A Project Report on

# Leveraging Transfer Learning to Analyse Software Developers Opinions for Enhanced Project Insights

*submitted in partial fulfillment of the requirement for the award of the Degree of*

### BACHELOR OF TECHNOLOGY

in

### G. PULLAIAH COLLEGE OF ENGINEERING AND TECHNOLOGY

**(Autonomous)**

(Approved by AICTE | NAAC Accreditation with 'A' Grade | Accredited by NBA (ECE & EEE) |Permanently Affiliated to JNTUA)

**By**

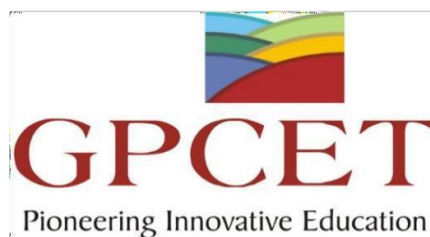**KASI HARISH (22AT5A0505)**

**BOLLE RAJKISHORE (21AT1A0516)**

**BELAGAM SYED MOHAMMAD ASHWAK (21AT1A0514)**

**CHITTAM DURGA CHARAN REDDY (21AT1A0524)**

**Under the Guidance of**

## Dr S Prem Kumar MTech., Ph.D.

**Professor**



### DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

### G. PULLAIAH COLLEGE OF ENGINEERING AND TECHNOLOGY

**(Autonomous)**

(Approved by AICTE | NAAC Accreditation with 'A' Grade | Accredited by NBA (ECE & EEE) Permanently Affiliated to JNTUA)

**2021-2025**

I

# G. PULLAIAH COLLEGE OF ENGINEERING AND TECHNOLOGY

## (Autonomous)

(Approved by AICTE | NAAC Accreditation with 'A' Grade | Accredited by NBA (ECE & EEE) |

Permanently Affiliated to JNTUA)

## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



## CERTIFICATE

This is to certify that the project report entitled **"LEVERAGING TRANSFER LEARNING TO ANALYSE SOFTWARE DEVELOPERS OPINIONS FOR ENHANCED PROJECT INSIGHTS"** being submitted by **KASI HARISH (22AT5A0505),BOLLE RAJ KISHORE (21AT1A0516),B.S.MOHAMMAD ASHWAK (21AT1A0514), C. DURGA CHARAN REDDY (21AT1A0524)** in partial fulfillment of the requirement for the award of the degree of Bachelor of Technology in Computer Science and Engineering of G. Pullaiah College of Engineering and Technology, Kurnool is a record of Bonafide work carried out by them under my guidance and supervision. The results embodied in this project report have not been submitted to any other university or institute for the award of any Degree or Diploma.

Dr S Prem Kumar MTech., Ph.D.

Mrs.K. LAKSHMI MTech., (Ph.D.)

Project Supervisor

Head of the Department

Date of Viva-Voce _____

INTERNAL EXAMINER

EXTERNAL EXAMINER

# ACKNOWLEDGEMENT

We thank our project supervisor**, Dr.S.Prem Kumar** MTech., Ph.D. for his guidance, valuable suggestions and support in the completion of the project.

We would like to express our deep sense of gratitude and our sincere thanks to HOD **Mrs.K.LAKSHMI** MTech.,(Ph.D.) . Department of Computer Science and Engineering, G. Pullaiah College of Engineering and Technology, Kurnool for providing the necessary facilities and encouragement towards the project work.

We owe indebtedness to our principal **Dr. C. SREENIVASA RAO** MTech., Ph.D. G. Pullaiah College of Engineering and Technology, Kurnool for providing us the required facilities.

We are extremely grateful to Chairman **Mr. G. V. M. MOHAN KUMAR** of G. Pullaiah College of Engineering and Technology, Kurnool, Andhra Pradesh for their good blessings.

We gratefully acknowledge and express our thanks to teaching and non teaching staff of CSE Department.

