

## PREDICTING HATE SPEECH IN SOCIAL MEDIA USING DEEP LEARNING

**Dr. Syed Shabbeer Ahmad<sup>1\*</sup>, Dr. Krishna Prasad K<sup>2</sup>**

<sup>1</sup>Post-Doctoral Research Scholar (23SUPDR16), Institute of Engineering and Technology, Srinivas University, Mangalore, Karnataka, India Orchid Id: 0000-0003-2931-8218,

Email Id: shabbeer.ahmad@mjcollege.ac.in

<sup>2</sup>Professor, HOD of Cyber Security and Cyber Forensics, Institute of Engineering and Technology, Srinivas University, Mangalore, India. Orchid Id: 0000-0003-1081-5820,

Email Id: krishnaprasadkcci@srinivasuniversity.edu.in

---

### ABSTRACT

Identifying hate speech detection on social media has become increasingly crucial for society. As the hate speech words are misused among the group of people. This rise of hate speech has led to conflicts and cases of cyber bullying. To overcome this problem deep learning and machine learning techniques are used such as LSTM (long short term memory), RF (random forest), SVM (support vector machine) and DT (decision tree). Among these techniques we have proposed LSTM and RF because of his best accuracy confidence level. The datasets or tweets are collected from twitter to make a separate dataset to be trained by ml techniques. These modules are trained datasets for machine learning techniques. These data files after analysis are stored as .csv files. These files go for data preprocessing and specifies every statement with IDs with the help of tokenizer. These modified csv files are applied for LSTM architecture layers 1) Embedding 2) LSTM 3) Dense. After applying the techniques the modified LSTM is saved as .hy file and the same process is done for RF and the modified file is stored as .pkl file. The app.py is a flask to construct a router interconnection between frontend and backend. To accept the file it use request method and to send it use render template method this app.py is a constructor to run the website to classify the speech is positive or negative. The precision rate of accuracy of proposed LSTM is 95.74% and for RF is 80.88% when compared with other techniques (SVM &DT).

**Key Words: Hate Speech, Social Media, Deep Learning**

### 1. INTRODUCTION

Social media networks (SMNs) are the fastest means of communication as messages are sent and received almost instantaneously. SMNs are the primary media for perpetrating hate speech nowadays. Cyber - hate crime has also grown from past few decades. More researches are being conducted to curb with the rising cases of hate speeches in social media (SM). Different methods have been made to SM providers to filter each comment before allowing it into the public domain. The impacts of hate crimes are already overwhelming due to widespread adoption of SM and the anonymity enjoyed by the online users. To detect such kind of hate speech in big era of data by manually is a time consuming and difficult process

to classify the text. Besides, the precision of the categorization of manual text can easily be influenced by human factors, such as exhaustion and competence. To achieve more accurate results, it is beneficial to use machine learning (ML) approaches to automate the text classification processes. There have been significant advancements in ML techniques from classical ML, ensemble and deep learning techniques for hate speech detection. Due to the unprecedented advancement in natural language processing (NLP), several machine learning methods have achieved superior outcomes.

To improve classification of SM texts as hate speech or non-hate speech, researchers and practitioners require an updated understanding of machine learning methodologies, which is fast evolving. Considerable effort has been spent on creating new and effective features that better capture hate speech on SM. Slangs and new vocabularies are also constantly evolving in the SM space. New and updated datasets are also available across different regions of the world. To bridge the gap, there is a need to review the literature and keep professionals, old and new researchers in the know of the current developments in this research area. On this note, this review becomes necessary to be conducted.

It is undeniable that social media has improved our lives in many ways, like allowing interactions with others all over the world and network expansion for businesses. However, there are detrimental effects of such accessibility, including the rapid spread of hate through offensive messages typically directed toward gender, religion, race, and disability, which can cause psychological harm. To address this problem of social media, many researchers have recently proposed various algorithms powered by machine learning (ML) and deep learning for the detection of hate speech. This work proposes a hate speech detection model based on long-short term memory (LSTM). Posters of hate speeches usually attack their targets using the following attributes: Religion, Race, political affiliation, gender, marital status, ethnicity, health status, disability and nationality. The data generated by SM sites are increasing in the geometrical proportion daily. About 7.7 billion population of the world following approximate population are actively connected on one social site or the other as shown in figure 1.1

