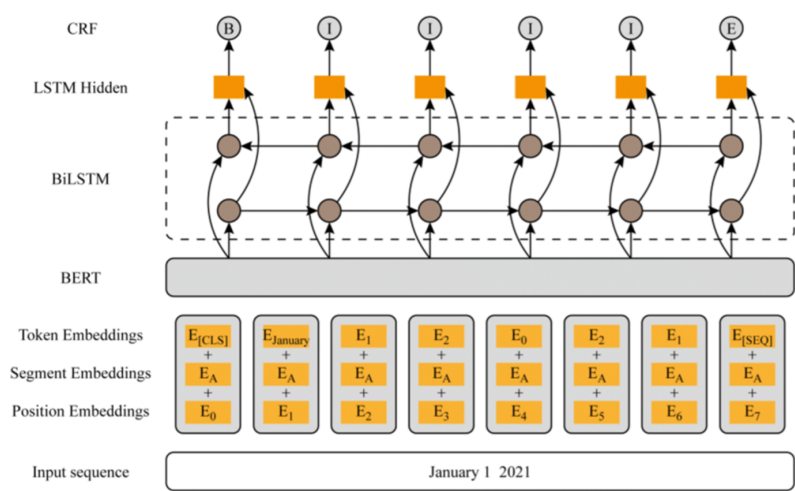


Figure 2: Overall architecture of BERT
(Source: Mochihashi, 2020)

2.4.2 BERT for Spam Entity Recognition

The entity recognition of spam in various E-Commerce websites are highly managed and exploded by the model of BERT. As per (Yaseen et al.; 2020), the deep learning practices with the intervention of pre-trained models and classifiers ensure to the transformation of the situational representation of languages. This helps in forecasting entity types while integrating the various significant classifiers and models that reflect upon a certain data set to detect the sentiment of the language to a huge extent. The various public comments that involve sarcasm, misleading information, dis-satisfactory experience and other significant approaches the significant analysis of categorising the sentiment of every comment with the use of various classifiers. The task-specific BERT Model helps in responding to various spam emails, sentiment analysis, question answering sessions and spam reviews to a huge extent. The validation of the utility of the BERT model and text extraction process. The entity recognition of the spam reviews is highly managed

by the NLP task while creating a significant BERT model for developing text-mining materials throughout a project. This is a baseline and foundation for proposing better responses to spam reviews and emails that are encountered in social media platforms to a huge extent. (Liu et al.; 2022) proposes how the superior task-specific BERT model analysis the original meaning and sentiment of the public reviews and comments on the e-commerce website of any organisation.

2.4.3 Multi-Modal BERT Approach for Sentiment Analysis

The multi-modal approach of sentiment analysis is highly driven by the BERT model. (Li et al.; 2020) opined that the fine-tuned BERT model generally manages the novel method of predicting spam reviews. This is also driven by the multimodal fusion concepts that update the internal nonverbal data and reviews of the customers on the various e-commerce websites as a whole. This approach with the multi-model framework targets to streamline the operational practices of sentiment analysis while signifying the need of addressing the spam review on various e-commerce websites (Wu et al.; 2019). This approach helps in classifying the various phrases into distinctive polarity sentiments. The fine tune BERT model generates a better inside and distinctive approach of the efforts and impact of algorithmic data sets. This also focuses on various complicated data sets that encompass picture representation, emojis, sarcasm comments and word embedding. Irrespective of the complexity of the datasets the outperformance of the BERT model helps in assuring the improvement of classification techniques of spam reviews (Paul et al.; 2019). This insurance targets the various significant approaches in managing the unprecedented change within the text content among various spam reviews from the public across the diversified social media platform.