**Project Associates**

**KASI HARISH (22AT5A0505)**

**BOLLE RAJKISHORE (21AT1A0516)**

**BELAGAM SYED MOHAMMAD ASHWAK (21AT1A0514)**

**CHITTAM DURGA CHARAN REDDY (21AT1A0524)**

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| Abbreviations | Definitions |
|---|---|
| RNN | Recurrent Neural Networks |
| DT | Decision Tree |
| NLP | Natural Language processing |
| ML | Machine Learning |
| AI | Artificial intelligence |
| LR | Logistic Regression |
| RF | Random Forest |
| SVM | Support Vector Machine |

# ABSTRACT

This project explores the application of multiple advanced machine learning and deep learning algorithms, including Random Forest, Multi-Layer Perceptron (MLP), Convolutional Neural Networks (CNN-1D), Bidirectional Encoder Representations from Transformers (BERT), and Generative Pretrained Transformers (GPT), to improve sentiment analysis of social media content. By leveraging the strengths of these diverse models, the research aims to enhance the accuracy and depth of sentiment classification in the context of software development projects. Random Forest is utilized for its ability to handle large datasets and uncover complex interactions within the text data. On the other hand, the neural network-based models like MLP, CNN-1D, and BERT capture more intricate features and contextual nuances of the textual content. GPT, a cutting-edge language model, contributes to refining the sentiment analysis by generating richer representations of language. The integration of these methods forms a powerful ensemble system, combining the advantages of traditional ensemble learning with modern deep learning techniques. This hybrid approach allows for a more detailed and accurate understanding of developer opinions, providing valuable insights that can inform software development processes and project management strategies. Ultimately, this project has the potential of advanced AI techniques to improve the quality and effectiveness of sentiment analysis for decision-making in software development.

# CHAPTER 1

# INTRODUCTION

## 1.1 General Introduction:

## 1. Background and Context

In recent years, social media platforms have revolutionized communication and information dissemination. Platforms like Twitter, Facebook, Instagram, and Reddit have become integral to everyday life, offering users a space to share opinions, experiences, and information in real time. This rapid exchange of content has created vast repositories of unstructured data that reflect public sentiment on a wide array of topics, from political events to consumer products. The ability to analyse and interpret this sentiment has profound implications for businesses, governments, and researchers, making sentiment analysis a critical area of study.

## 2. The Significance of Sentiment Analysis

Sentiment analysis involves using computational methods to determine the emotional tone behind a series of words. This process is crucial for various applications, including market research, customer feedback analysis, brand management, and social trend monitoring. By understanding public sentiment, organizations can make data-driven decisions, tailor their strategies, and respond proactively to emerging trends. For instance, companies can gauge consumer reactions to a new product, while policymakers can assess public opinion on legislative changes.

## 3. Challenges in Sentiment Analysis

Despite its importance, sentiment analysis is fraught with challenges. Social media texts are often informal, containing slang, abbreviations, and emojis, which can complicate analysis. The context in which sentiments are expressed can vary widely, leading to ambiguity and difficulty in accurately interpreting the emotional content. Additionally, sentiment expressions are often nuanced and may not fit neatly into predefined categories of positive, negative, or neutral. This variability necessitates sophisticated methods to achieve accurate results.

## 4. Overview of Sentiment Analysis

- **Definition and Importance**: Begin by defining sentiment analysis and explaining its significance in the broader context of data science, especially in text mining.

- Discuss how sentiment analysis helps organizations understand public opinion, consumer behaviour, and social media trends.

- **Sentiment Analysis in social media**: Highlight the rapid rise of social media as a rich source of user-generated content, and how it serves as an invaluable tool for opinion mining, particularly in assessing public sentiment towards products, services, and, in your case, software development projects.

- **Challenges in Sentiment Analysis**: Discuss the challenges such as ambiguity in natural language, context-dependent sentiment, and the noisy nature of social media data (e.g., slang, abbreviations, mixed emotions).

## Sentiment Analysis in Software Development

- **Importance of Developer Feedback**: Emphasize the importance of capturing the opinions and sentiments of software developers to gain insights into software projects. This could relate to understanding developer experiences, opinions on tools or frameworks, collaboration dynamics, and project-specific issues.

- **Impact on Project Management**: Explain how sentiment analysis can impact project management decisions, from improving team dynamics to identifying early signs of project distress or success.

- **Role of Sentiment Analysis in Improving Project Outcomes**: Discuss how effective sentiment analysis of developer opinions can guide project managers in making informed decisions, enhancing communication, and streamlining software development processes.

