

**Detecting stance in Code-Mixed Hindi-English Social Media  
Data Using Deep Learning**

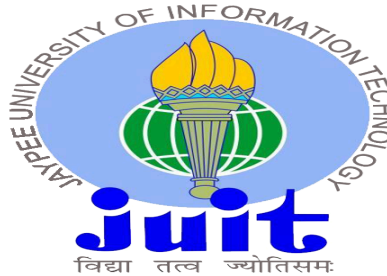
A major project report submitted in partial fulfillment of the  
requirement for the award of degree of  
**Bachelor of Technology**  
in  
**Computer Science & Engineering / Information Technology**  
*Submitted by*

**Puneet Katoch (201385)**

**Anubhav Thakur (201231)**

*Under the guidance & supervision of*

**Mr. Maneet Singh**



**Department of Computer Science & Engineering and  
Information Technology**  
**Jaypee University of Information Technology,**  
**Waknaghat, Solan - 173234 (India)**

# I

## Certificate

We hereby declare that the work presented in this report entitled “**Detecting stance in Code Mixed Hindi-English Social Media Data Using Deep Learning**” in partial fulfilment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering/Information Technology** under the supervision and guidance of **Dr. Maneet Singh (Associate Professor)** submitted in the department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology, Wazirpur, Lucknow is an authentic record of our own work carried out over a period from February 2024 to June 2024.

Puneet Katoch  
201385

Anubhav Thakur  
201231

This is to certify that the above statement made by the candidate is true to the best of my knowledge.

Dr. Maneet Singh  
Assistant Professor (SG)  
Computer Science & Engineering  
15/05/24

**JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT**  
**PLAGIARISM VERIFICATION REPORT**

Date: .....

Type of Document (Tick):  PhD Thesis  M.Tech Dissertation/ Report  B.Tech Project Report  Paper

Name: \_\_\_\_\_ Department: \_\_\_\_\_ Enrolment No \_\_\_\_\_

Contact No. \_\_\_\_\_ E-mail. \_\_\_\_\_

Name of the Supervisor: \_\_\_\_\_

Title of the Thesis/Dissertation/Project Report/Paper (In Capital letters): \_\_\_\_\_

**UNDERTAKING**

I undertake that I am aware of the plagiarism related norms/ regulations, if I found guilty of any plagiarism and copyright violations in the above thesis/report even after award of degree, the University reserves the rights to withdraw/ revoke my degree/report. Kindly allow me to avail Plagiarism verification report for the document mentioned above.

**Complete Thesis/Report Pages Detail:**

- Total No. of Pages =
- Total No. of Preliminary pages =
- Total No. of pages accommodate bibliography/references =

(Signature of Student)

**FOR DEPARTMENT USE**

We have checked the thesis/report as per norms and found **Similarity Index** at .....(%). Therefore, we are forwarding the complete thesis/report for final plagiarism check. The plagiarism verification report may be handed over to the candidate.

(Signature of Guide/Supervisor)

Signature of HOD

**FOR LRC USE**

The above document was scanned for plagiarism check. The outcome of the same is reported below:

Copy Received on	Excluded	Similarity Index (%)	Generated Plagiarism Report Details (Title, Abstract & Chapters)	
	<ul style="list-style-type: none"> <li>• All Preliminary Pages</li> <li>• Bibliography/Images/Quotes</li> <li>• 14 Words String</li> </ul>		Word Counts	
<b>Report Generated on</b>			Character Counts	
		<b>Submission ID</b>	Total Pages Scanned	
			File Size	

Checked by  
Name & Signature

Librarian

**Please send your complete thesis/report in (PDF) with Title Page, Abstract and Chapters in (Word File) through the supervisor at [plagcheck.juit@gmail.com](mailto:plagcheck.juit@gmail.com)**

# Candidate's Declaration

We hereby declare that the work presented in this report entitled “**Detecting stance in Code Mixed Hindi-English Social Media Data Using Deep Learning**’ in partial fulfilment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science & Engineering / Information Technology** submitted in the Department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology, Waknaghat is an authentic record of my own work carried out over a period from August 2023 to May 2024 under the supervision of **Dr. Maneet Singh**(Associate Professor, Department of Computer Science & Engineering and Information Technology).

The matter embodied in the report has not been submitted for the award of any other degree or diploma.

Student Name: Puneet Katoch

Roll No.: 201385

Student Name: Anubhav Thakur

Roll No.: 201231

This is to certify that the above statement made by the candidate is true to the best of my knowledge.

Supervisor Name: Dr. Maneet Singh

Designation: Associate Professor

Department: Computer Science & Engineering and Information Technology

## II

### ACKNOWLEDGEMENT

Firstly, we express our heartiest thanks and gratefulness to almighty God for His divine blessing makes it possible for us to complete the project work successfully.

We are really grateful and wish our profound indebtedness to Supervisor **Mr Maneet Singh** , **Associate Professor, Department of CSE Jaypee University of Information Technology, Wagnaghat**. Deep Knowledge & keen interest of our supervisor in the field of “Deep Learning to carry out this project. His endless patience, scholarly guidance, continual encouragement, constant and energetic supervision, constructive criticism, valuable advice, reading many inferior drafts and correcting them at all stages have made it possible to complete this project.

We would also generously welcome each one of those individuals who have helped us straight forwardly or in a roundabout way in making this project a win.

Finally, We must acknowledge with due respect the constant support and patience of our parents.

Puneet Katoch(201385)  
Anubhav Thakur(201231)

### III

## TABLE OF CONTENTS

<b>S. No.</b>	<b>Title</b>	<b>Page No.</b>
<b>1.</b>	<b>Declaration</b>	<b>I</b>
<b>2.</b>	<b>Certificate</b>	<b>II</b>
<b>3.</b>	<b>Acknowledgement</b>	<b>III</b>
<b>4.</b>	<b>Abstract</b>	<b>IV</b>
<b>5.</b>	<b>CHAPTER 1: INTRODUCTION</b>	<b>01 - 13</b>
<b>6.</b>	<b>CHAPTER 02: LITERATURE SURVEY</b>	<b>14 - 17</b>
<b>7.</b>	<b>CHAPTER 03: SYSTEM DEVELOPMENT</b>	<b>18 - 42</b>
<b>8.</b>	<b>CHAPTER 04: PERFORMANCE ANALYSIS</b>	<b>43 - 47</b>
<b>9.</b>	<b>CHAPTER 05: Conclusion</b>	<b>48 - 50</b>
<b>10.</b>	<b>References</b>	<b>51 - 53</b>

<b>Abbreviation</b>	<b>Name</b>
SVM	Support Vector Machine
NN	Neural Network
ANN	Artificial Neural Network
SC	Softmax Classifier
GDO	Gradient Descent Optimisation
BFGSO	Broyden–Fletcher–Goldfarb–Shanno Optimisation
PCA	Principal Component Analysis
ML	Machine Learning
SDLC	Software Development Life Cycle
VS	Visual Studio
EDA	Exploratory Data Analysis

## **ABSTRACT**

This research uses deep learning techniques to solve the problem of encoding mixed data, often related to single-word communication of multiple languages. We propose a new deep learning architecture that combines iterative processes and attention to preserve complex patterns and contexts in complex programming contexts.

In the age of global communication, social media platforms have become a great place where users can express their thoughts, feelings and positions in different languages. This research focuses on the challenging task of visual analysis in mixed Hindi-English social media documents where users combine both languages. Behavior analysis plays an important role in understanding the behavior and thoughts expressed by users of a topic or site.

The data for this study was selected from popular social media platforms and has coding examples where Hindi and English content are mixed in various ways. The first method is used to solve problems caused by mixed data, such as different word and text changes.

The proposed model is trained and evaluated using the test model containing the following compound language nuances. Experimental results demonstrate the effectiveness of deep learning in identifying user locations. The model's performance is compared to the baseline, demonstrating its ability to outperform traditional methods and highlighting the importance of using neural network architectures for understanding compound words.

This research contributes to the popularization of natural language processing, especially in the context of social analysis across different cultures. The design not only improves our understanding of working in a code hybrid context, but also paves the way for applications of sentiment analysis, sentiment mining, and content recommendations to accommodate aspects of code hybrid communication.



# 1. INTRODUCTION

## 1.1 INTRODUCTION

In the developing social environment where people from different languages come together and share their feelings and thoughts, mixing has become more and more important. Synthesis, which involves combining multiple languages into a single communicative example, poses a unique challenge for natural language processing (NLP). In a private hybrid environment, understanding the user's stance (i.e., their behavior or opinion on a topic or place) is important to recognize penalties for behaviors expressed on these platforms.

This research focuses on the task of detecting mixed language in Hindi-English social media profiles where users mix two different languages. Social media posts in Hindi and English demonstrate the variety of speech patterns found in various cultural contexts and provide difficult examples for statistical analysis. Word order changes, syntactic patterns, and inconsistent text transformations are just a few of the complex compound words in code that traditional NLP models struggle to capture.

To address these issues, we employ deep learning, a field that has shown remarkable success in numerous NLP projects. Deep learning models have the potential to capture subtleties and dependencies in data, especially when they combine listening techniques with repetition. Our goal is to create a deep learning system that can combine nuances in code to assess speech accuracy and detail.

The data set for the study, which covers a wide range of subjects and domains, was gathered from well-known social media platforms. To deliver a solid model across multiple languages, use

cutting-edge methods to handle intricate procedures such as tokenization, script analysis, and optimization.

As we delve deeper into our methods and experiments, we hope to demonstrate how well our deep learning approach captures subtle cues and patterns, expressing the user stance in the mixed content of the rules. Furthermore, we compare how well the model performs to the baseline approach, emphasizing the need for improved techniques to address specific issues caused by complex language processing.

As we delve deeper into our methods and experiments, we hope to demonstrate how well our deep learning approach captures subtle cues and patterns, expressing the user stance in the mixed content of the rules. Furthermore, we compare how well the model performs to the baseline approach, emphasizing the need for improved techniques to address specific issues caused by complex language processing.

This research advances NLP while also having broader implications for understanding communication across linguistic and cultural barriers. Our findings highlight the complexities of conflicting messages in social media policies, paving the way for additional cognitive analysis, opinion mining, and content recommendations (if any) to help us understand digital discourse in a variety of contexts.

## 1.1 PROBLEM STATEMENT

The widespread use of social media has changed the field of communication and allowed individuals to express different opinions and thoughts. In the context of multiculturalism, users often mix codes to combine different languages into a single communication instance. Understanding the nuances expressed in hash code is a challenge that cannot be effectively solved by natural language processing (NLP).

Code mixing reveals complex languages such as changing words, changing letters and syntactic structures. These challenges hinder the performance of traditional NLP models and make it difficult to accurately predict user stance in social media contexts where Hindi-English codes are mixed.

To deal with the complexity of complex language, there is an urgent need for decision models that can learn and represent the complex language present in social networks. Deep learning about the ability to capture hierarchical representations and dependency contexts holds promise for improving statistical analysis in mixed contexts.

Lack of knowledge of specific content specific to digital-hybrid Hindi-English social media hinders the development and evaluation of computational models. This shortcoming makes it difficult to train robust models that can be generalized to multiple languages.

This research addresses these issues by proposing and applying a deep learning approach suited to the complexity of mixed Hindi-English social media data. By doing so, it aims to contribute to advances in natural language processing in a multilingual environment and to foster a deeper understanding of user paths in discussion in complex contexts.

## 1.2 OBJECTIVE

### **Develop a deep learning model for visual search:**

Develop and implement a deep learning model tailored to the specific challenges of encoding mixed Hindi-English social media data. The model should take advantage of the nuances of presentation in mixed content.

### **Research Topics and Interests:**

Exploring and integrating Relational Neural Networks (RNN) and monitoring techniques in proposed models to detect the effects of physical and contextual nuances present in the mixed language of the programming process. This process should improve the model's ability to determine location in the dynamic language.

### **Long and Advanced Code Mix Collection:**

Writing and editing a collection of real-life codes with mixed content in Hindi and English from social media platforms. Use advanced techniques to solve problems such as letter substitutions, different word orders, and features unique to compound words.

### **Train and validate model:**

Train deep learning model datasets of selected models and optimize them to achieve high accuracy in pose detection. Follow the validation process that determines appropriate metrics for issues caused by mixing languages in your code, such as changes in text and language usage.

### **Comparison with baseline models:**

Evaluate the performance of deep learning models compared to baseline models, including NLP techniques. A comparative analysis was conducted to demonstrate the superiority of the deep learning method in handling differences in mixed social data. > Ensure that deep learning techniques are robust in capturing events across multiple contexts, reflecting the diversity of user-generated content.