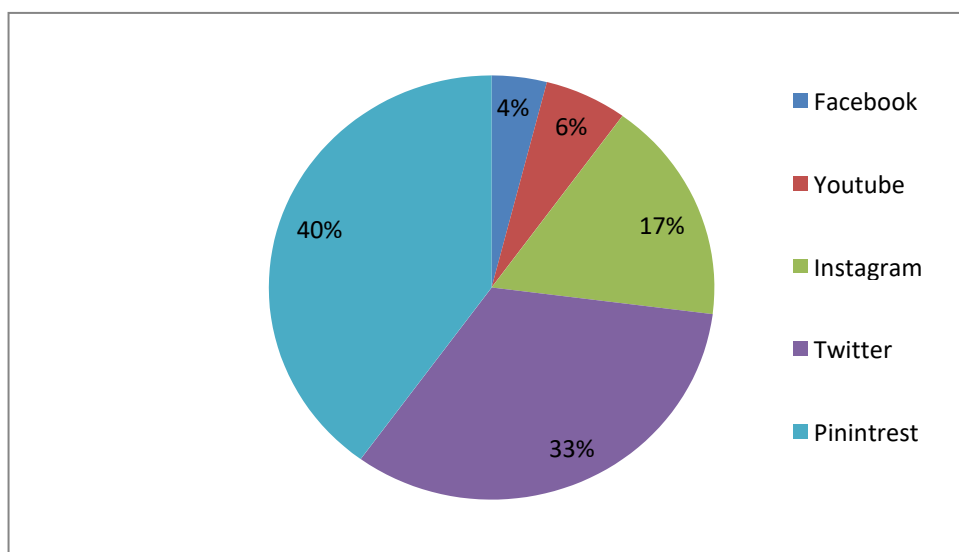


Fig 1.1: Percentages of active users

### 1.1 Objective

Social media is an extremely popular form of communication today. People offer their opinions and insights on a variety of topics, including politics, video games, and their personal lives. These platforms are occasionally used by some people to propagate false information about another person or group of people. The term “hate statement” refers to this kind of offensive material. One of the most well-known social media platforms is Twitter. But many people also use Twitter to disseminate offensive material. It is very hard to manually weed out abusive comments from the hundreds of millions of tweets that are generated every day on Twitter. Therefore, these offensive tweets ought to be automatically filtered out. In this study, we are developing an LSTM model for categorizing tweets as either containing hate content or not. To improve classification in social media texts such as hate speech or non-hate speech, researchers and practitioners require an updated understanding of machine learning methodologies, which is evolving rapidly. These methodologies help in creating a new and effective model that have better features to capture hate speech on social media. Slangs and new vocabularies are also evolving in social media. New and updated datasets are also available across the database. To bridge the gap, there is a need to review the old and traditional technologies into latest emerging and updated technologies.

### 1.2 Existing System

The automatic detection of hate speech using machine learning approaches is relatively new, and there are very limited review papers on techniques for automatic hate speech detection. The recent and related survey papers available on review of hate speech detection methods during this research work were few. The following were the available traditional literature review related to automatic detection of hate speech using ML algorithms.

ML algorithms have contributed immensely in hate speech detection and SM content analysis generally. Offensive comments such as HS and cyberbullying are the most researched areas in NLP in the past few decades. ML algorithms have been of great help in this direction in terms of SM data analysis for the identification and classification of offensive comments. The advances in ML algorithms researches have made significant impacts in many fields of endeavor which led to some important tools and models for analyzing a large amount of data in real-world problems like SMNs content analysis. In this survey conducted by the authors presented a brief review on eight hate speech detection techniques and approaches. These eight techniques include TF-IDF, dictionaries, N-gram, sentiment analyses, template-based approach, part of speech, Bag of the word, and rule-based approach. The limitation of the review is that techniques such as deep learning and ensemble approach were not considered in their work.

The authors offered a brief and critical analysis of the areas of automated hate speech detection in natural language processing. The authors also analyzed the features for hate speech detection in literature which includes: simple surface features, word generalization, sentiment analysis, lexical resources, linguistic features, knowledge-based features, meta-information and multimodal information.

The limitation of these two reviews is that techniques such as deep learning and ensemble approach are not considered in their work. The most significant step in text classification pipeline is selection of the best classifier. Therefore, the need to review all techniques is of essence. We intent to make this selection phase easier for researchers by reviewing more algorithms than the previous review work have covered. In this case, we reviewed techniques like deep learning, ensemble learning among others that have been employed for the automatic detection of hate speech in social media.

#### Existing System Disadvantages

1. The accuracy is less.
2. It is significantly slower due to an operation such as max pool.
3. Training of RNN models are difficult.

### 1.3 Proposed System

#### LSTM (long short term memory)

Long short-term memory (LSTM) network is a recurrent neural network (RNN), aimed to deal with the vanishing gradient problem present in traditional RNNs. Its relative insensitivity to gap length is its advantage over other RNNs, hidden Markov models and other sequence learning methods. It aims to provide a short-term memory for RNN that can last thousands of time steps, thus "long short-term memory" It is applicable to classification, processing and predicting data based on time series, such as in handwriting, speech recognition, machine translation, speech activity detection, robot control, video games and healthcare.