## Overview of Machine Learning and Deep Learning Techniques

- **Machine Learning in Sentiment Analysis**: Provide a brief explanation of machine learning (ML), particularly supervised learning, and how models like Random Forest and MLP are widely used for text classification tasks such as sentiment analysis.

- **Deep Learning for Sentiment Analysis**: Transition into the realm of deep learning, focusing on how neural networks (e.g., CNN-1D, BERT, and GPT) can capture more complex patterns, context, and nuances in textual data.

## 5. Previous Research and Developments

Prior research has explored various methods for sentiment analysis, including rule-based approaches and simpler machine learning models. Traditional techniques such as Support Vector Machines (SVM) and Naive Bayes classifiers have been commonly used, but they often struggle with the nuances of social media language. Recent advancements have seen the integration of deep learning approaches, such as Recurrent Neural Networks (RNNs) and Transformer models, which capture context and semantics more effectively. However, these models can be computationally intensive and require significant resources.

## 6. Importance of the Study

This study's importance lies in its potential to improve sentiment analysis accuracy and applicability. By refining the methods used to classify sentiment, the research contributes to the broader field of natural language processing and offers practical benefits for businesses, policymakers, and researchers. Enhanced sentiment analysis tools can lead to better decision-making, more targeted marketing strategies, and a deeper understanding of social dynamics.

## 1.2 Objectives :

The main objective of our project is,

- Create and implement a sentiment analysis model using Random Forest, MLP, BERT , CNN-1D and GPT approach to classify social media text into positive, negative, or neutral categories. Assess the accuracy, precision, recall, and F1-score of the ML and DL models in predicting sentiment compared to baseline methods and traditional techniques.

- Analyze and address the challenges associated with the informal, nuanced, and context-dependent nature of social media language, including slang, abbreviations, and emojis.

- Compare the effectiveness of all model and traditional sentiment analysis techniques to identify strengths and limitations of each approach.

# CHAPTER 2
# SYSTEM PROPOSAL

## 2.1 EXISTING SYSTEM :

In existing system, Sentiment analysis of social media has evolved through various approaches, each with its strengths and limitations. Traditional sentiment analysis systems often rely on rule-based methods or simpler machine learning algorithms. Rule-based systems use predefined lists of words and phrases associated with positive or negative sentiments. These systems, while straightforward and easy to implement, face significant limitations. They struggle with the informal and diverse nature of social media language, such as slang, abbreviations, and emojis, which can lead to inaccurate sentiment classification. Additionally, rule-based approaches lack the ability to understand the context in which words are used, making them less effective in handling nuanced sentiments. On the other hand, machine learning-based systems such as Support Vector Machines (SVM) and Naive Bayes classifiers have been employed to improve accuracy. These methods learn from labelled training data to classify sentiment and can handle a broader range of textual features compared to rule-based systems. Despite their advantages, these algorithms have limitations as well. SVMs, for instance, can be computationally expensive and may not perform well with very large datasets or complex feature interactions. Naive Bayes, while efficient, assumes independence between features, which is often not the case in textual data, leading to suboptimal performance in capturing sentiment nuances.

## 2.1.1 DISADVANTAGES :

- Struggle with informal language, slang, abbreviations, and emojis commonly used in social media.
- Lack the ability to understand the context or nuances of sentiment expressions, leading to potential misclassifications.
- Require extensive manual effort to create and maintain sentiment lexicons and rules.
- Difficulty in scaling to handle large volumes of diverse text data efficiently.
- Can be slow and resource-heavy, especially with large datasets or high-dimensional feature spaces.
- The decision boundary created by SVMs can be complex, making it hard to interpret how classifications are made.

- May struggle with very large datasets, leading to longer training times and decreased efficiency.
- Assumes that features are independent of each other, which is often not true in textual data, potentially leading to reduced accuracy.

