# Research on Negative Post Identification in the Regional Language (Hindi)

MSc Research Project
Msc Data Analytics

## Deborah Ebbu Kammu
Student ID: x20217561

School of Computing
National College of Ireland

Supervisor:    Anderson Simiscuka

15th December 2022

# National College of Ireland
# Project Submission Sheet
# School of Computing

| | |
|---|---|
| **Student Name:** | Deborah Ebbu Kammu |
| **Student ID:** | x20217561 |
| **Programme:** | MSc Data Analytics |
| **Year:** | 2022 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Anderson Simiscuka |
| **Submission Due Date:** | 15/12/2022 |
| **Project Title:** | Research on Negative Post Identification in the Regional Language (Hindi) |
| **Word Count:** | 5300 |
| **Page Count:** | 18 |

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 bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| | |
|---|---|
| **Signature:** | |
| **Date:** | 15/12/2022 |

## PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Research on Negative Post Identification in the Regional Language (Hindi)

Deborah Ebbu Kammu

x20217561

15th December 2022

**Abstract**

Identifying negative posts has become a primary requirement in the current era, where social media plays a significant role in the lives of Internet users. In this context, the paper makes an attempt to determine negative posts in the regional language. Therefore, the aim of the research paper is to identify different methods that can help in determining negative posts in the regional language. In order to fulfil this aim, the research conducts a significant evaluation of other literary works that help in understanding multilingual BERT, "K-Means and Naive Bayes algorithms", deep learning approaches and "Recurrent Neural Network-based Approach" as a method of detecting negative posts. Further, particular methodologies are used for creating a design specification technique that helps in elaborating the means of implementing the method in understanding negative posts in Hindi. A post detection technique using BERT training model has been used a best approach in this project. The main result of this project has been proposed by showing the percentage of accuracy level of different model, SVM and random forest are two most widely used model here. Keywords: Negative post, social media rumours, abusive text, Hate speech, Hindi language, multilingual BERT, Deep learning approaches, convolutional neural network, "K-Means and Naive Bayes algorithms".

# 1 Introduction

## 1.1 Background

Social media is one of the most popular internet use applications. Its increasing use is due to online interaction with practical mobile technology that is portable and allows the free expression of opinions. However, this platform has significantly been used for criminal activities in the real world in the form of bullying and others. The devastating impact of cyberbullying has to be detected on social media by classifying data into positive or negative Anusha and Shashirekha (2020) or Bully or not bully. Extreme negative posts are created intentionally to target a particular race, ethnicity, gender or even country. Such negative posts on the regional language Hindi are growing significantly

1

**Byreddy et al. (2022). Hindi is the third most popular Global Language and the official Indian language. The increased availability of desktop and mobile applications and Hindi keyboards has increased the use of Hindi on social media platforms leading to negative posts (Joshi et al., 2021)**

## 1.2    Motivations

Regardless of how social media provides a powerful platform for expressing their views and opinions, it expects the users to reflect a sense of responsibility while using their freedom of speech. Particularly when social media today has a significant impact on society. The motivation for the study is due to the increasing lack of sense of duty from individuals who end up sharing deliberate material with the intention of hurting people and defaming a nation or society. Particularly identifying negative posts on the regional language, Hindi is further important since this language is actively spoken by 637 million people Bhatnagar et al. (2021). Moreover, using a regional language like Hindi helps them connect better with the people in the same region due to the flavour and context captured by the cultural prevalence. The language is further used negatively for Malpractices and wrong traditions or beliefs. Since a large chunk of people understands Hindi, negative posts are targeted in this language for spreading rumours and hate speech.

## 1.3    Aims

The aim of the paper is to identify methods for detecting negative posts on social media in the regional language Hindi that can allow early prevention of hate speech, abusive context, rumours or non-hostile aspects.

## 1.4    Objectives

The following are the comprehensive objectives of the paper in the context of identifying negative post

- To evaluate hostile post-detection techniques using a multilingual BERT pre-trained model

- To synthesise means of detecting negative posts on social media in Hindi using "K-Means and Naive Bayes algorithms"

- To evaluate the possibility of using deep learning approaches for detecting negative posts.

## 1.5    Research question

RQ 1: How does the multilingual BERT pre-trained model be helpful in detecting negative posts considering its bidirectional approach?

RQ 2: By what means can the "K-Means and Naive Bayes algorithms" successfully classify the different posts in terms of negative or positive in the Hindi language?

RQ 3: How successfully is deep learning approaches such as CNN contributing towards differentiating negative and positive textual content on social media?

## 1.6 Literature gap

The prevalent literature gap present in the research is due to the emphasis on detecting negative posts in English or other commonly spoken languages, where Hindi is often neglected, considering it as a regional language, regardless of its national significance.