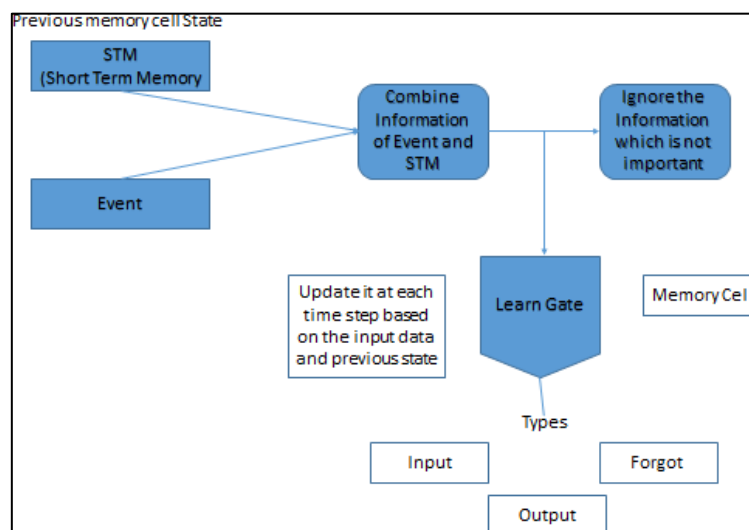


Fig 1.2: LSTM Architecture

A common LSTM unit is composed of a cell, an input gate, an output gate and a forget gate. The cell remembers values over arbitrary time intervals and the three gates regulate the flow of information into and out of the cell. Forget gates decide what information to discard from a previous state by assigning a previous state, compared to a current input, a value between 0 and 1. Value 1 means to keep the information, and a value of 0 means to discard it. Input gates decide which pieces of new information to store in the current state, using the same system as forget gates. Output gates control which pieces of information in the current

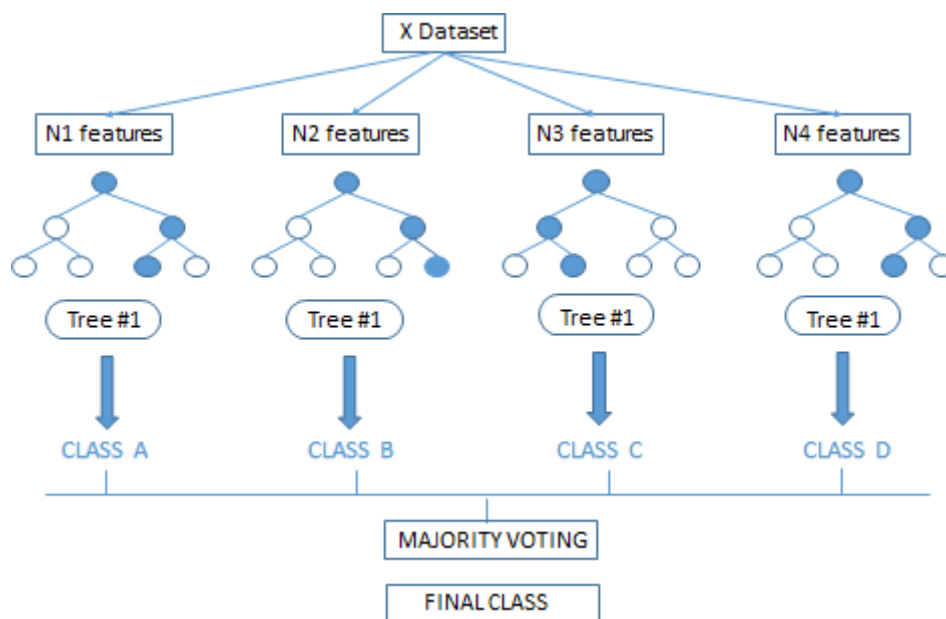
state to output by assigning a value from 0 to 1 to the information, considering the previous and current states. Selectively outputting relevant information from the current state allows the LSTM network to maintain useful, long-term dependencies to make predictions, both in current and future time-steps.

## Challenges

One of the main challenges of using LSTM for action recognition is the data requirements. LSTM needs a large amount of labeled video data to train effectively and generalize well to new scenarios. However, collecting and annotating video data is time-consuming, expensive, and prone to errors and inconsistencies

## RF (random forest)

The random forest algorithm is made up of a collection of decision trees, and each tree in the ensemble is comprised of a data sample drawn from a training set with replacement, called the bootstrap sample. Of that training sample, one-third of it is set aside as test data, known as the out-of-bag sample. Another instance of randomness is then injected through feature bagging, adding more diversity to the dataset and reducing the correlation among decision trees. Depending on the type of problem, the determination of the prediction will vary. For a regression task, the individual decision trees will be averaged, and for a classification task, a majority vote i.e. the most frequent categorical variable will yield the predicted class. Finally, the out-of-bag sample is then used for cross-validation, finalizing that prediction.



### Proposed System Advantages

1. Easy to predict.
2. It is very effective even with high dimensional data.
3. It can be used for both regression and classification problem.

## Challenges

Time-consuming process: Since random forest algorithms can handle large data sets, they can provide more accurate predictions, but can be slow to process data as they are computing data for each individual decision tree. Requires more resources as random forests possess larger data sets, they require more resources to store the data. More complex: The prediction of a single decision tree is easier to interpret when compared to more trees.