**Regarding ethical decisions:**

Know and address ethical issues related to social media analysis, ensuring research adheres to the principles of user privacy and data user responsibility. Take steps to reduce bias and maintain ethical standards when collecting and analyzing composite data.

**Provides insight into multilingual NLP:**

Provides insight into the broader field of multilingual processing by demonstrating the feasibility of interactive learning models, the powerful effect of keeping words together. It provides useful findings that may inform future research in the field of computational mathematics and NLP in various languages.

**Advances in social media understanding:**

Demonstrating the impact of real-time verification on code hybrid social media data. Explore powerful data sources like sentiment analysis, sentiment mining, and content recommendations that show real-world insights from deep learning models.

**Disseminate research results:**

Disseminate research results through academic publications, conferences and other communication methods. Contribute to the research community's understanding of mixed language processing in the programming process and lay the foundation for future advances in research in the multilingual community.

## 1.3 Significance and Motivation of the Project Work

The significance and motivation of a project on stance detection using Deep Learning addresses critical challenges in various domains and leveraging advanced technologies for practical applications. Here are key aspects of the project's significance and motivation:

### 1.) Understanding Multilingual Social Dynamics:

**Significance:** In multilingual societies, individuals often engage in code-mixing, reflecting the dynamic linguistic diversity of their communication. Understanding the complex social dynamics in these bilingual environments requires identifying stance in code-mixed content.

**Motivation:** This study employs deep learning techniques to decipher the complex web of opinions and attitudes expressed in code-mixed Hindi-English social media conversations in order to contribute to a better understanding of multilingual social dynamics.

### 2.) Enhancing Stance Detection Accuracy:

**Significance:** Understanding user perspectives on various topics requires the identification of stance. Traditional models struggle to capture subtleties in code-mixed languages.

**Motivation:** Deep learning offers a promising solution to address the limitations of existing methods and significantly improve the accuracy of stance in code-mixed social media data by modeling complex linguistic patterns.

### 3.) Facilitating Targeted Content Recommendations:

**Significance:** Accurate stance detection can aid in the recommendation of specific content in addition to aiding in the understanding of user opinions. This is crucial for social media platforms in particular, as providing tailored content is vital.

**Motivation:** The development of a deep learning model that accurately identifies stances in code-mixed content could lead to the creation of more complex content recommendation algorithms. On multilingual platforms, this will increase user satisfaction and engagement levels.

#### **4.)Enabling Nuanced Sentiment Analysis:**

**Significance:** Understanding user attitudes and emotions requires the use of sentiment analysis. Sentiment analysis models must be adjusted for code-mixed language due to the special difficulties it presents.

**Motivation:** Deep learning is used in this study to enable more nuanced sentiment analysis in code-mixed social media data by attempting to capture the subtle differences in user emotions and sentiments expressed in Hindi and English.

#### **5.)Insights into Societal Attitudes:**

**Significance:** It is reflective of social media attitudes and trends. Understanding how multilingual communities work together to solve problems can be improved by doing an analysis of mixed-language realities.

**Motivation:** This study uses deep learning to identify complex concepts in the code and thus obtain a deeper understanding of social behavior, with the goal of revealing the subtleties that shape public opinion.

**6.)Best Performance in a Versatile Environment:**

**Significance:**In order to gain a deeper understanding of social behavior and uncover the nuances that influence public opinion, this study employs deep learning to identify complex concepts in the code.

**Motivation:** This research helps create models that can be extended across languages, which promotes advancement, with a focus on Hindi-English hybrid content in particular.

**7.)According to the needs of digital communication platforms:**

**Significance:** Given the increasing popularity of digital communication platforms, it is critical to comprehend others, collaborate effectively, and use appropriate language.

**Motivation:** Create deep learning models to analyze content in response to changing digital communication needs, tools to improve user experience, and multilingual interactions.



## 2: LITERATURE SURVEY

### 2.1 Overview of Relevant Literature

S . N O	Paper Title (cite)	Journal / Confere nce (year)	Tools/ Techniques/ Dataset	Results	Limitations
1.	“Named Entity Recognition for Hindi-English Code-Mixed Social Media Text”	Language Technologies Research Centre (LTRC), IIT Hyderabad	LSTM CRF and NC4K dataset used.	Machine learning models which achieved the best f1-score of 0.95 with both CRF and LSTM.	Lexically similar words with different tags makes the learning phase of our model difficult and hence some incorrect tagging of the tokens
2.	“Sentiment identification in Code-mixed social media Text”	Rutgers University, New Brunswick, NJ 08901, US	NLP tools HECM DATASETS	System shows an accuracy of 68.5% . On instances of neutral polarity, we can achieve an accuracy of 55.2%.	As the dataset is relatively small, we would like to create a larger dataset in future.

3.	Code mixing in language style as communication through social media twitter	ICCS 2021 Meliani	CNN model BiLSTM model twitter extracted datasets	Through training different models we achieve an accuracy of 91.2%	The use of code-mixing do not determines by gender, age, or level strata. It naturally and intentionally occurs among the language speaker who can communicate in different languages
4.	“Sentiment analysis of code mixed social media data (SA-CM SMT) in Indian Languages	IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) in 2023	HECM Computer vision Deep learning MoCA dataset used.	MoCA DATASET gives us Accuracy of 95.6%	Evaluated on a single dataset, accuracy lower in real world scenarios.
5.	Sentiment Analysis for Hinglish Code mixed tweets by means of cross-lingual word embeddings	LREC 11-16 May 2020	SGD Algorithm DND datasets twitter extracted datasets	The transfer learning experiments result in an F1-score of 0.556 which is almost on par with the supervised settings and speak to the robustness of the cross-lingu	We believe the cross-lingual embeddings can still be improved, as the embeddings constructed now are generic and can be further tailored with domain information to increase performance.

				al embedding approach.	
6.	Annotate corpus creation for sentiment analysis in code mixed Hinglish social network data	IEEE Conference on Computer Vision and Pattern Recognition (CVPR) in 2022.	VecMAP toolkit BiLSTM model twitter extracted datasets	The authors have proposed two-way method for sentiment analysis of Hinglish Data: one by using Lexicon based approach with an accuracy of 86% and secondly by using machine learning techniques with an accuracy of 76%	The method can be extended to the combination of various languages such as Gujrati-English, Punjabi-English Bhojpuri-English and Marwari-English
7.	Stance Detection in Hindi Mixed Data Using Multi Task Learning	International Institute of Information Technology, Hyderabad, 2023	Module Enhancement Module (FEM) . Twitter scrapped datasets	best model achieved the result with a stance prediction accuracy of 63.2% which is a 4.5% overall accuracy improvem	Automated assessment f the stance can not be done

				ent compared to the current supervised classified	
--	--	--	--	--	--

## 2.2 Key Gaps in the Literature

### 1.) Limited Code-Mixed Datasets:

**Gap:** The availability of comprehensive and diverse datasets specifically tailored for code-mixed Hindi-English social media data may still be limited.

**Potential Research Direction:** Creating and sharing standardized datasets that cover a wide range of topics, domains, and linguistic variations in code-mixed content would be beneficial for the development and evaluation of deep learning models.

### 2.) Ethical Considerations and Bias:

**Gap:** There might be gaps in the literature regarding the ethical considerations and potential biases associated with the analysis of code-mixed social media data.

**Potential Research Direction:** Further exploration into ethical considerations, including privacy issues, potential biases in models, and the impact of cultural nuances on stance detection in code-mixed content, is essential.

### 3.) Interplay of Linguistic Features:

**Gap:** The interplay of linguistic features in code-mixed content, especially the interaction between Hindi and English linguistic elements, might not be fully explored.

**Potential Research Direction:** Investigating the intricate linguistic patterns, such as the influence of language switching, syntactic variations, and cultural expressions, could contribute to a more nuanced understanding of code-mixed language in social media.

### 4.) Model Interpretability in Code-Mixed Contexts:

**Gap:** The interpretability of deep learning models in the context of code-mixed language might not be well-explored.

**Potential Research Direction:** Research focused on developing interpretable deep learning models or methodologies to interpret model decisions in code-mixed contexts could enhance the trustworthiness and usability of these models.

### **5.)Generalization Across Social Media Platforms:**

**Gap:** The generalization of deep learning models for code-mixed stance detection across various social media platforms may not be thoroughly studied.

**Potential Research Direction:** Investigating the transferability and robustness of models trained on one platform to others, considering variations in user behavior, linguistic styles, and topics.

### **6.)Combining Linguistic and Socio-Cultural Contexts:**

**Gap:** Integrating socio-cultural context with linguistic features for better stance detection in code-mixed content might not be fully explored.

**Potential Research Direction:** Exploring the impact of socio-cultural factors, including cultural references, contextual cues, and regional variations, in enhancing the performance of deep learning models for stance detection.

### **7.)Real-Time Stance Detection:**

**Gap:** The literature may have limited exploration of real-time or near real-time detection of stances in dynamic social media conversations.

**Potential Research Direction:** Developing models or methodologies that can efficiently process and analyze code-mixed content in real-time, considering the fast-paced nature of social media interactions.

## **3: SYSTEM DEVELOPMENT**

### **3.1) REQUIREMENT AND ANALYSIS**

Several essential prerequisites and analytical methods must be followed in order to use deep learning to detect stance in Code-Mixed Hindi-English social media data. The process of ascertaining a speaker's attitude or viewpoint toward a specific subject, issue, or topic is known as "stance detection." In code-mixed data, two or more languages—in this case, Hindi and English—are combined. An overview of the requirements and analysis procedure is provided below:

#### **1.)Gathering of Data:**

Compile a varied dataset of social media posts that are Code-Mixed Hindi-English. Posts from Facebook, Twitter, and other pertinent social media outlets should be included in this.

Add stance labels to the dataset. Labels for stances may fall into groups like neutral, opposition, or support.

#### **2.)Prior to processing:**

Clean up and tokenize the text data. Take care of things like special characters, hashtags, mentions, and punctuation.

To distinguish between the Hindi and English portions of the code-mixed text, do language identification.

#### **3.)Incorporated Representation:**

For both Hindi and English, use word embeddings that have already been trained. You can use embeddings such as Word2Vec, GloVe, or FastText.

Examine methods for managing code-mixed embeddings to make sure the model is able to comprehend the subtle differences in language between different languages.

#### **4.)Model of Architecture:**

Develop a deep learning model for posture recognition. Common architectures include transformer-based models, like BERT for contextual embeddings, long short-term memory networks (LSTMs), and recurrent neural networks (RNNs).

Adjust the model's architecture to take into consideration the mixed coding in the data. You may consider using a language's specific attention processes.

### **5.)Advice:**

Utilize the dataset to create training, validation, and test sets.

Train the deep learning model on the training set of data using a suitable loss function (categorical cross-entropy, for example) for stance detection.

Make changes to the hyperparameters to optimize the model's performance.

### **6.)Mixed-Code Difficulties:**

Address issues like language ambiguity, code-switching patterns, and user proficiency differences that are unique to code-mixed data.

### **7.)Moral Aspects to Take into Account:**

Make sure the model and dataset account for moral issues pertaining to privacy, justice, and bias in social media data.

### **8.)Implementation:**

Use fresh, untrained data to apply the learned model for real-time stance detection.

### **9.)Constant Enhancement:**

Think about ongoing model monitoring and development based on user input and changing social media language trends.

### **10.)Assessment:**

-Analyze the model using the validation set, and make any necessary adjustments.

-Analyze the model's performance using metrics such as F1 score, accuracy, precision, and recall on the test set.

-Analyze the model's performance on the English and Hindi portions independently.