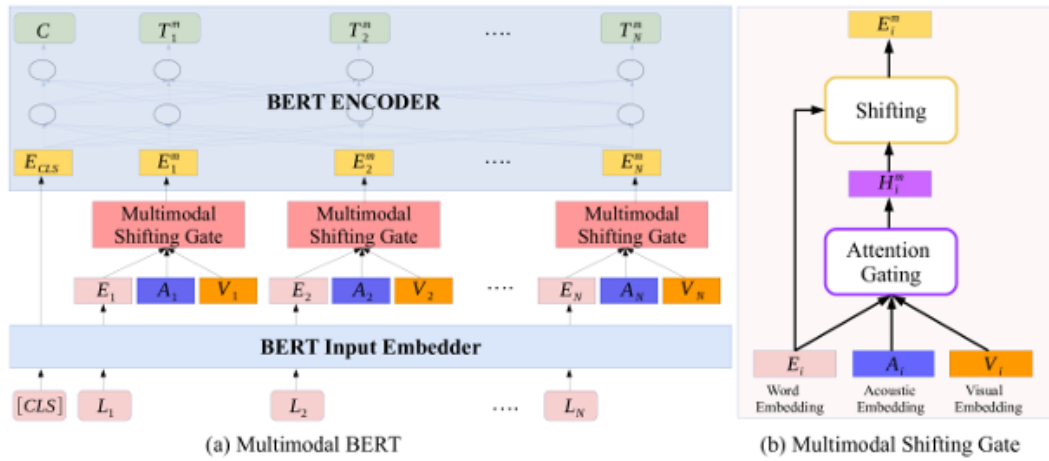


Figure 3: Multi-modal BERT
(Source: Saner, 2022)

2.5 Factors affecting Performance of BERT

The application of the BERT model in making significant changes in the way spam reviews are avoided and addressed on the e-commerce website is dependent on various factors. The architecture of this model utilises a deep neural network, which is more computationally costly than any other framework. The use of various management frameworks and

the other training contextual frameworks suggest how the BERT model integrates with ULMFit, ELMo, supervised sequence learning and generative contextual representation of the bidirectional BERT model. The various factors are highly dependent on the effective algorithmic intervention of its qualities. The various factors that affect the BERT model suggest some variances like sarcasm comments, picture representation and word embeddings. With the use of the learning methodologies, the BERT model determines the data set complexities while undergoing multiple studies on the impact of the same. As per (Picone; n.d.), the fine-tuning performance of the BERT model is not majorly impacted by the text categorization. However, it would be influenced majorly by the change in the text contents that further undergoes sentiment analysis. (eebda.org; n.d.) opined adoption of various algorithmic interventions like the optimizer of ADA-Boost, and ADAM helps in facilitating the algorithmic variety as a whole. This sheds light on the significant factor that widely affects the BERT performance while streamlining the free-trained management performances with the use of various fine tuning classifiers and place in the sentiment analysis of spam reviews. (Jeevantha et al.; 2021), opined that the critical approach of determining the various factors and their influence on the performance of BERT analysis also affects and degrades the spam review detection quality. Furthermore, (Huang et al.; 2020), contradicts the statement on how these factors are directly interconnected with the market approaches and other functional approaches. This model is developed by small data sets that ensure to categorise the various neural networks in order to employ a significant free trained model for driving better deep learning approaches for addressing the sentiment analysis practices on spam reviews.

2.6 BERT with SoftMax classifier facilitating the sentiment analysis on spam reviews

The user attitude finding inverse pitching the weights posted by the users and analysing it. Several relevant traits are fixed to extract the information of users from it. As per (Huang et al.; 2020) based on the details of users and the history of posted weight by every user is analysed to predict the attitude of users. Attitudes of users can be divided into three different categories that are users with an optimistic attitude, pessimistic attitude and neutral attitude. When implementing the pre-trained area in BERT several components can vary. It can use different types of approaches to enhance their performance and help choose a superior sentiment analysis for achieving a result. According to (Jeevantha et al.; 2021) Softmax function helps to examine the reliability of the BERT model. It is used as Loss Function and Cross Entropy Function to enhance the performance of the neural network. The software classifiers provide effectiveness to the sentiment analysis of spam reviews easily. In recent years with the rapid development of advanced technology, online shopping has become a mainstream way for users to purchase their desired products and services. As supported by (Jeevantha et al.; 2021), sentiment analysis of a wide number of users reviewed on the e-commerce platforms can effectively improve the satisfaction of users. With the help of a softmax classifier, the enhancement of the sentiment in the analysis of spam reviews can be supervised. As a result implementing the BERT model, which is an unsupervised technique, is not necessary. As opined by (Yagi et al.; 2021), the SoftMax function can be utilised as the activation function in the results of neural network models, which predict a probability distribution of multinomial. SoftMax is used as the function of activation for multi-class classification issues where class membership is needed on more than two class labels. According to (Saragih et al.; 2021) the softmax

can enable the strategy into a multi-class world. Sentiment analysis of a wide number of users reviewed.