## 1.7 Limitations

The persistent gap in the study in terms of academic and industrial perspectives is its contribution towards determining negative posts, which will be limited to the already identified means. The contribution of the research will be limited to theoretical understanding and without any practical implementation or implications.

## 1.8 Structure

The research paper has stated the desired problems that the research will address, along with the possible aims and objectives in the introduction section. Simultaneously the next section of related works demonstrates evidence of independent research on the classified topic by previous authors Dhagarra et al. (2020). It further provides an overview of the output gathered about different models for detecting negative posts. The following section enumerates the necessary methodology followed for creating a specified technique which is identified in the next section of the design's specification. The techniques are further implemented, and the results are thereby analysed critically to address the research questions Sharma et al. (2022). The last section discusses the solutions and addresses the problem concluding the discovery for future work.

# 2 Related Work

## 2.1 Multilingual BERT

Bhardwaj et al. (2020) proposed a pre-trained model using multilingual BERT ("Bi-directional Encoder Representations from Transformers"). The model was intended to consider annotation for hate speech and dimensions of the hostile text. The model is helpful for the computation of input embedding. Identifying such harmful post content on social media in the regional language Hindi requires conducting trials and various methods for better post identification that can classify offensive language, hate speech or neutral classes. The model requires extracting the last layer of the pre-trained model in the form of corresponding Hindi word embedding for every word within the sentence. Other sentences are presented as the average embedding of the constituents. It requires training models for both fine-grained and coarse-grained tasks employing sentence embeddings as the primary feature. The multilingual BERT-based neural model will be helpful in outperforming the existing baselines. It will help emerge as the state-of-the-art model for the problems associated with identifying negative posts in the regional language. Conducting extensive experiments such as pre-trained models, multiple pre-processing techniques, data sampling, architectural exploration, dimension reduction, and

hyperparameter tuning will be helpful De et al. (2021). This method would acknowledge the user's consensus and agreement in annotating offensive text, hate speech, bullying or any negative post.

Bhowmick et al. (2021) also proposed an efficient model with the use of the original BERT. It is referred to as the DistilBERT, which essentially processes English sentences. The model will evaluate and process mixed language sentences using the multilingual pre-trained model of XLM/Roberta. It will process Hindi-English mixed sentences supporting more than 100 4 languages that can outperform the existing BERT multilingual version. Further managing negative posts containing emoticons and emojis will be done by transforming text with such entities into corresponding text using the Python libraries Emot open source Shekhar et al. (2021). (2021) also suggested the use of the model for input post representation and further using the such pose to representation as input for machine learning models and artificial neural networks for multi-label and binary multiclass classification problems. The detected hostility development model by means of traditional machine learning algorithms such as "logistic regression, decision tree, support vector machine and random forest" on top of the embeddings of the post derived through a pre-trained BERT model.

The use of BERT as a bidirectional transformer language model that could consider the embeddings of Hindi, which required mono-lingual Corpus and post-text word embedding, allowed the use of minimal hyperparameter tuning on the validation set. As a result of this, the bidirectional model is trained on masking the next sentence prediction and token prediction task (Joshi et al., 2021). It evaluates the considered language Hindi using the original multilingual model of mBERT, which has proper training in 114 languages as launched by IndicBERT and Google, which has proper training in 10 different Indian languages launched by AI4BHARAT. Both these launched models are capable of targeting tasks as fine-tuned.

## 2.2   K-Means and Naive Bayes algorithms

Febriany and Utama (n.d.) proposed that another means of evaluating cyberbullying content and negative posts from Twitter and Facebook content is "K-Means and Naive Bayes algorithms". These methods help in classifying the sentiment into negative and positive classes. Comprehensive outcomes from evaluating such algorithms indicated that combining "K-Means and Naive Bayes algorithms" provides little less accuracy as compared to the use of Naive Bayes algorithms without applying the K-Means. The detection of cyberbullying content by means of classifying data using Naive Bayes algorithms demonstrated that the extracted data set from Facebook in the form of a new post was capable of 74 percent accuracy. It classified the data into Bully and non-bully classes Abdullah-Al-Kafi et al. (2021). The model shall be helpful in detecting negative posts that consist of different classes such as stupidity, psychology, animals, disabled persons, general pulling and attitude.

The text mining algorithm such as K means clustering along with Naive Bayes classification is helpful in deleting outlier data. K means clustering will help in grouping objects in a set in a manner that objects will be more similar to other objects in the similar group in comparison with those objects in other groups. On the other hand, the Naive Bayes algorithm acts as a simple probabilistic valuable classifier for calculating probability sets by counting the combination and frequency of values in a provided data set Sari et al. (2017). The analysis of posts on social media can also use sentiment analysis that

can determine the sentiments of the post using the "K Nearest Neighbours and Naive Bayes algorithm". This method of analysis will demonstrate classifiers extracting better outcomes for the reviews with 80 percent more accuracy and better results as compared to K-Nearest Neighbour's approach. 5 Therefore, detecting negative posts using Naive Bayes was intended to provide better accuracy; however, its combination with K means provided better results when combined with K-Nearest Neighbours Zul et al. (2018). The naive Bayes classification method established the learning of the actual level.