## 3.2 Project Design and Architecture

A project's design that uses deep learning to identify stance in Code-Mixed Hindi-English social media data must carefully take into account a number of factors. An overview of the project's architecture and design can be found below:

### 3.2.1) Design of the Project:

#### Establish the project's objectives:

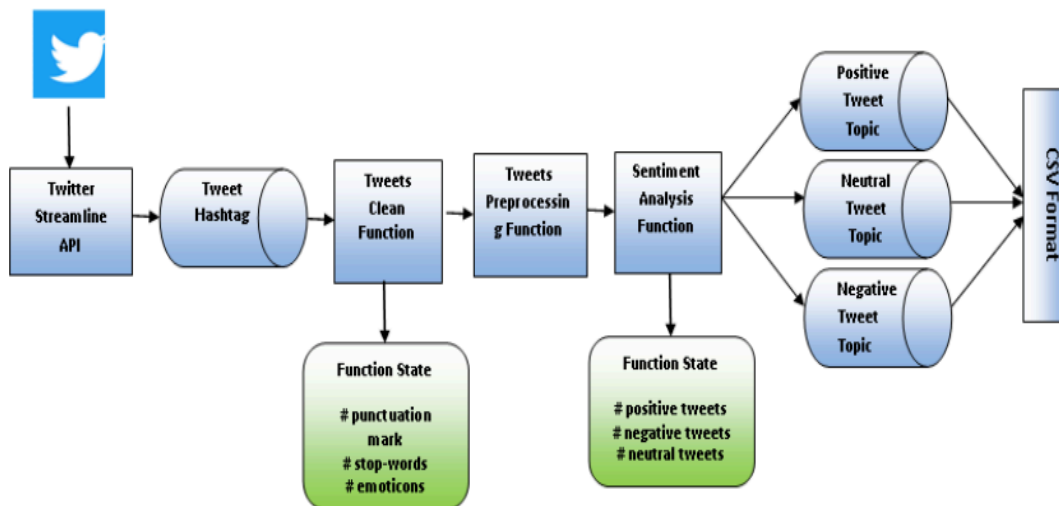
Make sure that the project's objective—to identify stance in code-mixed Hindi-English social media data—is clearly stated.

#### Recognize Stance Detection:

Learn about the literature and research that have already been done on stance detection, particularly with regard to social media data and code-mixing.

#### Determine the Stakeholders:

Determine who might be interested in sentiment analysis of Code-Mixed data: academics, social media analysts, or organizations.



Flowchart of project design

## **3.2.2)Project Architectures:**

### **Gathering and Preparing Data:**

Compile a varied dataset of social media posts that are Code-Mixed Hindi-English. Add stance labels to the dataset.

### **Prior to processing:**

Clean up and tokenize the text data.  
To differentiate between the Hindi and English portions, do language identification.

### **Incorporated Representation:**

For both Hindi and English, use word embeddings that have already been trained. Use code-mixed embeddings or other methods to deal with the data's mixed-language nature.

### **Model Architecture:**

Select an appropriate deep learning model for detection: Consider using a Transformer-based model such as Recurrent Neural Network (RNN), Long Short-Term Memory Network (LSTM), or BERT. Use special language listening techniques to decipher the jumbled code. Fine-tune the architectural model based on the characteristics of social media products.

### **Training:**

Divide the dataset into a training set, a validation set, and a testing set. Use appropriate loss to train the model for detection. Optimize hyperparameters for best performance.

### **Assessment:**

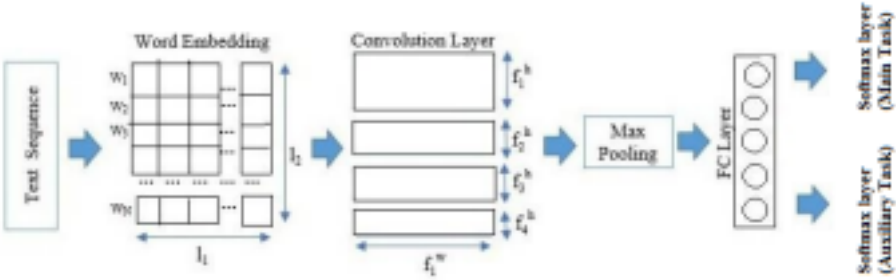
Evaluate the standard of the verification process and fine-tune as necessary. Evaluate the model's performance in the benchmark using metrics such as accuracy, precision, recall, and F1 score. Check the performance in Hindi and English section.

### **Ethical Considerations:**

Addresses issues specific to mixed data, such as ambiguous words and changing patterns. Consider ethical considerations regarding neutrality, honesty and confidentiality in social media profiles.

### **Information:**

Make sure to clearly document every step of the project, including preliminary steps, prototypes, and evaluations. Share policies and practices responsibly, considering consent and privacy issues.



CNN MTL Architecture

Model	Accuracy(%)
CNN	66.7
CNN + MTL	71.3

### **3.2.3)Data Preparation**

Data preparation is an important step in any machine learning process, especially when dealing with complex data like code-mixed Hindi-English social media texts.

#### **1. Data collection:**

Collection of different data on mixed Hindi-English social media posts.

To capture the diversity of social media content, ensure the data covers a wide range of topics, users, and different languages.

#### **2. Note:**

Use tags to fill in the records. Define categories based on the context of the project (for example, for, against, or intermediate projects).

#### **3. Text Cleaning:**

Remove noise and irrelevant information from text files. This may include:

Publishing unique characters, tags and URLs.

Get hashtags and comment accordingly.

Get emojis and emoticons.

#### **4. Grammar:**

Get grammar knowledge on differentiating Hindi and English in mixed text.

This is important to know the difference between the two languages and understand the context.

#### **5. Word segmentation:**

Split the text into words or sub-word units. This step is important to create a coherent strategy for deep learning models.

#### **6. Representational Embedding:**

Leveraging pre-trained word embeddings for Hindi and English. Embeddings such as Word2Vec, GloVe or FastText.

Discover ideas for manipulating number combinations or creating representations that include descriptions of letter combinations.

#### **7. Data classification:**

Divide the dataset into a training set, an evidence set and a testing set. The split will be 70-15-15%.

Make sure each group has an equal share of the list.

**8. Padding and segment length:**

Padding or trimming sequences to a long length to ensure size stability in deep learning models. Choose the appropriate length according to the average length of the material.

**9. Data curation (optional):**

Evaluate the data curation process, especially if your data is limited. Skills such as transcription or back-translation can help triangulate data. 10. Data analysis:

Check the distribution of records in the data. If there is a large gap, use strategies such as competition or failure to maintain class balance.

**11. Data format:**

Convert data into a format suitable for input into deep learning models. This is usually done to create a digital representation of the collected data.

**12. Data Analysis and Insights (optional):**

Search lighting data to gain insight into job classification, language patterns, and other relevant factors.

**13. Documentation:**

Record the entire data preparation process, including steps taken and any decisions made. This information is important to be reproduced and used in the future.

**14. Version Control (optional):**

If your documents need to change over time, use a control system to track differences between different documents.

By following the steps below, you will have well-prepared data to train deep learning models to identify verbs in mixed Hindi-English social media information. Be sure to review the data preparation process based on feedback from performance standards and outcome evaluations.

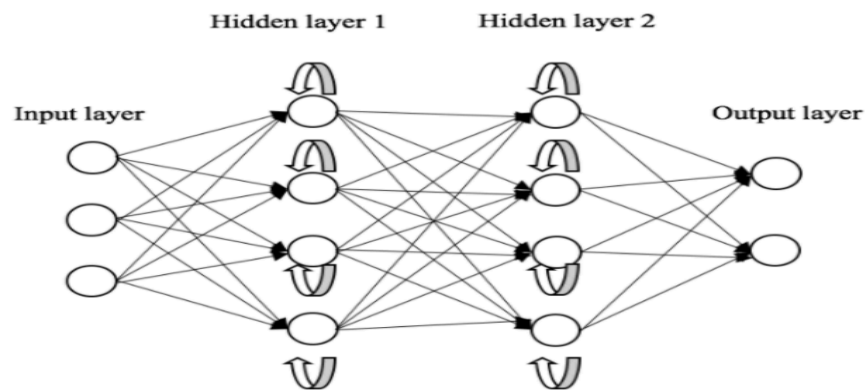
### 3.2.4)Implementation-

#### Tools and techniques-

#### Recurrent Neural Networks (RNNs):

RNNs, which include lengthy quick-time period reminiscence (LSTM) and Gated Recurrent Unit (GRU), can version sequential dependencies in code-blended text.

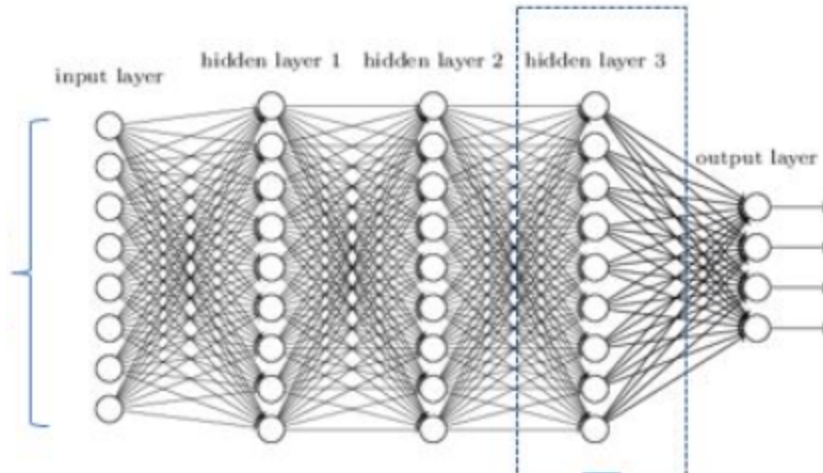
A Recurrent Neural community (RNN) is a sort of artificial neural community designed for sequence statistics and duties wherein the order of facts elements is vital. unlike traditional feedforward neural networks, RNNs have connections that form a directed cycle, allowing them to hold a hidden state that captures facts about previous inputs in the sequence.



#### Embeddings:

**Phrase Embeddings:** make use of pre-skilled word embeddings like Word2Vec, GloVe, or fastText to symbolize phrases in a non-stop vector area.

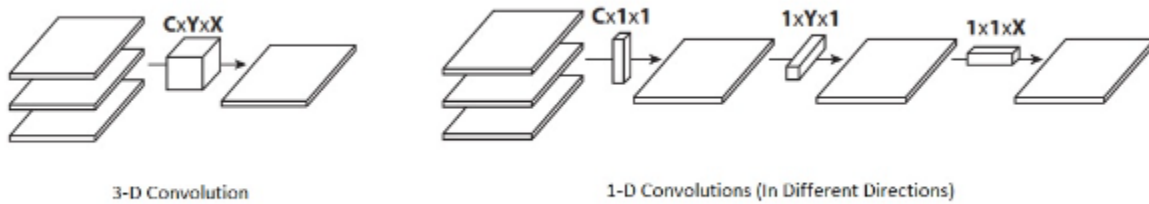
**Contextualized Embeddings:** appoint contextualized embeddings which include BERT (Bidirectional Encoder Representations from Transformers) or GPT (Generative Pre-educated Transformer) to seize contextual data.



**Convolutional Neural Networks (CNNs):**

A Convolutional Neural network (CNN) is a type of synthetic neural network designed for processing and studying visible facts, including snapshots. CNNs have been confirmed to be notably effective in computer imaginative and prescient responsibilities, such as image type, object detection, and photo segmentation. the key innovation of CNNs is the usage of convolutional layers, which allow the community to routinely study hierarchical representations of capabilities immediately from the raw pixel facts.

Apply CNNs to capture neighborhood patterns and dependencies in the input information, that's beneficial for figuring out stance in code-mixed text.



## **Hybrid models:**

Combine distinctive styles of neural networks, such as CNNs with LSTMs or GRUs, to capture both nearby and sequential styles.

### **-CNN-RNN Hybrids:**

Combining Convolutional Neural Networks (CNNs) with Recurrent Neural Networks (RNNs) is commonplace for tasks related to sequential statistics, inclusive of picture captioning or video evaluation. CNNs excel at extracting spatial functions from constant-length inputs like photographs, whilst RNNs are powerful in taking pictures temporal dependencies.

### **-CNN-LSTM Hybrids:**

long brief-term reminiscence (LSTM) networks, a type of RNN, may be included with CNNs for responsibilities where each spatial and temporal dependencies are vital. This combination is often utilized in video evaluation, movement recognition, and other time-collection duties.

### **CNN-interest Hybrids:**

Attention mechanisms, originally popularized in collection-to-sequence fashions, can be integrated with CNNs to allow the model to focus on particular regions of enter records. this is useful in responsibilities wherein positive components of the enter are greater relevant than others, such as photo captioning or visible query answering.

### **Transformer-based totally Hybrids:**

Transformers, which have gained prominence in herbal language processing, may be used in aggregate with different architectures. as an example, a Transformer can be included with a CNN for photograph classification or with an RNN for sequence-to-collection responsibilities.