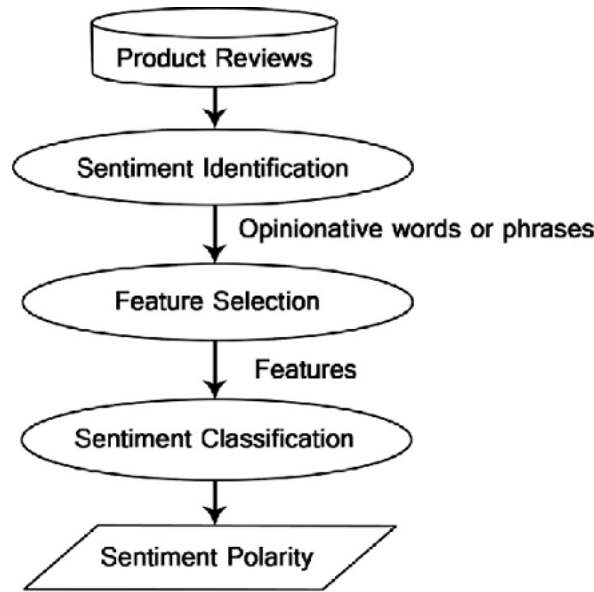


Figure 4: Sentiment analysis on product reviews
(Source: icadeis.com 2022)

2.7 Pre-training of Deep Bidirectional Transformers for Language Understanding

Several studies have been approached and used in order to interpret and predict the relevant features for detecting spam reviews, email and web spam. The various types of data categorization approaches demand the utilisation of high-level algorithms as a whole. The significant qualities of new algorithms throughout the sentiment analysis suggest the use of certain key sentiment phrases and sarcasm, which is evaluated by the BERT model. (Ye et al.; 2019), supports the approach again and introduces the need of improving the performance classification while considering the additional consonant classes, sentences, words, consonant clusters, and other ought later combinations. The significant lexical-based criteria are also analysed by this multi-task learning approach in deep learning concepts. It enables the supervised data sets to transfer into a share real form and understandable information that enhances the overall generalisation performance (Dikshit and Chandra; 2021). This reduces the threat of model overfitting for a huge number of data sets involving variations, noises and so on. The preprocessing technique that has been performed by the BERT model generally helps in task classification and feature extraction. This also involves segmentation disrupting the sentiment analysis while explaining the punctuations for better separation and segmentation of sentences. Adoption of RNN, CNN and other BERT models provides different outcomes.

2.8 Gaps in Literature

In order to carry out this research on determining the significance of the BERT Model in conducting the sentiment analysis, the researcher has gone through many scholarly

articles and journals. In recent times, the need for sentiment analysis suggests paramount importance in managing a successful customer experience in any business. However, a significant gap in the literature has been found that suggests the lack of relevance in the early research papers that are also published on the same topic. The rapid change in customer behaviour towards an e-commerce business also changes the perspective of sentiment analysis in addressing the spam reviews on any e-commerce website. (Devika et al.; 2021)