## 2.3 Deep learning approach

A deep learning approach for detecting hostility in the Hindi language was proposed on the basis of offensive language in tweets. It required the use of multiple long short-term memory classifiers considering the tweet content with user-specific characteristics in an ensemble setting. Other related approaches on the basis of support vector machine algorithms and long short-term memory are helpful in detecting negative posts (Joshi et al., 2021). This algorithm uses features on the basis of sentiment polarity, word embeddings, POS tags and various lexical characteristics for classifying textile content in the form of head speech. Various classifiers can be helpful in detecting hate speech, such as "gradient-boosted decision trees, support vector machines, random forest, Logistic regression" and various different "neural networks such as Long short-term memory and convolutional neural networks" Pitsilis et al. (2018).

Different Hindi text classifiers using different deep learning algorithms were based on models such as conventional neural networks and long short-term memory combining fast text word embeddings. This required use of resources in the form of models and data sets for Indian languages. In this context, Hindi acted as one of the primary mono-lingual corpus and fast text word embedding Kunchukuttan et al. (2020). The evaluation of various deep learning models for detecting negative posts demonstrates that the Hindi hostility detection data set has an input representation towards fundamental models. Comparison between pre-trained fast text embedding trend on Hindi Corpus and random initialisation of word vectors results in the identification of fine-tuned pre-trained fine text word embedding. The model using a multi-CNN algorithm with variations in word embeddings has a better advantage as compared to other basic models. The CNN-based model demonstrated better results of 11.84 percent in coarse-grained F1 scores Bhardwaj et al. (2020).

## 2.4 Recurrent Neural Network-based Approach (RNN)

The neural network solution consisting of several long short-term memory uses its powers of finding data representations that are relevant to classify data. RNN is a unique neural network type which is considered an extra loop to the architecture Pitsilis et al. (2018). It uses back propagation for training procedures that update the weight of the network in each layer.

A long short-term memory network is used in this proposed method of detecting negative posts, which is a powerful type of RNN. The model shall incorporate significant characteristics related to user information, such as the usual tendency towards sexism, racism, caste, political agenda or propaganda or others that can lead to negative posts. The approach has a critical benefit in dealing with short messages since it does not depend on pretrained vectors Badjatiya et al. (2017). Particularly users will 6 frequently

favour obfuscation of offensive terms creating new words, or using shorter slang words by inventing word concatenation or spellings.

**Limitations and Critical Analysis:**

The entire research is focusing on deep learning approach and other multiple algorithms for understanding the approach of negative post identification on regional language process. Here, in this literature five different aspects and algorithms of deep learning approach has been defined for understanding the research report approach completely. The information regarding the sentimental analysis should be provided more to explain the process completely. In this literature review section the information of sentimental analysis is quite missing. For identification of negative post the sentiment of the posts are one of the most important key factor to keep the track quite impressive.

# 3 Research Methodology

## 3.1 Research Type

For this research quantitative research methodology has been used, quantitative research methodology is a process of investigating and collecting quantifiable data and performing multiple statistical and computational techniques effectively to get the accurate result of particular research topic. Here, the research topic is based on the researching negative post identification in Hindi regional language where multiple quantifiable data and other stats are already available in multiple resources Joshi et al. (2021). To collecting and analysing all such possible resources quantitative research methodology has been used here. Social media has become an important platform for anyone to share their views on any particular topic or events. The active number of users of social media are increasing day-byday as well, where some of users are just here for creating any type of controversy. If anything happens around anyone, he or she post something related to that particular event whether he or she don't have any accurate idea about the event but still they post something the wants to. This type of activity creates a big reason for increasing negative impact on any particular post. The approach of text analysis and text identification is very important here for identifying any negative post in any language. Here in Hindi regional language both text analysis and text identification approach has been used to identify the negative post effectively. On social media sites people posts their opinion, suggestions and thinking in a free form and due to this a large number of text data be visible for interpretation Gupta et al. (2021). The method of text sentimental analysis (SA) is something which is very important in the negative post identification in any related language Jain et al. (2020).