## Ensemble models:

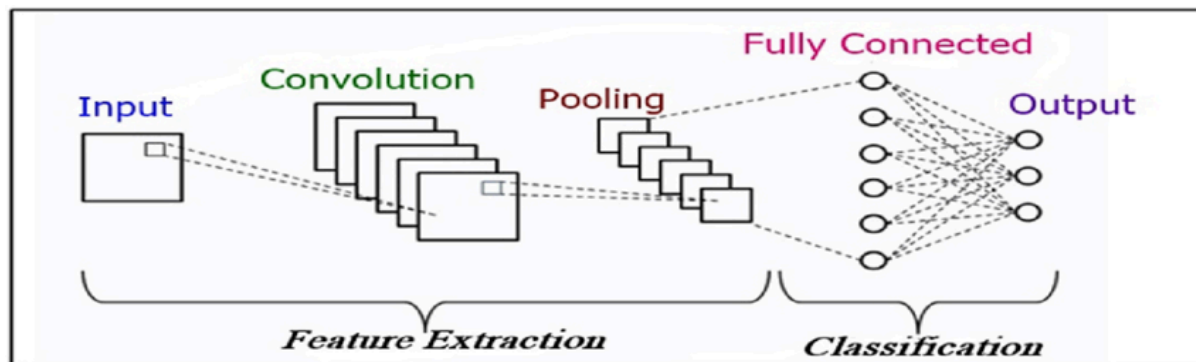
Ensembling involves combining predictions from a couple of fashions to improve universal performance. Hybrid ensembles can include a mixture of deep mastering models, conventional system learning models, or even models from exclusive domain names.

## Autoencoder-primarily based Hybrids:

Autoencoders, which can be used for unsupervised gaining knowledge of and function gaining knowledge of, can be included into larger architectures. as an instance, an autoencoder can be used for pre-education, and the found out representations may be high-quality-tuned for a particular challenge the usage of a distinctive neural community.

## Interest Mechanisms:

integrate interest mechanisms (e.g., self-attention or multi-head interest) to recognition on precise elements of the input collection, permitting the model to assign varying stages of significance to distinctive phrases.



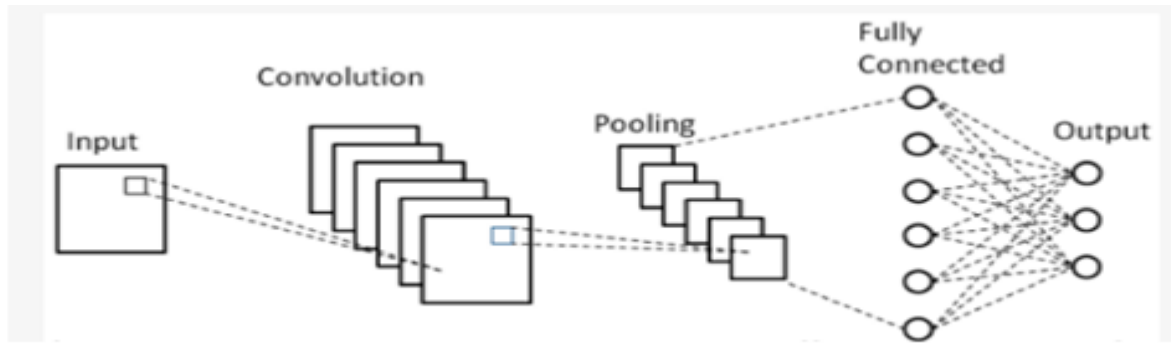


Fig. 1 Architecture of proposed System

### 3.1 )Algorithms

Deep learning (DL) is becoming an increasing number of tremendous in our each day lives. It has already had a big impact in fields together with voice popularity, self-using cars, precision medication, cancer detection, and forecasting. conventional mastering, class, and sample recognition strategies can not take care of huge-scale records units with their meticulously created characteristic extractors. the restrictions of previous shallow networks that impeded effective schooling and abstractions of hierarchical representations of multi-dimensional education statistics have been often solved by DL, relying at the complexity of the difficulty. multiple layers of gadgets with extraordinarily tuned algorithms and designs make up a deep neural community (DNN). in order to shorten training time and increase schooling accuracy, this examine explores a number of optimization techniques. We explore the arithmetic underlying the trendy deep network training techniques.

We define current issues, enhancements, and applications. The paper additionally discusses other deep designs, along with variational autoencoders, deep residual networks, deep convolutional networks, and recurrent neural networks.

With the introduction of the backpropagation mastering technique, DNN had a leap forward. It changed into first postulated in the Seventies , however it wasn't absolutely understood and implemented to neural networks till the center of the Eighties .

the subsequent classes may be used to institution distinct types of neural networks. 1.Feedforward neural community

#### **Recurrent Neural Networks (RNN)**

A neural network the usage of radial foundation features four Kohonen Self-Organizing Neural network

## **Modular Neural network.**

data in a feedforward neural network best actions in one route, thru any hidden nodes gift, from the input to the output layer. They do not create any loops or circles.

a selected form of multilayer feedforward neural network implementation with values and capabilities computed alongside the forward skip route is shown in figure 2a.  $Z$  denotes the non-linear activation characteristic  $f$  of  $Z$  at each layer, in which  $y$  represents the inputs' weighted total. The weights between the 2 devices within the next layers, represented through  $W$ , are the unfairness value of the unit is represented via the subscript letters and  $b$ .

The processing gadgets of RNN create a cycle, not like feedforward neural networks. A layer's output will become an input to the layer above it, that is often the only layer inside the community; as a end result, the layer's output turns into an enter to itself, creating a comments loop. This allows the community to save records about in advance states and use that records to have an effect on the present output. One essential end result of this difference is that, not like feedforward neural networks, RNNs can receive a series of inputs and generate a sequence of output values as nicely, making them extraordinarily beneficial for programs like speech reputation that call for the processing of a chain of time-phased input records.category of video frames-by way of-frame, etc. A RNN is shown unrolling in time in determine 2b.

As an illustration, the community would be opened up or unrolled 3 instances to create a 3-layer RNN if the enter have been a series of 3-word sentences.The diagram's mathematical justification is given beneath:The input at time  $t$  is represented by using  $x_t$ . The learnt parameters used by all phases are  $U$ ,  $V$ , and  $W$ . The output at time  $t$  is  $O_t$ . the following formula may be used to compute  $S_t$ , which stands for the state at time  $t$ , in which  $f$  is the activation feature, which includes ReLU.

several wonderful neural networks that have been controlled through an intermediate make up a modular neural community, which is an synthetic neural community.

each independent neural community acts as a module and tactics wonderful inputs to complete a particular element of the job the community is making an attempt to complete.

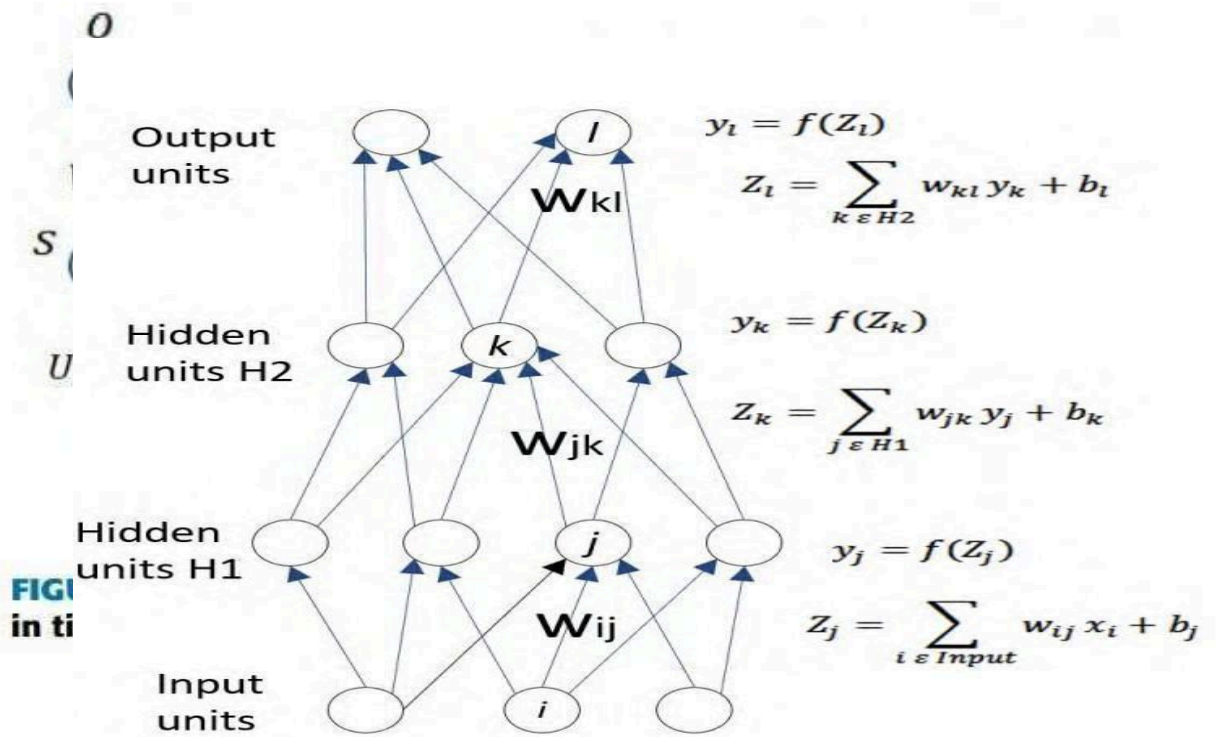


FIGURE  
in time

Radial basis methodology In category, carry out approximation, statistical prediction troubles, etc., neural networks are implemented. There are enter, hidden, and output layers in it. every node within the hidden layer is the cluster centre and consists of a radial foundation.

(carried out as a mathematician function). The output layer combines the outputs of the radial foundation perform and weight parameters to carry out classification or abstract concept once the network learns to assign the input to a centre .

Self-establishing Kohonen victimisation unattended learning, the neural community robotically arranges the network model into the computer record. it is created of associate diploma input layer companion degreed an output layer that at each completely related. The output layer is prepared up as a grid of two dimensions. The weights represent the houses (role) of the output layer node and there may be no activation carried out. Calculations are created referring to the weights to work out the geometrician distance among each output layer node and consequently the pc report. The method beneath updates the weights of the nearest node and its neighbours from the laptop record to convey them nearer to the laptop record.  $x(t)$  is that the computer report at time t,  $WI(t)$  is that the ith weight at time t, and Malaysia Militant institution is that the neighbourhood performs between the ith and jth nodes.

$$WI(t + 1) = WI(t) + \alpha(t)\eta_j * i(x(t) - WI(t))$$

huge networks are divided into smaller, freelance neural network modules victimisation preferred neural networks. The smaller networks do unique obligations which can be later mixed as one network output.

## 3.1 Model Development

### **Choose the model:**

You may select from a range of models looking on the goal you are attempting to realize. you'll use calculations for grouping, prediction, straight relapse, bunching, like k-means or K-Closest Neighbour, profound learning, additionally called brain organisations, bayesian, and then forth. Depending on the knowledge you may analyse, like pictures, sounds, text, and mathematical properties, there square measure many models that will be used. we'll inspect many models and their applications within the following table so you will use them in your comes.

### **Develop your machine model:**

You should originate the datasets to perform as anticipated and observe an identical rise within the foretold rate. As you train your model, the loads—which square measure the characteristics that replicate or have an effect on the linkages between the information sources and results—will inevitably vary. make sure to introduce your hundreds arbitrarily.

### **Assessment:**

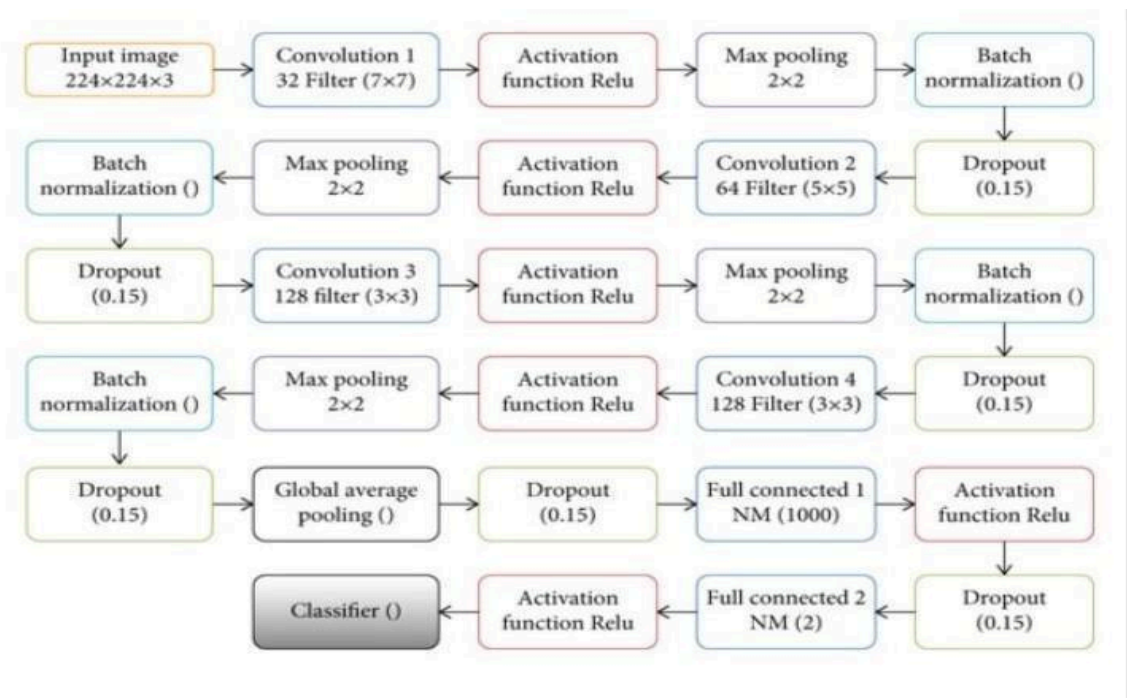
You should verify the accuracy of your totally made model by examining it in your analysis informative assortment, that includes inputs concerning that the model has no clue any. That approach will not be helpful as a result of it'd be comparable to moving a coin to create a choice, assuming the accuracy isn't precise or up to 0.5. you will place plenty of religion within the results the model predicts if you reach ninetieth or higher.