## 2. Literature Survey

| S.No. | Title/Description  | Objective   | Strategy  | Remarks   |
|-------|--|---|---|---|
| 1     | Advances in machine learning algorithms for hate speech detection in social media Nanlir Sallau Mullah.e.tal IEEE [2021]                         | Examining the basic baseline components of hate speech classification using ML algorithms.  | Different variants of ML techniques are reviewed which include classical ML, ensemble approach and deep learning methods.             | The limitation of this work is that no experiment was conducted with a given dataset  |
| 2     | Machine learning based automatic hate speech recognition system. P.William.et.al IEEE [2022]   | Examining datasets with different feature engineering techniques and machine learning algorithms.   | Datasets training with feature engineering techniques and machine learning algorithms.  | Support vector machine technique on testing showed best using bigram feature dataset.   |
| 3     | Hate speech detection in social media for the Kurdish language Ari M. Saeed.et.al Springer [2022]  | Detecting the hate speech in Kurdish language.  | Support vector machine, decision tree and naïve bays algorithms are implemented.  | Support vector machine showed the excellent result when compared with decision tree and naïve bays with 0.687.  |
| 4     | HCovBi-Caps hate speech detection using convolutional and bi-directional gated recurrent unit with capsule network Shakir Khan.et.al IEEE [2022] | This study presents a novel Convolutional, BiGRU, and Capsule network- based deep learningmodel, CovBi-Caps, to classify the hate speech. | Convolutional neural network, BiGRU and Capsule network- based deep learningmodel, HCovBi- Caps are implemented.                      | HCovBi-Caps show comparatively better performance over the unbalanceddataset with 0.93.   |
| 5     | Detection of hate speech texts using machine learning algorithm Mahamat Saleh Adoum Sanoussi.et.al IEEE [2022]                                   | Detection of hate speech for texts written in “lingua franca”, a mix of the local Chadianand French languages.                            | The four Machine Learning methods namely Logistic Regression, Support Vector Machine, Random Forest, andK-Nearest Neighbors are used. | The result showed that FastText features given as input to SVM classifier shown thebest accuracy of 95.4%.  |
| 6     | Deep learning for hate speech detectionin social media Ashwini Kumar.et.alIEEE [2021]  | Used a benchmark dataset of approximately 25 thousand annotated tweets to classify hate speech.   | Deep learning methods are implemented to classify the model.  | The deep learning methods are compared with traditional methods was measured in terms of f1 score and accuracy.   |
| 7     | Un-Compromised credibility social media based multi- class hate speech classification for textKhubaib Ahmed Khureshi.et.al IEEE [2021]           | A specific dataset availability and its high-performing supervised classifierfor text-based is addressed.                                 | Classification algorithms are usedto classify differentforms of datasets.   | Due to the application of latentsemantic analysis for dimensionality reduction the utilization of many complex and non- linear models and CAT Boost performed best. |



### 3. System Design

#### 3.1 Feature Categorization

Initially we have created the dataset of 14,640 annotated tweets which have been collected from twitter raw data. These tweets are being trained to classify the text speech as positive or negative i.e hate or non hate. These tweets are categorized according to its parameters and features of the speech. The tweets contain three columns namely class label, tweet id and text. As the system understand and implement the data as 0 or 1 i.e 0 represent positive and 1 represent negative. The dataset length is 200 and the size of dataset is 3093kb. The following table shows the dataset table

| airline_sentiment |          | text  |
|-------------------|----------|---|
| 0                 | neutral  | @VirginAmerica What @dhepburn said.               |
| 1                 | positive | @VirginAmerica plus you've added commercials t... |
| 2                 | neutral  | @VirginAmerica I didn't today... Must mean I n... |
| 3                 | negative | @VirginAmerica it's really aggressive to blast... |
| 4                 | negative | @VirginAmerica and it's a really big bad thing... |

Among 14,640 annotated tweets 2 are neutral, 2,264 are positive and 12,374 are negative. The following table depicts the categorization of dataset

| dataset | Positive | negative | neutral |
|---------|----------|----------|---------|
| 14,640  | 2,264    | 12,374   | 2       |

**Table 4.1 : dataset table**

The dataset is being iterated and tested with different ratios to show the best accuracy of the model. The following table shows the iteration levels of the model.