Talking about the procedures, text sentimental analysis as a procedure is used for this research topic because the entire practical scenario of this research is based on text identification and text analysis where sentiments or opinions matters a lot for identifying negative or positive post identification effectively. Sentiment analysis is a natural language processing (NLP) technique used to determine the accuracy of positive and negative data effectively. This type of analysis is usually performed on textual data to monitor the sentiment in any particular post feedbacks or something else. Sentiment analysis can be divided into multiple types such as graded sentiment analysis, emotion detection, multilingual sentiment analysis and many more Jahan et al. (2021). Here in this pro-

ject research the graded sentiment analysis is the main area of procedure that has been followed in the entire process effectively. Sentiment analysis is very important because it helps to analyse multiple feedbacks related to any posts like the opinions, negative or positive review, people reaction and many more. All such concepts are very much important and required to identify any positive and negative post identification inside a particular language scenario.

The suggested system accepts Regional text documents as input that are kept in a certain place. Following the selection of the input file, pre-processing is carried out, during which sentences are broken up into words, native Unicode values are stored in a matrices, and special letters and symbols, if any, are eliminated. The second step is machine translation, which involves looking for terms in Hindi that are in English. A

Figure 1: Architecture Methodology Diagram

pure Hindi word is used in place of any English words that are discovered. Sentiment analysis is very important because it helps to analyses multiple feedbacks related to any posts like the opinions, negative or positive review, people reaction and many more. All such concepts are very much important and required to identify any positive and negative post identification inside a particular language scenario.

## 3.2 Data Collection

Data collection is a process of collecting information from all the relevant sources to find a particular solution of the research problem. For making a better decision for the research or conducting a better result process the data collection approach is one of the most important aspects Rajalakshmi et al. (2022). As a data collection secondary data collection method has applied for this research to collect the data. Secondary data collection means the information of data that is already available and this type of data are previously collected and has undergone necessary statistical analysis as well. All such data in secondary data collection are to be collected from primary resources and later made available to everyone else to access the information of data.

In secondary data collection quantitative secondary data collection method has been used here. Quantitative data deal with numbers, statistics and many other technical information are connected with this type of data collection methods. Some of information related to the technical implementation like data sets and something else has been collected from using online resources for completing the entire practical implementation of this research project. Some of public libraries was accessible after using the secondary data collection method for the research project effectively, the data available in both government and non-government agencies has been successfully collected here by using this type of data collection method properly.

## 3.3   Data Analysis:

For the secondary data analysis three different steps has been followed such as development of research questions, identification of dataset and dataset evaluation. Here for analysing the data both data identification and dataset evaluation technique will be used. As this is clear that the data has been collected using secondary research and a complete code will be conduct to analyse this dataset effectively. A machine learning based approach has been used here by using python programming language. The main dataset is emotions.csv here in this project which is containing multiple information regarding negative post identification on the regional Hindi Language.

## 3.4   Techniques

For this negative post identification research project multiple pre-trained model has been used here such as BERT (bidirectional encoder representation from transformers) and the purpose of this model is to consider annotation for identifying negative speech or post in a particular data of text. Here the approach of text analysis and text identification is one of the most important aspects that has been followed for completing all such possible aspect in different manner (Granik and Mesyura, 2017). The entire identification process are to be done using social medial where more than one datasets are can be used for identify negative post effectively. The text mining algorithms has also been used in this research project to effectively identify negative posts from a set of texts on social media platforms Rudra et al. (2019). Text mining is also known as text analysis is a process of transforming unstructured text into structured data for easy text analysis process. Natural language processing (NLP) is to be used to perform the operation of text mining successfully. There are list of text mining algorithms such as:

- Support vector machine (SVM)

- Random Forest

- Neural Networks

- K-Nearest Neighbour (KNN)

In the practical implementation two main text mining algorithms has been used such as support vector machine and random forest for identifying and analysing negative post identification successfully. Support vector machine is a supervised machine learning 8 algorithm used for both classification and regression Kar et al. (2021). This type of algorithm is to be determine the best decision boundary between vectors to a particular group. A significant classification type is to be chosen for applying the text analysis process in SVM algorithm. Multiple features of the same category are to be used by SVM to predict its classification perfectly. In the practical solution file the resultant of SVM implementation has been successfully described. Random forest is one of the best classifiers which is most widely used for regression and classification data. In text mining approach the implementation of random forest (RF) is highly recommended and used effectively by applying different aspect of machine learning algorithms effectively. For providing the decisions accurately by using the decision tree algorithm the decision tree approach is to be used in random forest (RF) algorithmic approach in text mining process. The complete practical resultant has been effectively provided in the implementation file of this research project Keshri and Sahu (2022).

**Gap:**

The data preprocessing is one of the most important step here in this project, the purpose of data preprocessing is to prepare a raw data and to make is suitable for a machine learning model, for identifying a particular state the implementation of data preprocessing is very much required, that directly integrated with the machine learning model. Multiple steps are there of data preprocessing such as getting the dataset, library import, importing datasets, encoding data, feature scaling and many more. In this methodology section some gaps are there of data preprocessing but still the basic requirement of the practical part has been achieved successfully. Feature extraction is the process of transforming raw data into numerical feature that further can be processed while preserving the information in original data set. Here for this project a set of data has been integrated for describing the negative post identification process and that dataset is containing a list of data of Hindi texts in multiple order for achieving the process of data extraction by using feature extraction successfully.