### **Parameter standardisation:**

Before making a brand new style of boundaries for your model's boundaries, you ought to come to the preparation step if throughout the assessment you did not receive the high expectations you were hoping for and your accuracy is not what you needed. during this case, it's possible that you simply have overfitting or underfitting problems. you'll construct the ages at which you repeat your preparational info. The "learning rate," that is commonly a price that repeats the slope to bit by bit take it nearer to the world - or at the terribly least to limit the expense of the capability - is another vital boundary.

Increasing your quality by zero.1 units from zero.001 isn't adored; this could considerably have an effect on how quickly a model runs.

You can additionally demonstrate the most important error that was created by victimising your model. Your machine's preparation time will vary from many seconds to hours or maybe days. These limits square measure often stated as hyperparameters. As you investigate, this "tune" can become even more of a piece of art than a science. Their square measure typically has plenty of borders to cross, and once they square measure all at once, they'll trigger all of your choices. every calculation has its own parameters which will be altered. to supply another example, you ought to describe in Counterfeit Brain Organizations (ANNs) the amount of secret layers it'll have and steadily take a look at roughly and with the amount of neurons in every layer. This task would require exceptional perseverance and energy so as to supply glorious results.





## 3.1 Python Tools

### 1) SciPy

Many Python programmers had been growing gadget learning libraries in Python as the sector accelerated speedy, appreciably for use in clinical and analytical computing. In 2001, Travis Oliphant, Eric Jones, and Pearu Peterson made the choice to mix and unify the general public of those disparate codes. SciPy library turned into the name given to the very last library.

A free BSD licence is used to offer the SciPy library, which is currently being evolved and supported through a public developer network.

The SciPy library offers modules for fixing normal differential equations (ODEs), sign and picture processing, special capabilities, speedy Fourier remodel, integration interpolation, and other computing obligations in science and analytics.

SciPy uses a multidimensional array as its base records shape, that's made to be had via the NumPy bundle. For the array manipulation subroutines, SciPy depends on NumPy. The SciPy library offers user-pleasant and effective numerical capabilities similarly to helping NumPy arrays.

The reality that SciPy's capabilities may be utilized in math and other fields is one in every of its one of a kind characteristics. sign processing, statistical, and optimization functions are some of its often utilised features. It includes tools for calculating integrals' numerical solutions. so that you're capable of optimise

SciPy is one of the maximum extensively used system studying libraries because of the applicability within the following fields. processing of multidimensional pictures, solves differential equations and Fourier transforms

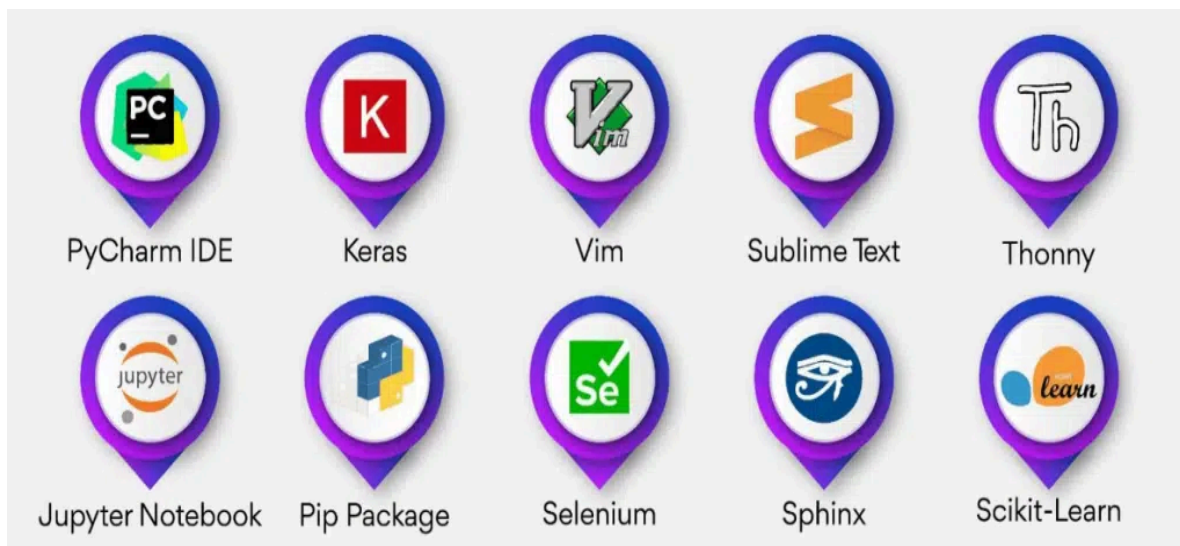
you could perform linear algebra calculations efficaciously and dependably way to its optimised algorithms.

### 2) Scikit-learn

David Cournapeau created the Scikit-analyze library as a issue of the Google summer time of Code project in 2007. INRIA participated in 2010 and accomplished the public launch in January 2010. Skikit-learn, the maximum broadly used Python system learning library for creating system studying algorithms, was constructed on top of Python libraries, NumPy and SciPy.

an expansion of supervised and unsupervised getting to know techniques are available in

Scikit-learn, which uses a standardised Python person interface. The library is likewise beneficial for statistics analysis and statistics mining. The Scikit-learn package deal is capable of coping with the following machine learning tasks: type, regression, clustering, dimensionality reduction, model selection, and preprocessing. Scikit-learn is a popular device for records scientists and ML aficionados.



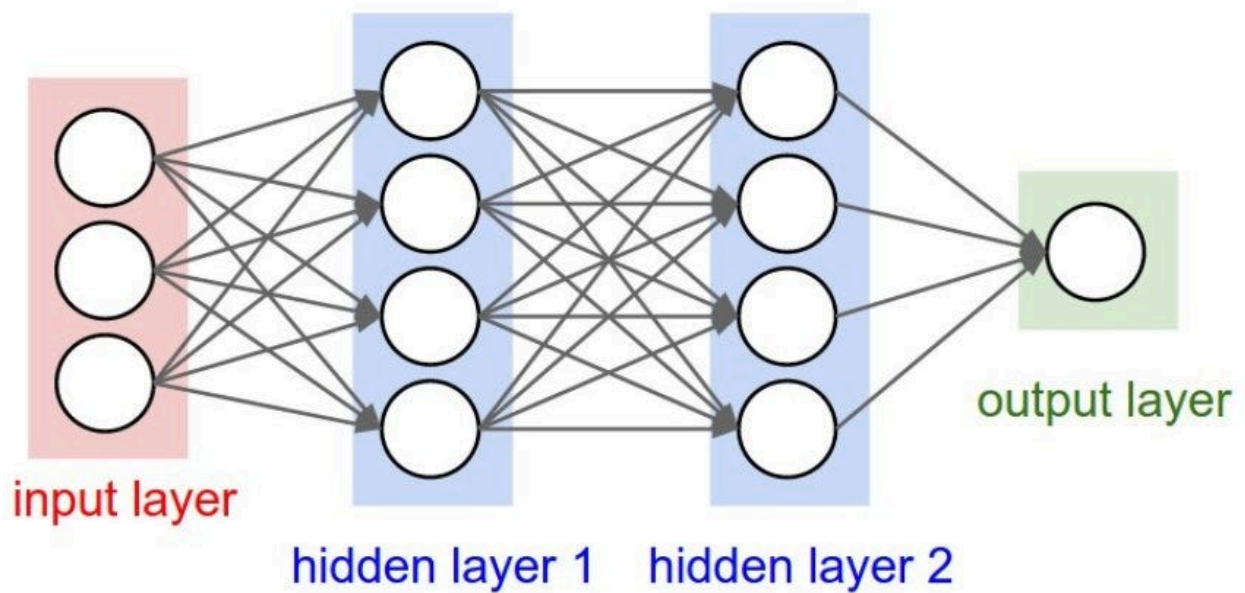
### 3) Theano

Theano turned into an open-source numerical computation library for Python that allowed developers to efficiently define, optimize, and evaluate mathematical expressions, in particular matrix-valued ones. However, as of September 28, 2017, the improvement and preservation of Theano were officially discontinued.

In case you are seeking to paintings with deep studying in Python, it is endorsed to apply greater current and actively maintained frameworks which includes TensorFlow or PyTorch. these frameworks offer big functionality, strong community guide, and integration with numerous excessive-stage APIs for constructing and schooling neural networks.

#### 4) Keras

Keras is an open-source deep learning framework for Python that gives a high-level neural networks API. It is designed to be consumer-friendly, modular, and extensible, making it a popular choice for both beginners and experienced deep learning practitioners. Keras acts as an interface for building and training neural networks on top of different numerical computation libraries, consisting of TensorFlow or Theano (although TensorFlow is now the default backend for Keras).



## 5) PyTorch

PyTorch is an open-source deep learning framework for Python that gives a flexible and dynamic computational graph, making it popular amongst researchers and practitioners within the discipline of synthetic intelligence. PyTorch is understood for its imperative programming style, which permits extra intuitive and dynamic model construction as compared to static graph frameworks like TensorFlow.



## 6) Pandas

Pandas is quickly becoming the maximum widely used Python library for facts evaluation due to the fact to its assist for brief, adaptable, and expressive information structures made to deal with each "relational" and "labelled" records. Pandas is a essential package deal for Python statistics analysis problems which are realistic and actual-global in nature these days. Pandas offers extremely dependable performance that is finely optimised. simplest C or Python is used to put in writing the backend code.

Pandas makes use of two number one classes of facts systems, which include:

collection (1-dimensional)

DataFrame (2-dimensional)

the majority of facts necessities and use instances throughout most sectors, along with science, facts, social work, finance, and, of route, analytics and different technological fields, may be treated via those two operating collectively.

the following forms of information are like minded with Pandas, and they also assist the others:

diverse information in a table's columns. don't forget the statistics in a square database or Excel spreadsheet, as an instance.

Time series records with and with out order. In comparison to different libraries and equipment, the frequency of time collection need now not be regular. Pandas is exceedingly capable at coping with time-collection facts that is unequal.

Heterogeneous or homogeneous kinds of facts may be gift within the rows and columns of an arbitrary matrix. any extra forms of observational or statistical facts sets. there may be virtually no want to label the statistics. without labels, the Pandas records shape can still handle it.