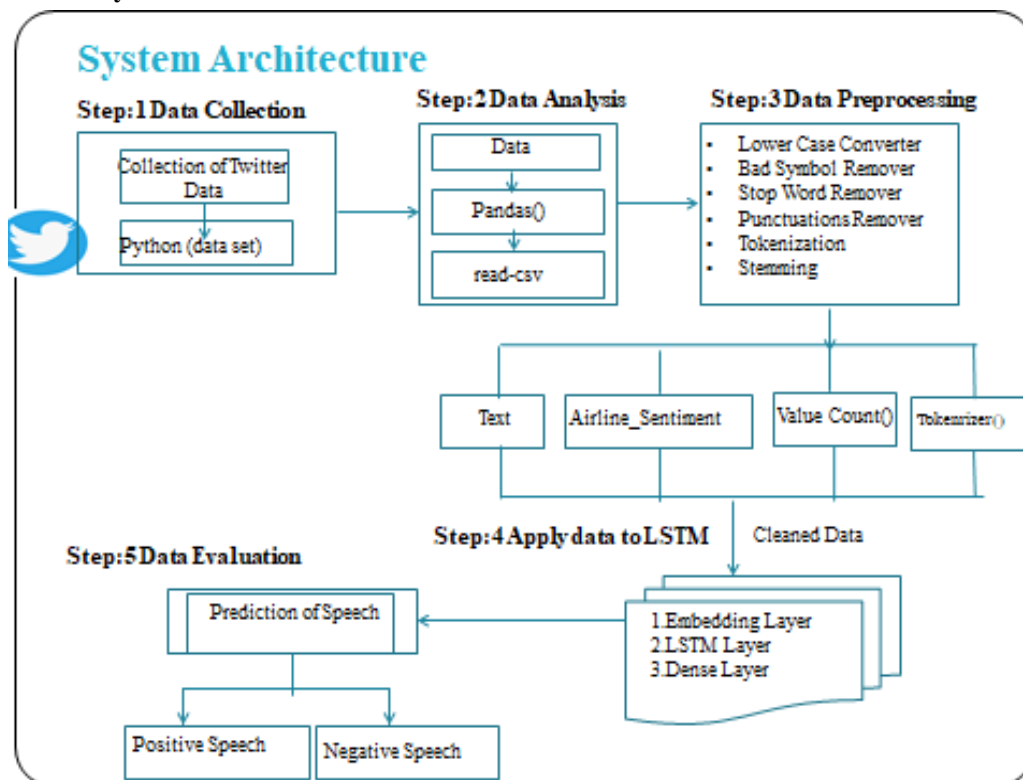
```

history = model.fit(padded_sequence,sentiment_label[0],validation_split=0.2, epochs=5, batch_
size=32) Epoch 1/5
289/289 [=====] - 43s 140ms/step - loss: 0.4068 - accuracy:
0.8310 - val_loss: 0.2246 - val_accuracy: 0.9092 Epoch 2/5
289/289 [=====] - 41s 142ms/step - loss: 0.2193 - accuracy:
0.9161 - val_loss: 0.1730 - val_accuracy: 0.9338 Epoch 3/5
289/289 [=====] - 40s 139ms/step - loss: 0.1616 - accuracy:
0.9392 - val_loss: 0.1690 - val_accuracy: 0.9381 Epoch 4/5
289/289 [=====] - 39s 135ms/step - loss: 0.1441 - accuracy:
0.9485 - val_loss: 0.1718 - val_accuracy: 0.9390 Epoch 5/5
289/289 [=====] - 38s 133ms/step - loss: 0.1191 - accuracy:
0.9569 - val_loss: 0.1679 - val_accuracy: 0.9438