# 4 Design Specification

A python programming based practical technical approach has been used here for this research project. Multiple techniques have been used for completing this research project for identifying negative post successfully such as deep learning approach, k-means and nave algorithm, recurrent neural network-based approach and others text mining approaches. The purpose of implementing all such possible techniques is to performing the text identification process and approaches effectively by using several technical description and implementation process. The machine learning and neural network deep learning network are most common technical fundamentals to be used in such kind of practical scenarios Mundra and Mittal (2022).



Figure 2: Architecture Diagram

This diagram is about the architecture diagram here, a basic architecture flow of this project is created like the data will be undergo in a training process by integrating a technical approach of machine learning programming language then a model approach based on the machine learning language will be integrated. Here, SVM and random forest two different models are used to complete the practical aspect effectively.

## 4.1 Text Sentiment Analysis in NLP:

Sentiments are some sort of expressions that are to be expressed by humans in every moment whether they are happy, sad or anything. Nowadays social medial platform has become a very common platform for showing the sentiments or emotions effectively to all over the world by posting some captions or posting some sort of images as well. Both ways are enough to show the sentiments effectively here. Sentiment analysis in NLP is about identifying or analysing such sentiments from texts that are spreading negative approaches or negative sentiments. For this object analysis and objective identification is a term that is to be used for performing all such possible operations for analysing such negative comments effectively. A set of datasets are required for performing such kind of operations here and that dataset should contain all possible related information in both positive and negative manner Roy et al. (2022). The data can be collected from multiple resources like social medial datasets of Twitter and many more. After collecting the dataset the further operation of loading data, implementing expletory data analysis, text cleaning, data set creation and evaluation metrics are to be performed for collecting or providing all possible result related to the dataset loaded.

Machine learning work with different types of textual information like social medial posts, emails, messages, spams and many more Mundra and Mittal (2022). The approach of

11

machine learning algorithm is very much required for discussing the operation related to the text mining process effectively. Here, the text mining approach of this project has been implemented both machine learning and natural language processing technique effectively by analysing multiple sources of data properly. The process of text analysis, text mining and text analytics are the terms that is very much important and required to understand here, text analysis and text mining are same but text analysis and text analytics are much different from each other. Text analysis works with the concept of meaning the text where this is be used to answer the questions asked in this dataset Mandl et al. (2021). Multiple machine learning algorithms are there that is to be used to perform all such possible aspects of practical scenarios.



Figure 3: Process Model Diagram

This given diagram is about the process model and here in this diagram all possible process regarding the negative post identification of this project has been identified properly. In this process model diagram the used approaches and techniques or process are defined.

# 5 Implementation

In the implementation part of this project here python programming language has been used for performing machine learning operation and natural language process operations as well. Multiple libraries of python programming language such as Pandas, NUMPY, and many more has been imported for performing all such possible operation of the code effectively Nayel and Shashirekha (2019). For writing and executing arbitrary python code using browser Google Colab has been used in this coding implementation part for performing the operation effectively. Several packages of this code execution has been imported for this practical for identifying negative post effectively.

Figure 4: Implementation

As this image provided a pip install run code for downloading the google coding runner and execution platform libraries for performing the browser supported coding and running execution style for the particular dataset. The reason of choosing this because here in this project social media posts will be analysed and identified for analysing the negative post and all such posts are to be posted on online platform so browser supporting is always required for this. This is the reason that particular library has been imported here effectively. Multiple 10 files has been imported here for understanding the emotions and sentiments of the posts Patwa et al. (2021). Every sentiment has a label to describe the type of sentiment like anger, happy and others. All possible aspects of machine learning and machine learning process like SVM, random forest has been implemented in this implementation process for understanding and analysing all different aspects of negative post identification process effectively Yadav et al. (2022).



Figure 5: Dataset

The problem occurs during this implementation was to clean the dataset effectively because in the dataset huge amount of data were available and this was very important and required to perform the data cleaning process to delete data duplicity effectively and

13

for these multiple other aspects were performed effectively. The process of data cleaning is very complex and important at a same point of time and for every dataset that is very much important and required to perform prediction and identification process this is important to clean such dataset properly Pareek et al. (2022). All the coding process were managed by a set of instructions followed by machine learning and neural learning network process to get the level of data accuracy effectively, now here some sort of process were implemented in the coding part count the length of the dataset, errors in dataset, duplicate entry in dataset and many more. All such possible aspects were performed effectively here in such a manner to collect the data perfectly to get the accurate answer. Before starting the coding part some sort of researcher was done to check the all possible aspects of this practical implementation part effectively, here a complete review of the project was done perfectly for managing the data and all other information effectively to understand the things properly Nanda et al. (2018).