## CODE SNIPPETS-

**Extracting tweets from twitter using snsrape tool in Google Collab and then converting all the tweets into a single dataset.**

```
!pip install -q snsrape==0.3.4

[ ] import os
import pandas as pd
from datetime import date

[ ] today = date.today()
end_date = today

[ ] search_term = 'Bhavesh Bhatt Data Scientist'
from_date = '2019-01-01'
```

### Extracting Exact Tweets

```
[ ] max_results = 100

extracted_tweets = "snsrape --format '{content!r}'+ f" --max-results {max_results} --since {from_date} twitter-search '{search_term} until:{end_date}' > extracted-tweets.txt"
os.system(extracted_tweets)
if os.stat("extracted-tweets.txt").st_size == 0:
    print('No Tweets found')
else:
    df = pd.read_csv('extracted-tweets.txt', names=['content'])
    for row in df['content'].iteritems():
        print(row)

(""" 🚨 Tech Talk of the Week alert! 🚨 Learn about TensorFlow Hub by joining the session hosted by Machine Learning GDE, ' and Data Scientist', " Bhavesh Bhatt. 🚨 This event will be held on Sunday, April 18th 🚨 Tech Talk Time! 🚨 From 8 to 9PM GST! 🚨 Learn about TFHub! 🚨 TensorFlow Hub is an open repository & library for reusable machine learning models. 🚨 Data Science Career | How to Transition to Data Science with Data Scientist Bhavesh Bhatt | GreyAtom https://t.co/oc7CP6LCEM""", nan, nan)
```

```
[ ] os.system(f"snsrape --since {from_date} twitter-search '{search_term} until:{end_date}' > result-tweets.txt")
if os.stat("result-tweets.txt").st_size == 0:
    counter = 0
else:
    df = pd.read_csv('result-tweets.txt', names=['link'])
    counter = df.size

print('Number Of Tweets : '+ str(counter))
```

```

# Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import classification_report, confusion_matrix
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, LSTM, Dense, Bidirectional
from tensorflow.keras.utils import to_categorical

# Load your dataset
# Replace 'your_dataset.csv' with the path or link to your actual dataset
data = pd.read_csv('misinformation.csv')

# Preprocess the data
# Assuming 'text' is the column containing the social media text, and 'stance' is the stance label
texts = data['text']
stances = data['stance']

# Convert stance labels to numerical format
label_encoder = LabelEncoder()
stances_encoded = label_encoder.fit_transform(stances)
stances_encoded_categorical = to_categorical(stances_encoded)

# Tokenize and pad sequences
tokenizer = Tokenizer()
tokenizer.fit_on_texts(texts)
sequences = tokenizer.texts_to_sequences(texts)
padded_sequences = pad_sequences(sequences)

```

Using python we will use the dataset and will perform stance detection using python tools in google collab.

### DATASET DOWNLOADED -

mid	mis	keywords	titles	news_urls fact_check_urls
13	A Democr	Trump	(ship OR shipped OR sent	<a href="https://www.washingtonpost.com/politics/2020/04/22/did-trump-ship-17-tons-american-masks-china/">https://www.washingtonpost.com/politics/2020/04/22/did-trump-ship-17-tons-american-masks-china/</a>
17	Nobel laur	Luc Monta	"The Coro	<a href="https://www.healthfeedback.com/claimreview/claim-by-nobel-laureate-luc-montagnier-that-the-novel-coronavirus-is-man-made-and-contains-genetic-material-from-hiv-is-inaccurate/">https://www.healthfeedback.com/claimreview/claim-by-nobel-laureate-luc-montagnier-that-the-novel-coronavirus-is-man-made-and-contains-genetic-material-from-hiv-is-inaccurate/</a>
26	A video arj	coronaviu.	"Programn	<a href="https://www.healthfeedback.com/viral-video-promotes-the-unsupported-hypothesis-that-sars-cov-2-is-a-bioengineered-virus-released-from-a-wuhan-research-laboratory/">https://www.healthfeedback.com/viral-video-promotes-the-unsupported-hypothesis-that-sars-cov-2-is-a-bioengineered-virus-released-from-a-wuhan-research-laboratory/</a>
46	Michigan C	Michigan L	"Michigan	<a href="https://www.leadstories.com/hoax-alert/2020/04/Fact-Check-Michigan-Governor-Did-NOT-Ban-Sale-Of-American-Flags-During-Coronavirus-Outbreak.html">https://www.leadstories.com/hoax-alert/2020/04/Fact-Check-Michigan-Governor-Did-NOT-Ban-Sale-Of-American-Flags-During-Coronavirus-Outbreak.html</a>
51	Michigan C	Michigan L	"Michigan	<a href="https://www.politifact.com/factchecks/2020/apr/15/facebook-posts/covid-order-doesnt-ban-gardening-or-sale-seeds-and/">https://www.politifact.com/factchecks/2020/apr/15/facebook-posts/covid-order-doesnt-ban-gardening-or-sale-seeds-and/</a>
67	The CDC ci	CDC lying	"The CDC	<a href="https://www.leadstories.com/hoax-alert/2020/04/fact-check-cdc-did-not-confess-to-lying-about-covid-19-death-numbers.html">https://www.leadstories.com/hoax-alert/2020/04/fact-check-cdc-did-not-confess-to-lying-about-covid-19-death-numbers.html</a>
68	A top Ger	man w	"A top Ger	<a href="https://www.leadstories.com/hoax-alert/2020/04/fact-check-a-top-german-doctor-recommends-whiskey-to-protect-against-covid-19.html">https://www.leadstories.com/hoax-alert/2020/04/fact-check-a-top-german-doctor-recommends-whiskey-to-protect-against-covid-19.html</a>
69	Bill Gates i	Bill Gates	vaccine tra	<a href="https://www.factcheck.org/2020/04/conspiracy-theory-misinterprets-goals-of-gates-foundation/">https://www.factcheck.org/2020/04/conspiracy-theory-misinterprets-goals-of-gates-foundation/</a>
72	COVID-19	COVID-19	"Did a COV	<a href="https://www.factcheck.org/2020/04/false-claim-of-deadly-coronavirus-vaccine-trial-in-africa/">https://www.factcheck.org/2020/04/false-claim-of-deadly-coronavirus-vaccine-trial-in-africa/</a>
80	The numb	Birx	(covid-19 OR cov	<a href="https://www.leadstories.com/hoax-alert/2020/04/Fact-Check-US-Coronavirus-Death-Tolls-NOT-Being-Inflated.html">https://www.leadstories.com/hoax-alert/2020/04/Fact-Check-US-Coronavirus-Death-Tolls-NOT-Being-Inflated.html</a>
84	SG exposi	5G hemogl	"Can 5G e	<a href="https://www.healthfeedback.com/claimreview/conspiracy-theorists-claim-that-5g-increases-vulnerability-to-covid-19-with-baseless-theory-that-it-affects-hemoglobin/">https://www.healthfeedback.com/claimreview/conspiracy-theorists-claim-that-5g-increases-vulnerability-to-covid-19-with-baseless-theory-that-it-affects-hemoglobin/</a>
99	Pat Robert	Pat Robertson	coron	<a href="https://www.politifact.com/factchecks/2020/apr/13/tweets/no-pat-robertson-didnt-say-covid-19-caused-oral-se/">https://www.politifact.com/factchecks/2020/apr/13/tweets/no-pat-robertson-didnt-say-covid-19-caused-oral-se/</a>
109	COVID-19	covid-19	death exagg	<a href="https://www.factcheck.org/2020/04/social-media-posts-make-baseless-claim-on-covid-19-death-toll/">https://www.factcheck.org/2020/04/social-media-posts-make-baseless-claim-on-covid-19-death-toll/</a>
110	President	Trump	hydroxychloroquine	<a href="https://www.washingtonpost.com/politics/2020/04/13/how-false-hope-spread-about-hydroxychloroquine-its-consequences/">https://www.washingtonpost.com/politics/2020/04/13/how-false-hope-spread-about-hydroxychloroquine-its-consequences/</a>
115	New data	21 million	"21 million	<a href="https://www.leadstories.com/hoax-alert/2020/04/fact-check-new-data-does-not-reveal-that-21-million-chinese-died-of-coronavirus.html">https://www.leadstories.com/hoax-alert/2020/04/fact-check-new-data-does-not-reveal-that-21-million-chinese-died-of-coronavirus.html</a>
125	The Chine	Chinese Communist P	https://tw	<a href="https://www.politifact.com/factchecks/2020/apr/07/charlie-kirk/china-spying-you-through-zoom-charlie-kirk-oversta/">https://www.politifact.com/factchecks/2020/apr/07/charlie-kirk/china-spying-you-through-zoom-charlie-kirk-oversta/</a>
159	Trump saic	Trump tes	"Remarks	<a href="https://www.politifact.com/factchecks/2020/apr/06/donald-trump/donald-trump-wrong-about-covid-19-testing-airplane/">https://www.politifact.com/factchecks/2020/apr/06/donald-trump/donald-trump-wrong-about-covid-19-testing-airplane/</a>
183	Microwavi	microwavi	"Does Mic	<a href="https://www.leadstories.com/hoax-alert/2020/04/Fact-Check-Microwaving-Masks-In-Plastic-Bags-Is-NOT-Safe-Way-To-Sterilize-Them.html">https://www.leadstories.com/hoax-alert/2020/04/Fact-Check-Microwaving-Masks-In-Plastic-Bags-Is-NOT-Safe-Way-To-Sterilize-Them.html</a>
188	Queen Eliz	Queen Eliz	"Queen	<a href="https://www.leadstories.com/hoax-alert/2020/03/Fact-Check-Coronavirus-'Patient-Zero'-Is-NOT-A-Man-Who-Had-Sex-With-A-Bat.html">https://www.leadstories.com/hoax-alert/2020/03/Fact-Check-Coronavirus-'Patient-Zero'-Is-NOT-A-Man-Who-Had-Sex-With-A-Bat.html</a>
200	COVID-19:	patient zer	"COVID-15	<a href="https://www.healthfeedback.com/claimreview/claim-by-nobel-laureate-luc-montagnier-that-the-novel-coronavirus-is-man-made-and-contains-genetic-material-from-hiv-is-inaccurate/">https://www.healthfeedback.com/claimreview/claim-by-nobel-laureate-luc-montagnier-that-the-novel-coronavirus-is-man-made-and-contains-genetic-material-from-hiv-is-inaccurate/</a>
201	HereaE	"s no specifi	"Manufact	<a href="https://www.healthfeedback.com/claimreview/claim-by-nobel-laureate-luc-montagnier-that-the-novel-coronavirus-is-man-made-and-contains-genetic-material-from-hiv-is-inaccurate/">https://www.healthfeedback.com/claimreview/claim-by-nobel-laureate-luc-montagnier-that-the-novel-coronavirus-is-man-made-and-contains-genetic-material-from-hiv-is-inaccurate/</a>
215	Biden calle	Biden Trun	"EXCLUSIV	<a href="https://www.politifact.com/factchecks/2020/mar/27/donald-trump/fact-checking-whether-biden-called-trump-xenophobic/">https://www.politifact.com/factchecks/2020/mar/27/donald-trump/fact-checking-whether-biden-called-trump-xenophobic/</a>
218	Ibuprofen	Ibuprofen	worsen coronavi	<a href="https://www.factcheck.org/2020/03/no-evidence-to-back-covid-19-ibuprofen-concerns/">https://www.factcheck.org/2020/03/no-evidence-to-back-covid-19-ibuprofen-concerns/</a>
234	NY Gov.	O Cuomo	rel	<a href="https://www.politifact.com/factchecks/2020/mar/26/facebook-posts/theres-no-evidence-covid-19-can-survive-surfaces-1/">https://www.politifact.com/factchecks/2020/mar/26/facebook-posts/theres-no-evidence-covid-19-can-survive-surfaces-1/</a>
253	The CDC n	COVID-19	"CDC says	<a href="https://www.politifact.com/factchecks/2020/mar/26/facebook-posts/theres-no-evidence-covid-19-can-survive-surfaces-1/">https://www.politifact.com/factchecks/2020/mar/26/facebook-posts/theres-no-evidence-covid-19-can-survive-surfaces-1/</a>
297	Belgium H	Belgium se	"Belgium F	<a href="https://www.leadstories.com/hoax-alert/2020/03/fact-check-belgium-health-minister-did-NOT-put-ban-on-non-essential-sexual-activities-of-persons-3-or-great.html">https://www.leadstories.com/hoax-alert/2020/03/fact-check-belgium-health-minister-did-NOT-put-ban-on-non-essential-sexual-activities-of-persons-3-or-great.html</a>

```
[ ] # Assuming Cluster 0 is in favor, Cluster 1 is against, and Cluster 2 is neutral
    stance_mapping = {
        0: 'in_favor',
        1: 'against',
        2: 'neutral'
    }

    df['predicted_stance'] = df['cluster'].map(stance_mapping)

[ ] print(df[['text', 'predicted_stance']])
```

```
[ ] import numpy as np

def kmeans(X, k, max_iters=100):
    centroids = X[np.random.choice(X.shape[0], k, replace=False)]

    for _ in range(max_iters):

        distances = np.linalg.norm(X[:, np.newaxis] - centroids, axis=2)
        labels = np.argmin(distances, axis=1)

        new_centroids = np.array([X[labels == i].mean(axis=0) for i in range(k)])

        if np.all(centroids == new_centroids):
            break

        centroids = new_centroids

    return centroids, labels

np.random.seed(42)

data = np.random.rand(100, 2)

num_clusters = 3

centroids, labels = kmeans(data, num_clusters)
```



```

padded_sequences = pad_sequences(sequences)

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(padded_sequences, stances_encoded_categorical, test_size=0.2, random_state=42)

# Build the deep learning model
model = Sequential()
model.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=100, input_length=len(padded_sequences[0])))
model.add(Bidirectional(LSTM(100)))
model.add(Dense(3, activation='softmax')) # Assuming 3 classes for stance detection

# Compile the model
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])

# Train the model
model.fit(X_train, y_train, epochs=5, validation_data=(X_test, y_test))

# Evaluate the model
y_pred = model.predict(X_test)
y_pred_classes = np.argmax(y_pred, axis=1)
y_test_classes = np.argmax(y_test, axis=1)

print("\nClassification Report:\n", classification_report(y_test_classes, y_pred_classes))
print("\nConfusion Matrix:\n", confusion_matrix(y_test_classes, y_pred_classes))