2.9 Conceptual framework

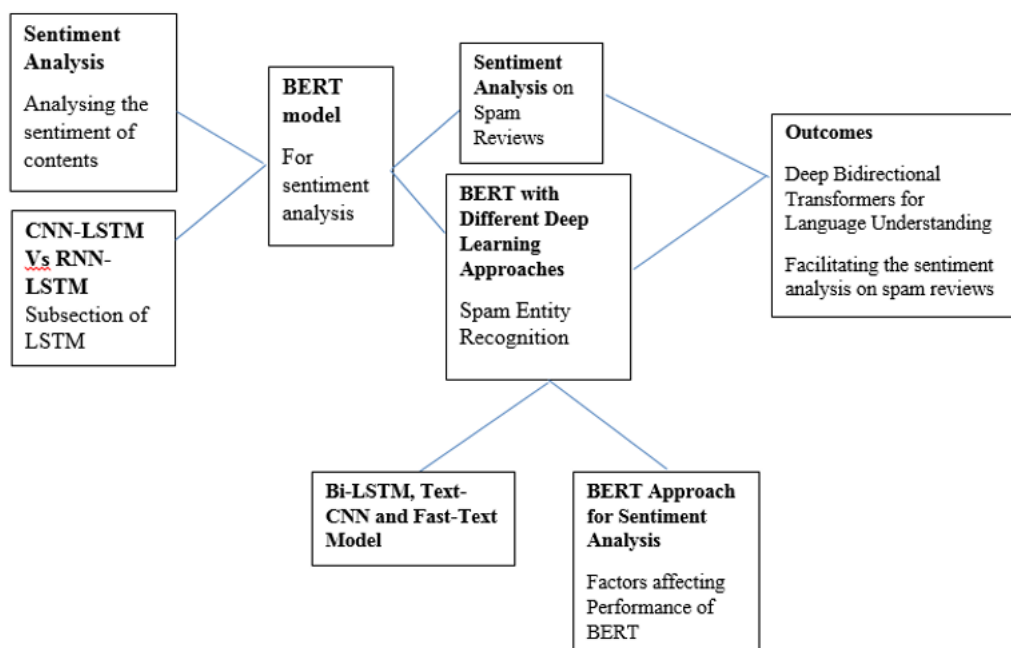


Figure 5: Conceptual framework
(Source: As created by author)

2.10 Summary

This chapter discusses the significant approach of adapting the BERT model to help in conducting a sentiment analysis and behalf of any e-commerce website. The utilization of various classifiers like SoftMax help in facilitating and managing the whole operational practices as a whole. The huge requirement of dealing with the negative sentiments of the customers in the public union drives better instances in categorising the public reviews on the diversified social media platform. The application of this model further caters to the various need of the process as a whole. This chapter sheds light on the various factors that affect the BERT performance and discusses the multi-modal approach of the BERT framework to help attain the optimized approach. Consist of the utilisation of various algorithmic models like ADAM and ADA-boost, which categorizes the contents of various texts while analysing the spam emails and reviews to a huge extent. This

also ensures drives the quality detection of BERT concepts while making better efforts in signifying the inner meaning of the spam reviews.

3 Methodology

This section discusses the methodology that has been implemented in the study. Following figure 6 depicts the steps in the methodology that has been followed.



Figure 6: Methodology Process

3.1 Dataset

The dataset under consideration contains reviews on Amazon e-commerce platform regarding musical instruments. The dataset has been acquired from Kaggle repository. The dataset contains in all 10261 entries of samples with 9 attributes. These attributes are listed in table 1 below.

Sr. No.	Attribute	Attribute Info	Values
1	<u>reviewerID</u>	ID of the reviewer	Alphanumeric
2	<u>asin</u>	ID of the product	Numeric
3	<u>reviewerName</u>	Name of the reviewer	String
4	helpful	Denotes how helpful the review is	Numeric
5	<u>reviewText</u>	Text in the review	String
6	overall	Overall rating of the product	Numeric
7	summary	Summary of the review	String
8	<u>unixReviewTime</u>	Time of the review	Numeric
9	<u>reviewTime</u>	Raw time of the review	Date

Table 1: Dataset Attributes

3.2 Exploratory Data Analysis (EDA)

Exploratory data analysis (EDA) is the process of visualizing the dataset in order to identify useful features as well as discard unwanted ones. Exploratory data analysis can be done using visual techniques and mathematical techniques. It is an important tool

for an analyst to get started with processing the data. It helps analysts get insights into the dataset at hand. With the help of EDA, an analyst can choose which model to use in the analysis. It also helps the analysts to select the proper methodology to go about during the analysis of the data. It helps identify the correlation between features in the dataset. EDA hence can be considered to be one of the most important processes in data analysis study.

3.3 Natural Language Processing

From the dataset, it is clear that the data this study aims to work on is present in natural language. The machines however do not understand the context of the text by themselves. Hence a text has to be processed such that unnecessary part of it are removed. The remaining are used for further processing. This is known as parts of speech tagging. The natural language processing is done through following steps.