# 6   Evaluation

This negative post identification project has been evaluated on an intelligent machine learning based system that is to be developed and implemented by python programming language along with neural language processing as well. Different types of machine learning algorithms and text mining algorithms has been performed here to identify and analyse the negative post identification effectively Padmaja et al. (2020) A set of collection counter was applied in this implementation part of the project that has been evaluated successfully here in this project. Apart from this multiple python libraries has also been attached and provided in this research project for effectively identify the negative post on social media platform effectively. A graph has been plotted here for showing the accuracy of negative comments identification on the posts.



Figure 6: Evaluation

Here, top 50 and bottom learned features of the negative posts has been identified in this coding scenario effectively by implementing multiple machine learning features effectively. The validation score has also been provided here for this dataset for understanding a cross validation score value effectively.

```
plt.plot(history.history['acc'])
plt.plot(history.history['val_acc'])
plt.title('Model accuracy')
plt.ylabel('Accuracy')
plt.xlabel('Epoch')
plt.legend(['Train', 'Test'], loc='upper left')
plt.show()
```
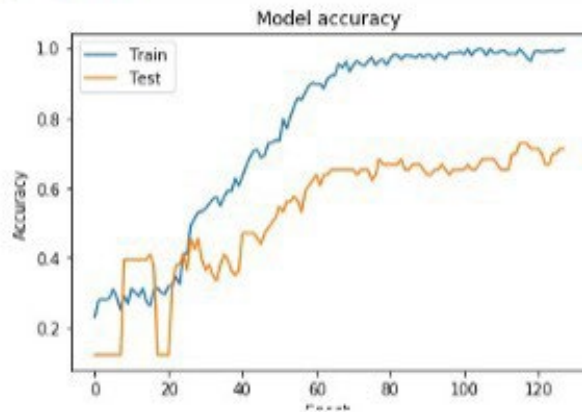
Figure 7: Graph Plotting

Here, in this graph both train and test data has been tested effectively for identifying the particular negative post identification effectively by considering multiple technical implementation and accuracy effectively. The result of logistics regression algorithm has also been used here there is providing the number of accuracy for both training and test data effectively in this practical implementation process Pitsilis et al. (2018).

## 6.1   Experiment / Case Study 1

The experimental screenshots are provided here:

Figure 8: SVM model

```
from sklearn.linear_model import LogisticRegression
model1 = LogisticRegression(C=0.05, max_iter=10000, solver='newton-cg', multi_class='multinomial')
model1.fit(X_train_vectorized, y_train)
```

```
                    LogisticRegression
LogisticRegression(C=0.05, max_iter=10000, multi_class='multinomial',
                   solver='newton-cg')
```

```
from sklearn.metrics import accuracy_score
X_test_transformed = vect.transform(X_test)
y_pred_train = model1.predict(X_train_vectorized)
y_pred_test = model1.predict(X_test_transformed)
print('Train accuracy = ', accuracy_score(y_train, y_pred_train))
print('Test accuracy = ', accuracy_score(y_test, y_pred_test))

Train accuracy =  0.7440347071583514
Test accuracy =  0.6346153846153846
```

Figure 9: Logistic Regression Model

## Random Forest

```
from sklearn.ensemble import RandomForestClassifier
classifier= RandomForestClassifier(n_estimators= 10, criterion="entropy")
classifier.fit(X_vectorized, y)
y_pred= classifier.predict(X_vectorized)
print(classification_report(y_pred, y))

              precision    recall  f1-score   support

       angry       1.00      0.93      0.96       140
       happy       0.97      1.00      0.98       146
     neutral       0.96      0.98      0.97       125
         sad       0.97      0.99      0.98       102

    accuracy                           0.97       513
   macro avg       0.97      0.98      0.97       513
weighted avg       0.98      0.97      0.97       513
```

Figure 10: Random Forest Model

## 6.2  Discussion

Here, in this implementation part all possible algorithms and text mining approaches has been performed effectively and the resultant of the research has been provided in the coding screenshot area also. The accuracy percentage of every used algorithm is quite different like for LSTM-128 around 65 percent accuracy with 4 mins is describing, for logistic regression around 74 percent accuracy with 20 sec is describing, for SVM with grid search around 95 percent accuracy with 1.5 mins is showing. The support vector machine is used for classification and for regression problems as well, the basic purpose and goal of support vector machine is to create the best line for a decision boundary and this boundary can provide a dimension space into classes with easy data point accessibility process. Random forest is one of the most popular machine learning algorithm which is used for supervised learning techniques and is also used for both classification and regression model as well. It is basically a classifier that contains a number of decision tree on various subsets to improve the productivity process.