```

```

[ ] from sklearn.decomposition import PCA
    import matplotlib.pyplot as plt
    import seaborn as sns

[ ] X_clustered = X_tfidf.toarray()

[ ] pca = PCA(n_components=2)
    X_pca = pca.fit_transform(X_clustered)

[ ] import matplotlib.pyplot as plt
    import seaborn as sns

[ ] cluster_counts = df['cluster'].value_counts().sort_index()

    stance_counts = df.groupby(['cluster', 'predicted_stance']).size().unstack(fill_value=0)

[ ] plt.figure(figsize=(12, 8))
    sns.set(style="whitegrid")
    stance_counts.plot(kind='bar', stacked=True, colormap='viridis', ax=plt.gca())
    plt.title('Distribution of Predicted Stances in Each Cluster')
    plt.xlabel('Cluster')
    plt.ylabel('Number of Samples')
    plt.legend(title='Predicted Stance', bbox_to_anchor=(1.05, 1), loc='upper left')

```

### **3.5) Key Challenges-**

#### **Confined categorized Datasets:**

Availability of labeled datasets for Code-combined Hindi-English stance detection may be constrained. Deep learning fashions generally require large quantities of annotated records for powerful training.

#### **Code-mixing Complexity:**

Social media customers frequently blend Hindi and English within the identical sentence or submit, main to code-blending demanding situations. The version wishes to recognize and correctly process this combined-language content material.

#### **Embedding challenges:**

Pre-skilled phrase embeddings won't capture the nuances of Code-blended language efficaciously. creating appropriate embeddings that represent the combined language context is critical.

#### **Ambiguity and Sarcasm:**

Social media content material regularly includes ambiguity, sarcasm, and nuanced expressions. information the sentiment and stance in such cases can be hard for a deep gaining knowledge of model.

#### **Multilingual Lexicons:**

Hindi and English have one of a kind linguistic systems and lexicons. Incorporating a suitable multilingual lexicon that understands both languages is vital for accurate stance detection.

#### **Class Imbalance:**

Stance detection obligations might have imbalances in magnificence distribution. a few stances or sentiments can be more standard than others, main to biased fashions if no longer addressed well.

#### **Contextual understanding:**

Deep studying models want to capture the contextual statistics in Code-combined textual content. know-how the context is essential for accurate stance detection, because the meaning of words can

exchange primarily based on the surrounding text.

### **Model Generalization:**

ensuring that the model generalizes nicely to diverse Code-blended facts past the training set is a undertaking. Overfitting or underfitting may additionally arise if the model isn't strong sufficient.

### **Useful resource Constraints:**

schooling deep getting to know fashions may be useful resource-extensive, mainly for large datasets and complicated architectures. confined computational sources might also avert the improvement and education of powerful models.

### **Ethical considerations:**

Stance detection in social media records may additionally involve touchy topics or biased content material. making sure moral concerns in schooling data and keeping off biases in model predictions is important.

## 4) Testing

### 4.1) Testing strategy

Detecting stance in Code-mixed Hindi-English social media facts using deep learning knowledge includes several steps. Stance detection refers to determining the mindset or angle expressed toward a particular target or topic in text. Here is a tried approach that will help you evaluate the performance of your deep learning version:

#### 1. Statistics Preprocessing:

Preprocess the code-mixed Hindi-English social media data, inclusive of text cleansing, tokenization, and handling special characters.

Convert the textual content into numerical representations appropriate for deep learning models.

#### 2. Statistics Splitting:

Divide your dataset into training, validation, and test sets. Not unusual splits are 70-15-15 or 80-10-10.

#### 3. Embedding Layer:

Use pre-trained word embeddings (Word2Vec, GloVe, or FastText) for each Hindi and English word. Implement an embedding layer which can deal with code-mixed textual content.

#### 4. Version architecture:

Design a deep learning architecture for stance detection. A recurrent neural network (RNN) or transformer-based architecture can be effective.

Ensure that the version can cope with code-mixed text efficiently.

#### 5. Training:

Teach your deep getting to know version on the training set.

make use of the validation set to song hyperparameters and save you overfitting.

## **6. Evaluation Metrics:**

pick out suitable assessment metrics inclusive of accuracy, precision, remember, and F1 rating for stance detection. consider using confusion matrices to understand the version's performance throughout different lessons.

## **7. Code-blending challenges:**

compare how well your version handles code-blending demanding situations, including varying language ratios, inconsistent language switches, and the presence of casual language.

## **8. Move-Validation:**

perform okay-fold go-validation to ensure robustness and generalize the model's overall performance.

## **9. Handling Imbalanced data:**

take a look at for class imbalances and rent techniques like oversampling, undersampling, or the use of class weights to address imbalanced information.

## **10. First-rate-Tuning:**

remember exceptional-tuning your model primarily based on unique challenges discovered at some point of trying out.

## **11. Interpretability:**

include methods for interpreting the model's predictions, which includes attention mechanisms, to apprehend which elements of the input make contributions to the output.

## **12. Deployment concerns:**

If applicable, do not forget deployment considerations which includes version length, latency, and resource necessities.

### **13. Post-Deployment tracking:**

installation monitoring mechanisms to song the version's overall performance in a real-global placing and deal with any drift or degradation.

### **14. Documentation:**

report your testing process, such as model structure, hyperparameters, and any unique issues for code-blended textual content.

### **15. Iterative development:**

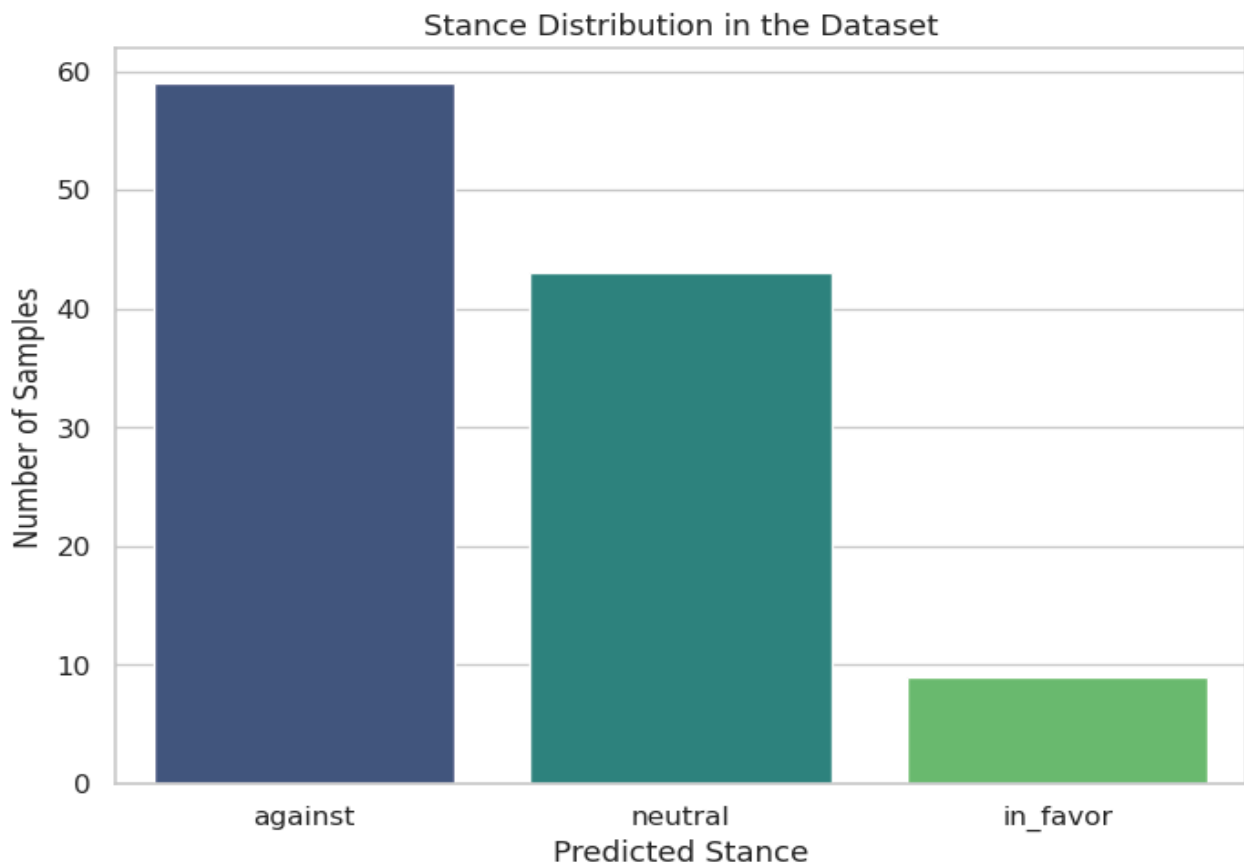
Iterate at the version based on testing outcomes and feedback. non-stop improvement is critical for adapting to evolving social media language traits.

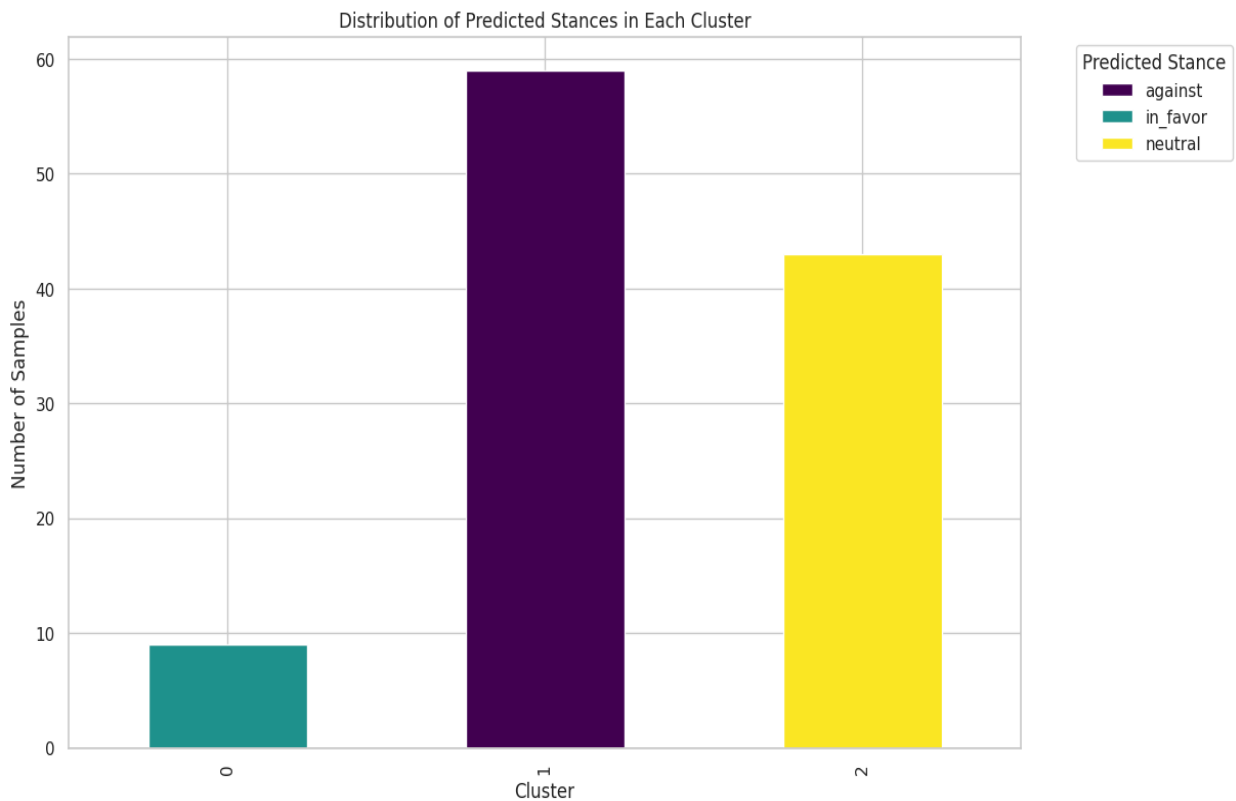
## 5.)Results and Evaluation

### 5.1) Results (presentation of findings, interpretation of the results, etc.)

Model	Accuracy(%)
RBF Kernel SVM*	58.7
Random Forest*	54.7
Linear SVM*	56.6
CNN	<b>61.4</b>
CNN + MTL	<b>63.2</b>