```

**Fig 4.1: Iteration of the model**

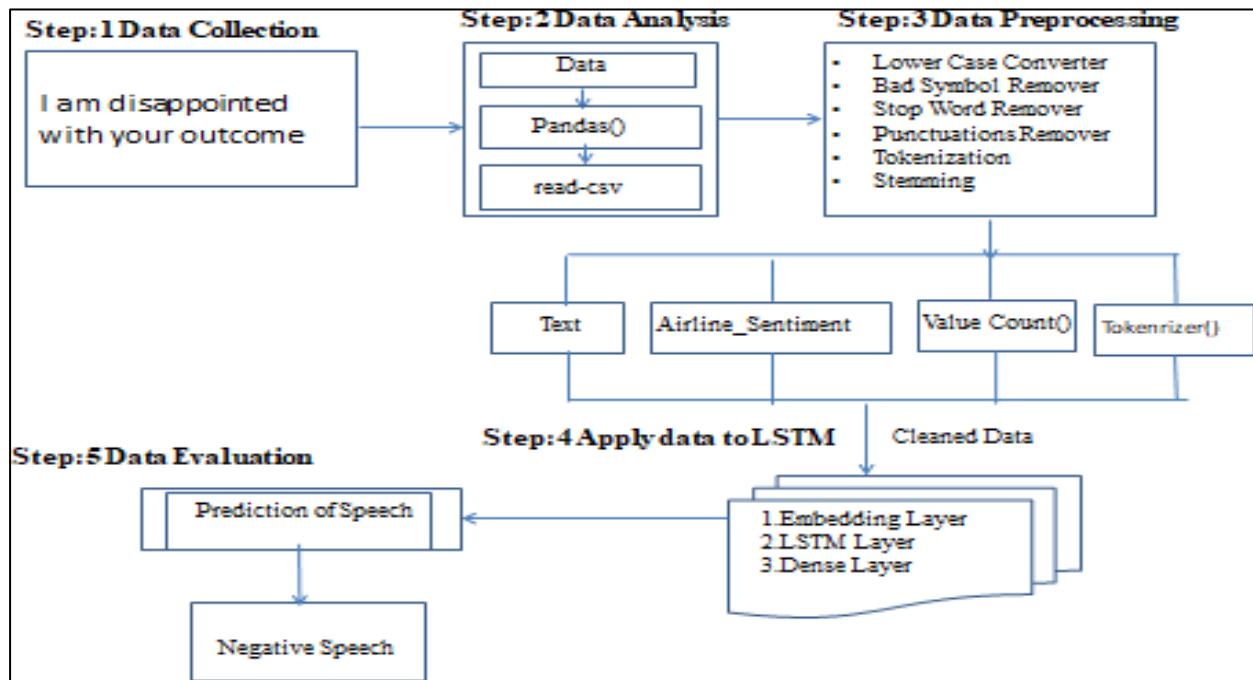
### 3.2 System Architecture



**Fig 4.2: system Architecture**

Initially the data is collected from social media twitter (raw data) and the data is transferred to python files through interface and create a python dataset file. This dataset is analyzed by data analysis in the second step, in this stage the data is able to read the python csv (comma separated value) files through method pandas (). After the data analysis the data is preprocessed i.e the data is cleaned or remove any unnecessary data or bad symbol remover, punctuation errors, lower case, upper case etc, Tokenization and stemming is also done in data preprocessing. It also prints the text, airline sentiment, value count and tweet id. In the next stage the cleaned data is applied to the LSTM in which it consists of three layers in which the first layer is embedding layer where the size of data is compressed to reduce errors and easy to implement then the next layer is LSTM layer where LSTMs are typically used with a 3-layer architecture: an input layer, an LSTM layer, and an output layer. The input layer converts the input sequence into a vector representation. The LSTM layer learns to identify the important features in the input sequence and to learn long-term dependencies between them. The output layer outputs a prediction for the task at hand, such as the next word in a sentence or the class of a document and simultaneously control, update and store the data. The next layer is dense layer in which the percentage of accuracy is calculated between each level. Then at last the data is evaluated and the speech is predicted as positive or negative.





**Fig 4.3: Example of system architecture**

1) The cleaned data is then applied to LSTM in which it consists of 3 layers namely 1) Embedding layer-which compresses the data, 2) LSTM layer-In which data is modified and compared with trained files and 3) Dense layer- This data calculates the ratio internally levels.

2) Finally the data is evaluated and predicted the speech either positive or negative

### 3.3 Pseudo code

Input Twitter dataset

Output Prediction of hate speech

1. Read twitter dataset through python files, pandas // Read python csv files
2. {
3. Import Pandas()
4. Load dataset
5. Normalize the dataset into values from 0 to 1 // the dataset is converted into integer values with tokenizer() module
6. Split dataset into train and test sets
7. Set input units, output units, lstm units and optimizer // the lstm model take input strings optimize the data and produce the output throughgates
8. **for** epochs and batch size **do**
9. Train the LSTM network
10. **End for**
11. Make predictions // prediction of speech either positive or negative

12. Calculate the Accuracy, Precision, F1 Score, Recall

13. }

#### 4. Results

##### 4.1 Evaluation Measures

We must ensure the integrity of our outcomes. As a result, we conducted a quantitative analysis of the frequency of data collection, data release mode, and data types. The obtained data is meticulously annotated before being entered into the constructed model for prediction. So we used an evaluation measure to examine the data and outcomes to see how well the model worked, how skewed the results were, and how generalizable our findings were. We employed the evaluation metric we devised throughout the trial. The idea of positive and negative affects all of the numbers we assess for accuracy, precision, recall, and F1 scores. Negative speech is defined as hate speech, whereas Positive speech is defined as speech that is not hateful. Figure shows the definitions of True Positives (TP), True Negatives (TN), False Positives (FP) and False Negatives (FN).

**PRECISION (Pr)** Precision is the ratio of true positive and total predictions. The following researchers made use of precision to evaluate their model performance. This can be represented mathematically stated as:

$$\text{Pr} = \text{TP} / \text{TP} + \text{FP} \quad (1)$$

Pr is a short for precision for the purpose of this study. Precision simply means a fraction of positive classifications that was correctly identified by the model. For example, the proportion of actual positives that were identified correctly from the example above is 4. Then the model precision is  $4/6$  (true positives / all positives) = 0.67. TP is a short for true positive. From the scenario above, TP is 4. Out of 5 hate speech tweets, the model was able to correctly identify 4 as hate speech. FP means false positive. This refers to non-hate speech tweets that were classified as hate speech. From the scenario above, 2 tweets were missed classified as hate speech tweets and in the real sense, they were non-hate speech tweets.

##### **RECALL (Rc )**

Rc is the ratio of the number of correct predictions and all correct observation in the sample space made use of recall for their evaluation. Mathematically stated as:

$$\text{Rc} = \text{TP} / \text{TP} + \text{FN} \quad (2)$$

Rc stands for Recall in this paper. This means the proportion of real positives that were established correctly. From the scenario, recall is  $4/5$  (true positives / all positives) = 0.8. This means the model was able to correctly identify 80% of the hate tweets. FN stands for false negative for the purpose of this study. This refers to those hate speech tweets that were not identified by the model as hate speech. The model considered them as non-hate while they were hate tweets in the real sense. In the example above, only one tweet was misclassified as non-hate and was actually hate speech.

### F-MEASURE

F-measure (F) or F1-score (F) is simply the weighted harmonic mean of precision and recall. This evaluation metric is normally employed when the dataset is unbalanced. It was employed to evaluate performance of hate speech prediction model. Mathematically stated as:

$$F = 2 * Pr * Rc / Pr + Rc \quad (3)$$

F is short for F-measure or F1-score and is used to test the model's performance with an imbalanced class distribution. In most real-life text classification tasks, imbalanced class distribution occurs and hence F1-score is a smarter metric to test a model. From example above,  $F = 2(0.67*0.8)/(0.67 + 0.8) = 1.072/1.47 = 0.72$ . This simply means that the F1-measure of the model is 72%.