| Algortihms Used | N-Grams |
|---|---|
| LSTM - 128 EPOCHS | 65% accuracy with 4 mins |
| LSTM -300 EPOCHS  Adadelta. | 72% accuracy with 7 mins |
| Logistic Regression | 74% accuracy with 20 sec |
| Logistic Rgeression with Cross Validation | 79% with 30 sec |
| SVM | 86 % accuracy with 55 secs |
| SVM with Grid Search | 95 % accuracy with 1 .5 mins |
| Random Forest | 97 % accuracy with 80 secs |

Figure 11: Algorithms and Results

# 7    Conclusion and Future Work

Here, in this research project all possible approaches related to machine learning algorithms and natural language processing has been performed effectively. The main motive of this research is to get the knowledge of both practical and theory for identifying the negative post on Hindi regional language by using a technical implementation part. The identification of negative post using social media platform is something very important that is highly required to be implement here for removing the risk of negative sentiment.

The approach of sentiment analysis has been used here in this project for identifying the negative post sentiment on social media. The sentiment analysis is very important to understand the sentiment any post posted on social media platform and for understanding the sentiment of every available posted post this is very much important to implement a set of technical practical code to rectify the entire process effectively. Here text mining approach, machine learning and natural language process related multiple algorithms has been implemented successfully for describing the negative post identification effectively. The coding part of this research is showing the accurate result based on code analysis. In further work some other advanced research technical approaches will be used that will provide a notification about negative sentiment on the time of posting and this will be a massive further change.

# References

Abdullah-Al-Kafi, M., Tasnova, I. J., Wadud Islam, M. and Banshal, S. K. (2021). Performances of different approaches for fake news classification: An analytical study, *International Conference on Advanced Network Technologies and Intelligent Computing*, Springer, pp. 700–714.

Anusha, M. and Shashirekha, H. (2020). An ensemble model for hate speech and offensive content identification in indo-european languages., *FIRE (Working Notes)*, pp. 253–259.

Badjatiya, P., Gupta, S., Gupta, M. and Varma, V. (2017). Deep learning for hate speech

detection in tweets, *Proceedings of the 26th international conference on World Wide Web companion*, pp. 759–760.

Bhardwaj, M., Akhtar, M. S., Ekbal, A., Das, A. and Chakraborty, T. (2020). Hostility detection dataset in hindi, *arXiv preprint arXiv:2011.03588* .

Bhatnagar, V., Kumar, P., Moghili, S. and Bhattacharyya, P. (2021). Divide and conquer: an ensemble approach for hostile post detection in hindi, *International Workshop on Combating On line Ho st ile Posts in Regional Languages dur ing Emerge ncy Si tuation*, Springer, pp. 244–255.

Bhowmick, R. S., Ganguli, I., Paul, J. and Sil, J. (2021). A multimodal deep framework for derogatory social media post identification of a recognized person, *Transactions on Asian and Low-Resource Language Information Processing* **21**(1): 1–19.

Byreddy, R. R., Malladi, S., Srikanth, B. and Battula, V. (2022). Analysis of different methodologies for sentiment in hindi language, *Smart Intelligent Computing and Applications, Volume 1*, Springer, pp. 561–567.

De, A., Elangovan, V., Maurya, K. K. and Desarkar, M. S. (2021). Coarse and fine-grained hostility detection in hindi posts using fine tuned multilingual embeddings, *International Workshop on Combating On line Ho st ile Posts in Regional Languages dur ing Emerge ncy Si tuation*, Springer, pp. 201–212.

Dhagarra, D., Goswami, M. and Kumar, G. (2020). Impact of trust and privacy concerns on technology acceptance in healthcare: an indian perspective, *International journal of medical informatics* **141**: 104164.

Febriany, A. and Utama, D. N. (n.d.). Analysis model for identifying negative posts based on social media.

Gupta, V., Jain, N., Shubham, S., Madan, A., Chaudhary, A. and Xin, Q. (2021). Toward integrated cnn-based sentiment analysis of tweets for scarce-resource language—hindi, *Transactions on Asian and Low-Resource Language Information Processing* **20**(5): 1–23.

Jahan, M. S., Oussalah, M., Mim, J. and Islam, M. (2021). Offensive language identification using hindi-english code-mixed tweets, and code-mixed data augmentation, *Forum for Information Retrieval Evaluation (Working Notes)(FIRE), CEUR-WS. org*.

Jain, D., Kumar, A. and Garg, G. (2020). Sarcasm detection in mash-up language using soft-attention based bi-directional lstm and feature-rich cnn, *Applied Soft Computing* **91**: 106198.