Step 1: Word deconstruction

Step 2: Punctuation removal

Step 3: Stopwords removal

Step 4: Stemming

Step 5: Lemmatization

Step 6: Tokenization

3.4 Dimension Reduction

Dimensionality reduction is the method of reducing the random variables under consideration. In data analysis it plays a vital role. Using dimensionality reduction the data in higher dimension is reduced to a lower dimension making it easy for the machine to process the data as it reduces the processing requirements of the machine. This is done by finding the highly correlated features. The features that are highly correlated can be then reduced into a smaller number of features. Dimensionality reduction can be achieved using various techniques such as

1. Principal Component Analysis (PCA)
2. Independent Component Analysis (ICA)
3. Generalized Discriminant Analysis (GDA)

This study utilizes PCA as a dimensionality reduction technique. In PCA, the higher dimension data is converted to lower dimension data with high variance. This is done through following steps.

PCA may cause data loss during reduction but most of the variances are retained using eigenvectors.

- Step 1.** The covariance matrix of the data is calculated
- Step 2.** Eigenvectors of the matrix is calculated
- Step 3.** The largest eigenvalues given by the eigenvectors are used to reconstruct the large variance feature set
- Step 4.** Principal Components are used for further analysis

3.5 Bert Tokenization

Bidirectional Encoder Representations from Transformers (BERT) Tokenization is the process of tokenizing words which then can be interpreted by the machines. Tokenization in other words encodes textual data into numerical values. This help in representing the text data that is suitable for the machines to process. BERT makes use of Long Short Term Memory (LSTM) neural network model for encoding the data.

BERT is used to extract vectors representing the words that are present in the data. This is done by checking if the word is present in the vocabulary. If the word is not present in the vocabulary, the word is deconstructed into subwords. These subwords are then again checked for their presence in the vocabulary. If these words are present then the vectors representing these words are returned. These vectors as they contain numerical values, are then used for further processing.

BERT can be implemented using two techniques that are listed below.

1. Large BERT
2. Base BERT

The difference between the two models is the number of stacked layers that are responsible for encoding. In Large BERT model, there are 24 encoder layers while in Base BERT model the number of encoder layers stacked are 12. After the words in the review are tokenized, this data then can be used for modelling the classifiers.

3.6 Evaluation

The models that have been utilized in the study are evaluated based on the accuracy of the classification that has been achieved. This is done by splitting the data into two sets viz. Testing Set and Training Set. The testing set data is used for training the models whereas the testing set of data which acts as unknown data for the models is used for the classification. After the classification of the data, the predicted class of the data samples are compared with the original samples to find out the overall accuracy of the models.

Accuracy can be defined as the proportion of the predictions that are correct out of all the predictions. Mathematically, it can be written as,

4 Implementation, Evaluation, and Results

This section of the paper discusses the implementation of the models used in the study. The study is implemented in Python using Jupyter notebook. Following are the python libraries that have been used in the study.

vectorization, the vectors are normalized using the `StandardScaler()` function of “sklearn”. As these vectors tend to be of huge size, the study utilized PCA to identify 13 principal components of these vectors. This is done using `pca()` function of sklearn. The final feature vector is obtained by transforming the data with principal components.

The final feature vector is then subjected to classification using Naïve Bayes, Support Vector Machine, and Long Short Term Memory neural network. These models are trained on the training data and are evaluated on the testing data. The training data and the testing data are obtained from the feature vector by splitting it such that the training data contains 1000 samples while the testing data contains the remaining samples of the feature vector. Their implementation, evaluation, and results are as follows.

4.1 Implementation, Evaluation and Results of Naïve Bayes Classifier

Naïve Bayes (NB) classifier works by calculating the conditional probability of a given some other sample is already classified.

Implementation: The Gaussian version of the Naïve Bayes classifier is implemented in a grid manner that identifies the classifier configuration in such a way that the classifier performance is the best. It makes use of cross-validation to identify the best hyper-parameters of a classifier in order to get the best performance out of it. This is known as hyper-parameter tuning. It is done using the `GridSearchCV()` function of the sklearn library.

Evaluation and Results: Through the Grid Search, the best model is applied to the feature vector. The highest accuracy the model could achieve is 67.46 percentage for Base BERT and 66.6 percentage for Large BERT.

4.2 Implementation, Evaluation and Results of Support Vector Machine Classifier

Support vector machine works by identifying a hyperplane that distinguishes between the classes efficiently.

Implementation: The support vector classifier is implemented using the `SVC()` function of the ‘sklearn’ library. The value of hyper-parameter C of SVM is chose to be 10 and that of the gamma is set to auto.