### ACCURACY (A)

Accuracy is the ratio of correct prediction and total observations. Accuracy of a model is considered best if and only if we have symmetric dataset in which the value of FP and FN are almost equal for the two-class problem. Accuracy is not the best option in multiple and imbalanced data sets, hence, other evaluation parameters may be considered, like F1-score. In the following researches, accuracy was used. Mathematically, accuracy (A) can be expressed as:

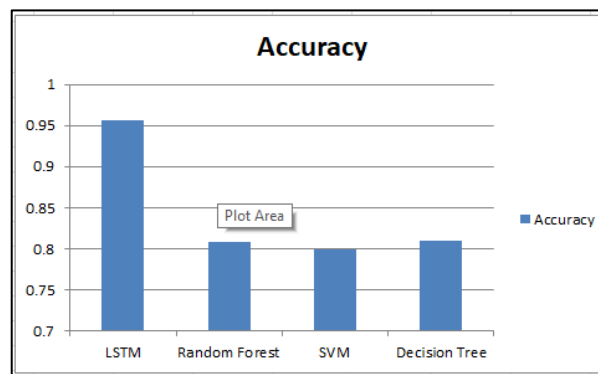
$$A = TP + TN / TP + FP + FN + TN \quad (4)$$

After implementing the evaluation measures we get the following percentages in the below table.

| Techniques    | Accuracy | Accuracy% |
|---------------|----------|-----------|
| LSTM          | 0.9569   | 95.69     |
| Random Forest | 0.8078   | 80.78     |
| SVM           | 0.7992   | 79.92     |
| Decision Tree | 0.8096   | 80.96     |

**Table 6.1: accuracy of models**

### 4.2 Comparison Analysis



**Fig 6.1: Comparison of models**

Figure 6.1 depicts the comparison analysis of models in accuracy in proposed model. In my analysis I get best results in LSTM model when compared with other models (RF, SVM & DT).

## 5. Conclusion

In the computing domain we have proposed LSTM (long short term memory) and RF (random forest) to create a website that classify the speeches into positive or negative speech. This can be achieved by creating the dataset from twitter and training the dataset by different methods. These data is analyzed and create files through pandas. These files are saved as .csv (comma separated value) after the data preprocessing it prints all columns by columns method but only two columns are taken according to the project text and airline sentiment. These datasets when trained will predict positive or negative and every word is given unique ids by tokenizer. These csv files or dataset or tweets are applied to LSTM architecture layers 1) Embedding 2) LSTM 3) Dense. This trained data or modified LSTM is saved as .h5 file and the same process is done for RF and this trained data is saved as .pkl file. The app.py is a flask to construct a router interconnection between frontend and backend. To accept the file it use request method and to send it use render\_template method this app.py is a constructor to run the website to classify the speech is positive or negative. The precision rate of accuracy of proposed LSTM is 95.74% and for RF is 80.88% when compared with other techniques.

## 6. Future Scope

In the future scope, the application of machine learning for automatic hate speech detection on social media needs to be encouraged and supported. The hate speech variables based on each country is an issue that needs more researchers' attention. Each country or region has different variables for hate speech. We should work more on advancement features on automatic hate speech detection to classify different forms of text, images, audios, videos, graphs etc. We are also looking for different kinds of languages and datasets of any form to be trained to obtain good accuracy with any computing language. We could also plan for some robotic machines to justify the classification speech of any language. etc; In the computing domain we have proposed LSTM (long short term memory) and RF (random forest) to create a website that classify the speeches into positive or negative speech. This can be achieved by creating the dataset from twitter and training the dataset by different methods. These data is analyzed and create files through pandas. These files are saved as .csv (comma separated value) after the data preprocessing it prints all columns by columns method but only two columns are taken according to the project text and airline sentiment. These datasets when trained will predict positive or negative and every word is given unique ids by tokenize. These csv files or dataset or tweets are applied to LSTM architecture layers 1) Embedding 2) LSTM 3) Dense. This trained data or modified LSTM is saved as .h5 file and the same process is done for RF and this trained data is saved as .pkl file. The app.py is a flask to construct a router interconnection between frontend and backend. To accept the file it use request method and to send it use render template method this app.py is a constructor to run the website to classify the speech is positive or negative. The precision rate of accuracy of proposed LSTM is 95.74% and for RF is 80.88% when compared with other techniques.