Stance	Sentiment of Tweet (%)		
	Positive	Negative	Other
Favor	22.36	62.16	15.48
Against	5.45	87.97	6.59
Neither	12.96	50.47	36.57





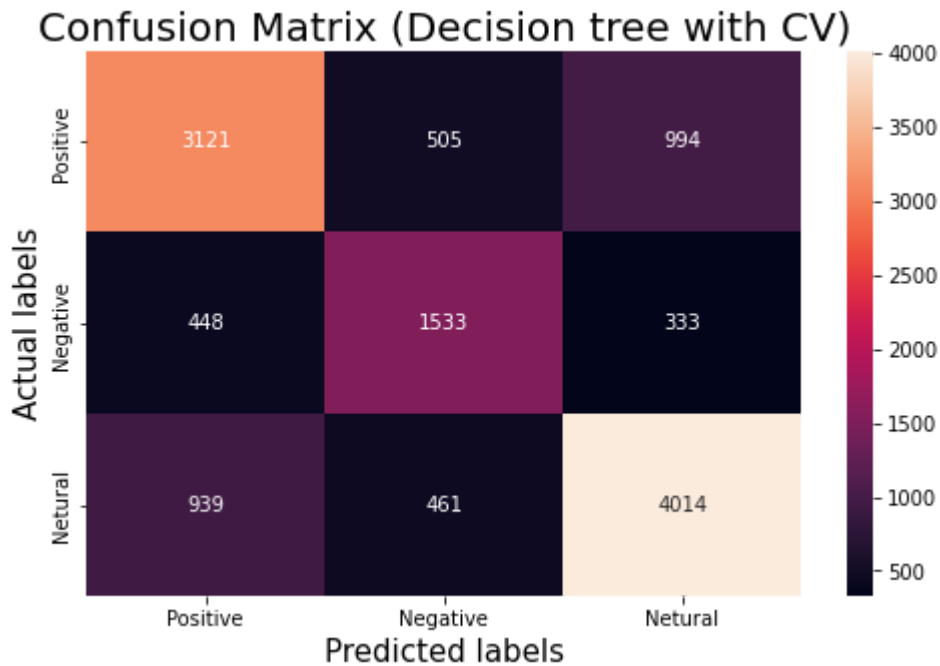
### Evaluation Metrics:

Compute trendy assessment metrics including accuracy, precision, keep in mind, and F1 rating. these metrics will offer a comprehensive view of your version's performance.

### 2. Confusion Matrix:

analyze the confusion matrix to recognize how well your model is appearing throughout distinctive stance lessons. discover areas where the version might be suffering.





**3. Magnificence Imbalances:**

If there are imbalances in the dataset, consider comparing the version the usage of metrics like precision-recollect curves or vicinity below the precision-take into account curve (AUC-PR) for a more nuanced evaluation.

**4. Code-blending demanding situations:**

look at how well your version handles code-blending demanding situations. examine its overall performance on instances with various language ratios and casual language utilization.

**5. Cross-Validation outcomes:**

in case you executed okay-fold move-validation, analyze the effects from every fold to make sure consistency in performance.

**6. Interpretability:**

in case you integrated interpretability features like attention mechanisms, analyze them to advantage insights into how the model makes predictions. this could help pick out styles and capability regions for improvement.

## **7. Put up-Deployment tracking:**

in case your version is deployed, screen its performance in a actual-world putting. take a look at for any go with the flow or degradation in performance through the years and make important modifications.

## **8. Qualitative assessment:**

Manually evaluate misclassified times to apprehend the character of errors. This qualitative evaluation can provide precious insights into regions in which the version struggles.

## **9. Person remarks:**

If feasible, collect remarks from cease-users or area specialists. Their insights may be treasured for refining the model and addressing real-world issues.

## **10. Comparative evaluation:**

examine your model's performance with baseline models or present approaches to assess its effectiveness.

## **11. Iterative development:**

based totally at the evaluation consequences, iteratively enhance your model by using pleasant-tuning, adjusting hyperparameters, or incorporating extra functions.

## **12. Documentation and Reporting:**

Record the assessment results comprehensively. provide a clean report that consists of key metrics, visualizations, and insights won from the assessment process.

No.	Tweet	Label	NS Prediction	DAN Prediction
1	@brad_dickson My son teaches in Japan. They wear masks because they are a polite society. School closed asap in Feb. Did remote learning. But as of early June, back to in school learning due to so few cases. Masks work.	FAVOR	FAVOR	FAVOR
2	Hell no to your mask mandate	AGAINST	AGAINST	AGAINST
3	If, 6 months later, you're still wearing a mask.....you might as well wear one the rest of your life.	AGAINST	FAVOR	FAVOR
4	People tweeting from their smart phones about how masks are a form of government control is hilarious to me.	FAVOR	AGAINST	AGAINST
5	Thank goodness Trump wasn't there to greet the astronauts after splashdown. I'm sure he would have shown up with no mask! #SplashDown #SpaceX	FAVOR	FAVOR	AGAINST
6	Some of ya'll couldn't dissect a frog in high school but you know more than health professionals about the Coronavirus!?!? :man_facepalming_dark_skin_tone: #COVID19	FAVOR	NONE	FAVOR
7	Small local grocery store did not have sign requiring mask per state mandate and was pretty busy. Only about half of customers wearing masks. The dairy section looked almost empty. Love it and they will continue to get my business.	AGAINST	FAVOR	AGAINST
8	@simondolan @SaltySeaDog7 By not wearing a mask you are giving the children of the COVID generation a chance to go to school, play sports, and have real childhoods	AGAINST	FAVOR	AGAINST

	Precision	Recall	F-measure
<b>Class Positive</b>	0.641	0.481	0.55
<b>Class Neutral</b>	0.716	0.857	0.78
<b>Class Negative</b>	0.588	0.455	0.513

Table 6: Precision, Recall and F-measure.

## **6.) Conclusion and Future Scope-**

### **6.1 )Conclusion (summarize key findings, limitations and contributions to the field).**

We've built a COVID-19- Stance dataset that can be used to in addition the research on stance detection, particularly inside the context of COVID-19 pandemic. similarly to the dataset, we have established baselines using numerous supervised models used in previous works on stance detection, and additionally two fashions which can employ unlabeled data and information from a previous stance detection venture, respectively. As a part of future paintings, we plan to look at the benefits of the opinion and sentiment data that we annotated closer to the stance detection.

#### **Key findings**

1. The tweet explicitly expresses opinion/sentiment approximately the goal.
2. The tweet expresses opinion/sentiment approximately some thing/a person aside from the target.
3. The tweet isn't always expressing opinion/sentiment about whatever.

#### **Limitations and Challenges-**

##### **Annotation challenges:**

Stance detection often calls for classified statistics for education system studying models. Annotating code-combined statistics for stance may be subjective, and disagreements amongst annotators may also get up because of the inherent complexity of mixed-language content.

##### **Code-mixing Variability:**

Code-blending styles can vary widely across distinctive social media platforms, consumer demographics, and even person users. Adapting a version to handle this variability may be difficult.

### **Contextual Ambiguity:**

Code-mixed content frequently is based heavily on contextual cues. the ambiguity in context, specially in brief social media posts, can make it hard for models to accurately determine the stance of a declaration.

### **Aid Constraints:**

Deep studying models, in particular complicated ones, may require massive computational assets. education deep learning fashions for code-mixed information is probably aid-intensive, proscribing accessibility for researchers with confined computational talents.

### **Generalization troubles:**

models educated on a selected dataset may conflict to generalize to new or unseen code-blended statistics. Adapting fashions to specific social media platforms or evolving language trends can be necessary.

### **Moral concerns:**

Studying social media facts, in particular in the context of sentiment or stance, increases ethical concerns. ensuring privateness and responsible use of user-generated content material is crucial.

### **Evaluation Metrics:**

choosing appropriate assessment metrics for stance detection in code-blended information can be difficult. not unusual metrics might not fully seize the nuances of language mixing and stance expression in this context.

## 5.2 Future Scope

### **Multilingual Stance Detection:**

increasing the scope to locate stance in a broader variety of languages beyond Hindi and English, especially in code-combined contexts, can increase the applicability of the fashions. this will involve incorporating extra language pairs or even multilingual fashions.

### **Nice-Grained Stance evaluation:**

Going past binary stance detection, destiny research may explore best-grained or multi-class stance analysis. this could involve categorizing stances into extra nuanced classes or figuring out specific factors of stance, imparting a deeper information of person opinions.

### **Context-conscious fashions:**

growing fashions which are greater touchy to contextual cues and might higher understand the subtleties of code-mixing in specific contexts. Context-conscious fashions can also improve accuracy in ambiguous or evolving language conditions.

### **Ethical issues and Bias Mitigation:**

research have to preserve to cope with moral considerations associated with reading social media content. Efforts to mitigate bias, ensure person privateness, and sell responsible use of AI in social media evaluation are crucial.

### **Consumer-Centric procedures:**

Considering the person's attitude in stance detection fashions can be important. research might also explore user-centric features, possibilities, and feedback mechanisms to increase the interpretability and person-friendliness of stance detection systems.

## References

- 1.) Agarwal, A., & Choudhury, M. (2018). "Stance class in code-mixed tweets: Unsupervised bootstrapping vs. semi-supervised learning." In complaints of the 11th worldwide convention on Language sources and evaluation (LREC).
- 2.) Choudhury, M., Gambäck, B., & Palmer, M. (2007). "in the direction of a wellknown named entity recognizer for Indian languages." In proceedings of the IJCNLP-08.
- 3.) Gupta, A., & Choudhury, M. (2018). "Code-blending: A assignment for Language expertise within the Social Media generation." ACM Transactions on Asian and low-aid Language statistics Processing (TALLIP), 17(three), 14.
- 4.) Liu, B. (2012). "Sentiment analysis and Opinion Mining." Synthesis Lectures on Human Language technologies.
- 5.) Pennington, J., Socher, R., & Manning, C. (2014). "GloVe: international Vectors for word representation." complaints of the 2014 convention on Empirical strategies in natural Language Processing (EMNLP).
- 6.) Ruder, S., & Plank, B. (2018). "Linguistic style and social network analysis for authorship attribution of tweets." In court cases of the 2018 convention of the North American chapter of the association for Computational Linguistics: Human Language technology.
- 7.) Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C., Ng, A., & Potts, C. (2013). "Recursive deep fashions for semantic compositionality over a sentiment treebank." In proceedings of the 2013 conference on Empirical methods in natural Language Processing (EMNLP).
- 8.) Liu, B. (2012). "Sentiment analysis and Opinion Mining." Synthesis Lectures on Human Language technologies

# PLAGIARISM REPORT

201385@juitsolan.in

## ORIGINALITY REPORT

**12** %

SIMILARITY INDEX

**4** %

INTERNET SOURCES

**5** %

PUBLICATIONS

**3** %

STUDENT PAPERS

## PRIMARY SOURCES

<b>1</b>	<b>ir.juit.ac.in:8080</b> Internet Source	<b>4</b> %
<b>2</b>	<b>Submitted to The Robert Gordon University</b> Student Paper	<b>1</b> %
<b>3</b>	<b>aclanthology.org</b> Internet Source	<b>1</b> %
<b>4</b>	<b>indjst.org</b> Internet Source	<b>1</b> %
<b>5</b>	<b>Submitted to Baker College</b> Student Paper	<b>1</b> %
<b>6</b>	<b>Submitted to Glasgow Caledonian University</b> Student Paper	<b>1</b> %
<b>7</b>	<b>www.aclweb.org</b> Internet Source	<b>1</b> %
<b>8</b>	<b>Submitted to University of Hertfordshire</b> Student Paper	<b>&lt;1</b> %
<b>9</b>	<b>"Green Internet of Things and Machine Learning", Wiley, 2021</b> Publication	<b>&lt;1</b> %