Evaluation and Results: With the given values of the hyper-parameters, SVM could achieve the highest accuracy of 65.50 percentage and 65.27 percentage with Large BERT and Base BERT respectively.

4.3 Implementation, Evaluation and Results of Long Short Term Memory Model

LSTM model that is developed in the study consists of 3 layers. The activation function for these layers is Softmax. The loss parameter of the model is chosen to be ‘binary crossentropy’, the metrics for learning is chosen to be ‘accuracy’, and the optimizer is

chosen to be ‘Stochastic Gradient Descent with learning rate of 0.075’. LSTM model also consists of a dropout layer. This layer randomly zeroes the input at each step of learning at the given by the ‘rate’ value. This avoids overfitting in the model. The rate of dropout has been chosen as ‘0.01’.

Implementation: The implementation of the LSTM is done using the Keras library. It is trained on the training data for 10 epochs.

```

Model: "sequential_1"

```

Layer (type)	Output Shape	Param #
lstm_1 (LSTM)	(None, 256)	264192
dense_2 (Dense)	(None, 64)	16448
dropout_1 (Dropout)	(None, 64)	0
dense_3 (Dense)	(None, 1)	65

```

Total params: 280,705
Trainable params: 280,705
Non-trainable params: 0

```

Figure 8: LSTM Model Summary

Figure 8 above depicts the summary of the LSTM model. From the figure it can be seen that there are 3 layers and 1 dropout layer. The input layer has 256 neurons, hidden layer has 64 neurons and the output layer has 1 neuron. Overall 280,705 parameters are trained by the model.

Evaluation and Results: The implemented model could achieve the accuracy of 67.68 percentage for Base BERT and 66.8 percentage for Large BERT. From the analysis, it is observed that the validation accuracy of the model did not change with the epoch although the validation loss has been found to change. This may be due to the limitation associated with SoftMax activation function.

Table 2 below compares the modeling results for implemented classifiers.

Tokenization Model	Classifier	Accuracy (%)
Large BERT	Naïve Bayes	66.6
	SVM	65.5
	LSTM	66.8
Base BERT	Naïve Bayes	67.46
	SVM	65.26
	LSTM	67.68

Table 2: Modeling Results

From the table it can be seen that the LSTM model has performed better than the other two. All the models except for SVM showed improved results of accuracy for Base BERT for which the accuracy reduced from 65.5 percentage to 65.26 percentage.

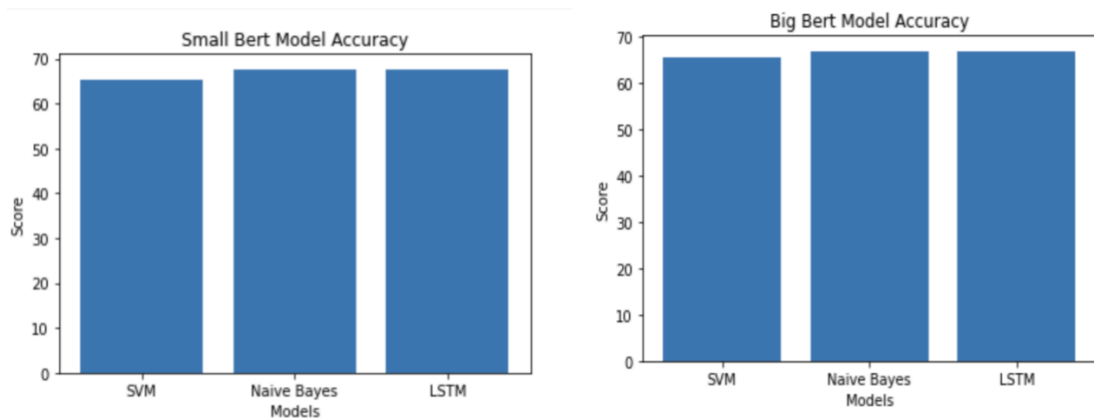


Figure 9: Comparing accuracies of the classifier for Small and Large BERT model

The Naïve Bayes model being a probabilistic model is useful in processing natural language text. The results obtained from the study provide evidence about the same as it is able to classify the sentiment better compared to the SVM model. The Naïve Bayes model also benefitted from the hyperparameter tuning performed in the study. LSTM being a sequential model is able to retain the classification information for the previous data helping it retain the sentence generation. This has helped the model to achieve the highest classification accuracy for sentiment detection. SVM model in the study shows the lowest accuracies for both of the BERT architectures. This can be attributed to the absence of hyperparameter tuning.