## References

- 1) M. Sajjad, F. Zulifqar, M. U. G. Khan and M. Azeem, "Hate Speech Detection using Fusion Approach," 2019 International Conference on Applied and Engineering Mathematics (ICAEM), 2019, pp. 251-255, doi: 10.1109/ICAEM.2019.8853762.
- 2) Das, A. K., Al Asif, A., Paul, A., & Hossain, M. N. (2021). Bangla hate speech detection on social media using attentionbased recurrent neural network. *Journal of Intelligent Systems*, 30(1), 578-591.
- 3) Nugroho, K., Noersasongko, E., Fanani, A. Z., & Basuki, R. S. (2019, July). Improving random forest method to detect hatespeech and offensive word. In 2019 International Conference on Information and Communications Technology (ICOIACT) (pp. 514-518). IEEE.
- 4) Alzubaidi, L., Zhang, J., Humaidi, A. J., Al-Dujaili, A., Duan, Y., Al-Shamma, O., ... & Farhan, L. (2021). Review of deep learning: Concepts, CNN architectures, challenges, applications, future directions. *Journal of big Data*, 8(1), 1- 474 A. Verma et al. / GMSARN International Journal 17 (2023) 468-474 74.
- 5) Abro, S., Shaikh, S., Khand, Z. H., Zafar, A., Khan, S., & Mujtaba, G. (2020). Automatic hate speech detection using machine learning: A comparative study. *International Journal of Advanced Computer Science and Applications*, 11(8).
- 6) Davidson, T., Warmley, D., Macy, M., & Weber, I. (2017, May). Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media* (Vol. 11, No. 1, pp. 512-515).
- 7) Sazany, E., & Budi, I. (2018, September). Deep learningbased implementation of hate speech identification on texts in indonesian: Preliminary study. In 2018 International Conference on Applied Information Technology and Innovation (ICAITI) (pp. 114-117). IEEE.
- 8) Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.
- 9) Zhang, Z., Robinson, D., & Tepper, J. (2018, June). Detecting hate speech on twitter using a convolution-gru based deep neural network. In *European semantic web conference* (pp. 745-760). Springer, Cham..
- 10) Joachims, T. (1998, April). Text categorization with support vector machines: Learning with many relevant features. In *European conference on machine learning* (pp. 137-142). Springer, Berlin, Heidelberg..
- 11) Severyn, A., & Moschitti, A. (2015, August). Twitter sentiment analysis with deep convolutional neural networks. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval* (pp. 959-962)..
- 12) Krishnan, H., Elayidom, M. S., & Santhanakrishnan, T. (2017). Emotion detection of tweets using naive bayes classifier. *Emotion*, 4(11), 457-62..
- 13) Fatahillah, N. R., Suryati, P., & Haryawan, C. (2017, November). Implementation of Naive Bayes classifier algorithm on social media (Twitter) to the teaching of Indonesian hate speech. In 2017 International Conference on Sustainable Information Engineering and Technology (SIET) (pp. 128-131). IEEE.
- 14) Kiilu, K. K., Okeyo, G., Rimiru, R., & Ogada, K. (2018). Using Naïve Bayes algorithm in detection of hate tweets. *International Journal of Scientific and Research Publications*, 8(3), 99-107..
- 15) Malmasi, S., & Zampieri, M. (2017). Detecting hate speech in social media. *arXiv preprint arXiv:1712.06427*.
- 16) Chandrasekharan, E., Samory, M., Srinivasan, A., & Gilbert, E. (2017, May). The bag of communities: Identifying abusive behavior online with preexisting internet data. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (pp. 3175-3187).

- 17) Burnap, P., & Williams, M. L. (2015). Cyber hate speech on twitter: An application of machine classification and statistical modeling for policy and decision making. *Policy & internet*, 7(2), 223-242.
- 18) Waseem, Z., & Hovy, D. (2016, June). Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop* (pp. 88-93).
- 19) Badjatiya, P., Gupta, S., Gupta, M., & Varma, V. (2017, April). Deep learning for hate speech detection in tweets. In *Proceedings of the 26th international conference on World Wide Web companion* (pp. 759-760).
- 20) Uoc, N. Q., Duong, N. T., Son, L. A., & Thanh, B. D. A (2022) Novel Automatic Detecting System for Cucumber Disease Based on the Convolution Neural Network Algorithm. *GMSARN International Journal* 16 (2022) 295- 301
- 21) KS, A. K., Sarita, K., Kumar, S., Saket, R. K., & Swami, A.(2022) Machine Learning-based Approach for Prevention of COVID-19 using Steam Vaporizer *GMSARN International Journal* 16 (2022) 399-404.
- 22) Gupta, C. L., Bihari, A., & Tripathi, S. (2021). Protein Classification Using Machine Learning and Statistical Techniques. *Recent Advances in Computer Science and Communications (Formerly: Recent Patents on Computer Science)*, 14(5), 1616-1632.
- 23) Ali, M.M., Qaseem, M.S., Ahmad, S.S. (2023). Rumour Detection Model for Political Tweets Using ANN. In: Kumar, A., Ghinea, G., Merugu, S. (eds) *Proceedings of the 2nd International Conference on Cognitive and Intelligent Computing. ICCIC 2022. Cognitive Science and Technology*. Springer, Singapore. [https://doi.org/10.1007/978-981-99-2742-5\\_15](https://doi.org/10.1007/978-981-99-2742-5_15)
- 24) Syed Shabbeer Ahmad, & Shreyas Jagadeep Shete. (2022). Innovative Deep Learning-Based Medical Report Analysis for Timely Diagnosis and Improved Healthcare. *Sparklinglight Transactions on Artificial Intelligence and Quantum Computing (STAIQC)*, 2(2), 16–28. Retrieved from <https://www.sparklinglightpublisher.com/index.php/slp/article/view/53>