Joshi, R., Karnavat, R., Jirapure, K. and Joshi, R. (2021). Evaluation of deep learning models for hostility detection in hindi text, *2021 6th International Conference for Convergence in Technology (I2CT)*, IEEE, pp. 1–5.

Kar, D., Bhardwaj, M., Samanta, S. and Azad, A. P. (2021). No rumours please! a multi-indic-lingual approach for covid fake-tweet detection, *2021 Grace Hopper Celebration India (GHCI)*, IEEE, pp. 1–5.

Keshri, S. P. and Sahu, N. (2022). Sentiment analysis of public expressions in hindi language on media platform, *Global Journal of Computer Science and Technology* .

Kunchukuttan, A., Kakwani, D., Golla, S., Bhattacharyya, A., Khapra, M. M., Kumar, P. et al. (2020). Ai4bharat-indicnlp corpus: Monolingual corpora and word embeddings for indic languages, *arXiv preprint arXiv:2005.00085* .

Mandl, T., Modha, S., Shahi, G. K., Madhu, H., Satapara, S., Majumder, P., Schäfer, J., Ranasinghe, T., Zampieri, M., Nandini, D. et al. (2021). Overview of the hasoc subtrack at fire 2021: Hate speech and offensive content identification in english and indo-aryan languages, *arXiv preprint arXiv:2112.09301* .

Mundra, S. and Mittal, N. (2022). Cmhe-an: Code mixed hybrid embedding based attention network for aggression identification in hindi english code-mixed text, *Multimedia Tools and Applications* pp. 1–28.

Nanda, C., Dua, M. and Nanda, G. (2018). Sentiment analysis of movie reviews in hindi language using machine learning, *2018 International Conference on Communication and Signal Processing (ICCSP)*, IEEE, pp. 1069–1072.

Nayel, H. A. and Shashirekha, H. (2019). Deep at hasoc2019: A machine learning framework for hate speech and offensive language detection., *FIRE (Working Notes)*, pp. 336–343.

Padmaja, S., Bandu, S. and Fatima, S. S. (2020). Text processing of telugu–english code mixed languages, *Advances in Decision Sciences, Image Processing, Security and Computer Vision*, Springer, pp. 147–155.

Pareek, K., Choudhary, A., Tripathi, A. and Mishra, K. (2022). Comparative analysis of social media hate detection over code mixed hindi-english language, *Advances in Data and Information Sciences*, Springer, pp. 551–561.

Patwa, P., Bhardwaj, M., Guptha, V., Kumari, G., Sharma, S., Pykl, S., Das, A., Ekbal, A., Akhtar, M. S. and Chakraborty, T. (2021). Overview of constraint 2021 shared tasks: Detecting english covid-19 fake news and hindi hostile posts, *International Workshop on Combating On line Ho st ile Posts in Regional Languages dur ing Emerge ncy Si tuation*, Springer, pp. 42–53.

Pitsilis, G. K., Ramampiaro, H. and Langseth, H. (2018). Detecting offensive language in tweets using deep learning, *arXiv preprint arXiv:1801.04433* .

Rajalakshmi, R., Reddy, P., Khare, S. and Ganganwar, V. (2022). Sentimental analysis of code-mixed hindi language, *Congress on Intelligent Systems*, Springer, pp. 739–751.

Roy, P. K., Bhawal, S. and Subalalitha, C. N. (2022). Hate speech and offensive language detection in dravidian languages using deep ensemble framework, *Computer Speech & Language* **75**: 101386.

Rudra, K., Sharma, A., Bali, K., Choudhury, M. and Ganguly, N. (2019). Identifying and analyzing different aspects of english-hindi code-switching in twitter, *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* **18**(3): 1–28.

Sari, H. L., Suranti, D. and Zulita, L. N. (2017). Implementation of k-means clustering method for electronic learning model, *Journal of Physics: Conference Series*, Vol. 930, IOP Publishing, p. 012021.

Sharma, A., Kabra, A. and Jain, M. (2022). Ceasing hate with moh: Hate speech detection in hindi–english code-switched language, *Information Processing & Management* **59**(1): 102760.

Shekhar, C., Bagla, B., Maurya, K. K. and Desarkar, M. S. (2021). Walk in wild: An ensemble approach for hostility detection in hindi posts, *arXiv preprint arXiv:2101.06004* .

Yadav, V., Verma, P. and Katiyar, V. (2022). Long short term memory (lstm) model for sentiment analysis in social data for e-commerce products reviews in hindi languages, *International Journal of Information Technology* pp. 1–14.

Zul, M. I., Yulia, F. and Nurmalasari, D. (2018). Social media sentiment analysis using k-means and na¨ıve bayes algorithm, *2018 2nd International Conference on Electrical Engineering and Informatics (ICon EEI)*, IEEE, pp. 24–29.