Table below compares the results obtained from the study with the studies in the literature review.

Author(s)	Model	Results Obtained	Comparison
Liu, Chen and Liu (2022)	Sentiment Analysis using Sentiment Dictionary, TF-IDF and SVM	Achieved accuracy of 82.9%	Better accuracy compared to the developed model owing to the Sentiment Dictionary
Rayan and Taloba (2021)	Spam Email detection using Random Forest	Improvement in accuracy	Inconclusive results of how much improvement achieved with the presented model
Su , Yu and Luo (2020)	Aspect based sentiment analysis using Capsule Network with BERT model	Accuracy of 87.41%	Model in the study achieved better classification owing to the use of more sophisticated Capsule Networks
Wu et al (2019)	Sentiment Analysis using Multiple Sentiment Dictionary, Semantic Rule Sets	Achieved higher values for precision, recall and f1-score compared to single sentiment dictionary	No machine learning modalities used sentiment is calculated using mathematical equation

5 Conclusions and Future Work

Modern e-commerce websites are trying to implement sentiment analysis in order to understand what customers are thinking about the products that are enlisted on the website. This is because negative sentiments about a product may hamper the reputation of the website itself. Going through each and every review that are posted about products on the platform is a tedious and time-consuming task for human that also requires a large size of human resources. The other way to achieve this is to make the machines do the task. However, researches are being undertaken across the globe that are dedicated towards finding the context of the text data correctly. But there are a very small number of systems that can achieve the most reliable context identification. This limits the ability of most of the models to work in this domain. This study however tried to implement the models to classify processed natural language using BERT tokenizer.

The study implemented Naïve Bayes, SVM, and LSTM models along with Big and

Base BERT tokenizer models. From the study it was observed that the accuracies of these models have been comparable and almost equal. With the accuracy of 67 percentage, LSTM model of Neural Network outperformed the other models. This can be due to the nature of the processing that LSTM utilizes. LSTM is essentially a sequential neural network model. This nature of LSTM makes it suitable for processing sequential data such as speech or natural language. This has been evident from the results obtained from the study. Although the accuracy is highest, LSTM model cannot be implemented as an automated system of sentiment analysis that can be deployed in the website. This is because, there is a 33 percentage chance of misclassification of the data into opposite class.

Naïve Bayes model has been the second best in accurately identify the sentiment embedded in the text. Naïve Bayes model is a probabilistic model. In other words, it calculates the probability that a word is present in the text. The classification is then done using the probabilities of the words in the text data. Naïve Bayes has been found to perform very well in sentiment analysis applications with high accuracy. But however, the classifier has not performed well in this study. This may be due to the overlapping data in both the negative and positive classes of the samples. To improve the accuracy of Naïve Bayes or in fact all the classifiers in the study, the context of the text must be known beforehand to remove the overlapping words such as guitar or pedal that are visible in the wordclouds in figures 7 and 8. The other probabilistic model that has been used in the study is SVM. SVM has achieved the highest accuracy of 66 percentage with the Large BERT tokenizer. SVM creates hyperplane boundary to differentiate between the classes. To achieve the most reliable separation using the hyperplane, the SVM hyperparameters need to be tuned. This process is long and time-consuming process as it involves training and testing the classifier using multiple values of these hyperparameters and obtaining their values that correspond to the highest accuracy. Because of this the tuning has not been performed in the study. This must have hampered the accuracy of the model.

To conclude it can be said that the LSTM based models are the models of the future that will be implemented in such sentiment analysis systems.

Future Work: The accuracy of the models in this study have been low. The accuracy of the LSTM model can be improved by using deep architecture of the model that involves higher number of hidden layers. Higher number of hidden layers help LSTM model to retain larger information with better accuracy. The models should also be tested across different datasets to check for their reliability. Effect of other activation functions such as relu or sigmoid can also be tested on the BERT models. Advanced neural networks can also be implemented for sentiment analysis of the reviews.

References

- Asonam (2022). asonam 2022 - the 2022 ieee/acm international conference on advances in social networks analysis and mining, *Asonam 2022 - the 2022 IEEE/ACM International Conference on advances in social networks analysis and Mining* **7**: 30220–30233.
URL: <https://asonam.cpsc.ucalgary.ca/2022/>
- Devika, P., Veena, A., Srilakshmi, E., Reddy, A. R. and Praveen, E. (2021). Detection of fake reviews using nlp sentiment analysis, *2021 6th International Conference on Communication and Electronics Systems (ICCES)*, pp. 1534–1537.

- Dikshit, P. and Chandra, B. (2021). Evaluating sentiments in social media comments on tax transformation in india using deep learning, *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pp. 1779–1782.
- eebda.org (n.d.). <http://www.eebda.org>.
URL: *eebda.org*
- Huang, Z., Epps, J., Joachim, D. and Sethu, V. (2020). Natural language processing methods for acoustic and landmark event-based features in speech-based depression detection, *IEEE Journal of Selected Topics in Signal Processing* **14**(2): 435–448.
- Jeewantha, H. C. R., Gajasinghe, A. N., Naidabadu, N. I., Rajapaksha, T. N., Kasthuri-rathna, D. and Karunasena, A. (2021). English language trainer for non-native speakers using audio signal processing, reinforcement learning, and deep learning, *2021 21st International Conference on Advances in ICT for Emerging Regions (ICTer)*, pp. 117–122.
- Li, Z., Li, R. and Jin, G. (2020). Sentiment analysis of danmaku videos based on naïve bayes and sentiment dictionary, *Ieee Access* **8**: 75073–75084.
- Liu, H., Chen, X. and Liu, X. (2022). A study of the application of weight distributing method combining sentiment dictionary and tf-idf for text sentiment analysis, *IEEE Access* **10**: 32280–32289.
- Mochihashi, D. (2020). Robotics, grounding and natural language processing, *Journal of Natural Language Processing* **27**: 963–968.
- Paul, R., Bosu, A. and Sultana, K. Z. (2019). Expressions of sentiments during code reviews: Male vs. female, *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 26–37.
- Picone, J. (n.d.). The 2022 iee signal processing in medicine and biology symposium.
URL: <http://www.ieeespmb.org/2022>
- Saragih, P. S., Witasryah, D., Hamami, F. and Machado, J. M. (2021). Sentiment analysis of social media twitter with case of large scale social restriction in jakarta using support vector machine algorithm, *2021 International Conference Advancement in Data Science, E-learning and Information Systems (ICADEIS)*, pp. 1–6.
- Selvapandian, D., Mary, R. U. N. and Karthikeyan, C. (2019). Artificial intelligence in online shopping using natural language processing (nlp), *Journal of Critical Reviews* **7**(4): 2020.
- Su, J., shan Yu, S. and Luo, D. (2020). Enhancing aspect-based sentiment analysis with capsule network, *IEEE Access* **8**: 100551–100561.
- Tang, H., Liu, H., Xiao, W. and Sebe, N. (2020). When dictionary learning meets deep learning: Deep dictionary learning and coding network for image recognition with limited data, *IEEE transactions on neural networks and learning systems* **32**(5): 2129–2141.
- Wu, J., Lu, K., Su, S. and Wang, S. (2019). Chinese micro-blog sentiment analysis based on multiple sentiment dictionaries and semantic rule sets, *IEEE Access* **7**: 183924–183939.

- Yagi, S. M., Mansour, Y., Kamalov, F. and Elnagar, A. (2021). Evaluation of arabic-based contextualized word embedding models, *2021 International Conference on Asian Language Processing (IALP)*, pp. 200–206.
- Yaseen, Y. K., Abbas, A. K. and Sana, A. M. (2020). Image spam detection using machine learning and natural language processing, *Journal of Southwest Jiaotong University* **55**(2).
- Ye, J., Peng, X., Qiao, Y., Xing, H., Li, J. and Ji, R. (2019). Visual-textual sentiment analysis in product reviews, *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 869–873.
- Zhong, G., Zhang, K., Wei, H., Zheng, Y. and Dong, J. (2019). Marginal deep architecture: Stacking feature learning modules to build deep learning models, *IEEE Access* **7**: 30220–30233.