## 2.1 PROPOSED SYSTEM :

In proposed system, the input data as sentiment dataset is taken from dataset repository. In pre-processing, we can check missing values for avoid wrong prediction and label encoding for convert the strings into numeric integer value. Then, we can implement the NLP techniques for cleaning the text such as stop words, stem words, and remove punctuations, tokenization and padding. After that, we can implement the different vectorization techniques such as count vectorization. Then, we can split the cleaned text into test data and train data. Test data is used for prediction and train data is used for evaluation. The splitted data is carried out to ML and DL algorithms such as BERT, MLP, RF, CNN-1D and GPT.Finally, the system can estimate some performance metrics such as accuracy, precision, recall, f1-score and error rate. The effectiveness of the proposed method was confirmed by comparing accuracy improvement. Here, we can predict the sentiment from user's input such as negative or positive or neutral.

## 2.2.1 ADVANTAGES :

- The ML model scan handle large datasets efficiently.
- The experimental result is high when compared with existing system.
- Time consumption is low.
- Lack of ability to be spatially invariant to the input data.
- Implementing text cleaning techniques like removing stop words, stemming, punctuation removal, tokenization, and padding enhances the quality of the text data, making it more suitable for analysis.
- Integration of Hybrid Models: By combining Decision Tree with Logistic Regression, the proposed system benefits from the strengths of both models. Decision Trees handle non-linear relationships and complex feature interactions well, while Logistic Regression adds robustness to classification tasks, leading to more accurate and reliable predictions.

### 2.2  LITERATURE SURVEY :

**1. Title: "Sentiment Analysis of Social Media Data Using Machine Learning Algorithms"**

**Year** : 2022

**Author(s) :** A. K. Sharma, B. Kumar

**Methodology** :

This paper explores sentiment analysis on social media platforms by employing various machine learning algorithms, including Random Forest, Support Vector Machines (SVM), and Naive Bayes. The authors utilized a combination of feature extraction techniques such as Term Frequency-Inverse Document Frequency (TF-IDF) and word embeddings to convert text into numerical format. Data preprocessing included handling missing values and text normalization. The models were trained on a labeled dataset and evaluated using accuracy, precision, recall, and F1-score metrics.

**Demerits** :

- The study showed that SVM performed well in precision but was computationally intensive.
- Naive Bayes had lower recall, and the feature extraction methods were not sufficiently adaptive to handle the evolving slang and abbreviations in social media text.

**2. Title: "Hybrid Approaches for Sentiment Analysis Using Deep Learning Techniques"**

**Year** : 2022

**Author(s) :** M. R. Patel, S. K. Gupta

**Methodology** :

This paper introduces a hybrid model combining Convolutional Neural Networks (CNNs) and Long Short-Term Memory networks (LSTMs) for sentiment analysis. The approach leverages CNNs for feature extraction from text and LSTMs for capturing long-term dependencies and contextual information. Data preprocessing involved tokenization, stemming, and padding. The model was evaluated on various benchmarks, and its performance was compared with traditional machine learning models.

**Demerits** :

- Additionally, the complexity of deep learning models made them difficult to interpret.

- While the hybrid model showed significant improvements in accuracy and F1-score, it required extensive computational resources and long training times.

## 3. Title: "Enhancing Sentiment Analysis with Transformer Models and Attention Mechanisms"

**Year** : 2022

**Author**(s) : J. S. Lee, Y. T. Zhang

**Methodology** :

The paper investigates the application of Transformer models, specifically BERT (Bidirectional Encoder Representations from Transformers), for sentiment analysis. The study incorporates attention mechanisms to better capture context and dependencies in the text. Preprocessing steps included stop-word removal, tokenization, and BERT-specific tokenization.

**Demerits** :

- Although the Transformer model achieved high accuracy, it was computationally expensive and required large amounts of memory.
- The model's complexity also posed challenges for real-time sentiment analysis.

## 4. Title: "A Comparative Study of Ensemble Methods for Sentiment Classification"

**Year** : 2022

**Author**(s) : R. N. Sharma, P. A. Singh

**Methodology** :

This study compares various ensemble methods for sentiment classification, including Random Forest, Gradient Boosting, and Voting Classifiers. The authors employed a range of preprocessing techniques such as lemmatization, stop-word removal, and feature selection. The models were evaluated on sentiment datasets using cross-validation techniques to assess their performance in terms of accuracy and F1-score.

**Demerits :**

- Ensemble methods, while improving accuracy, were found to be prone to overfitting .
- Additionally, managing multiple models increased computational overhead.

# CHAPTER 3

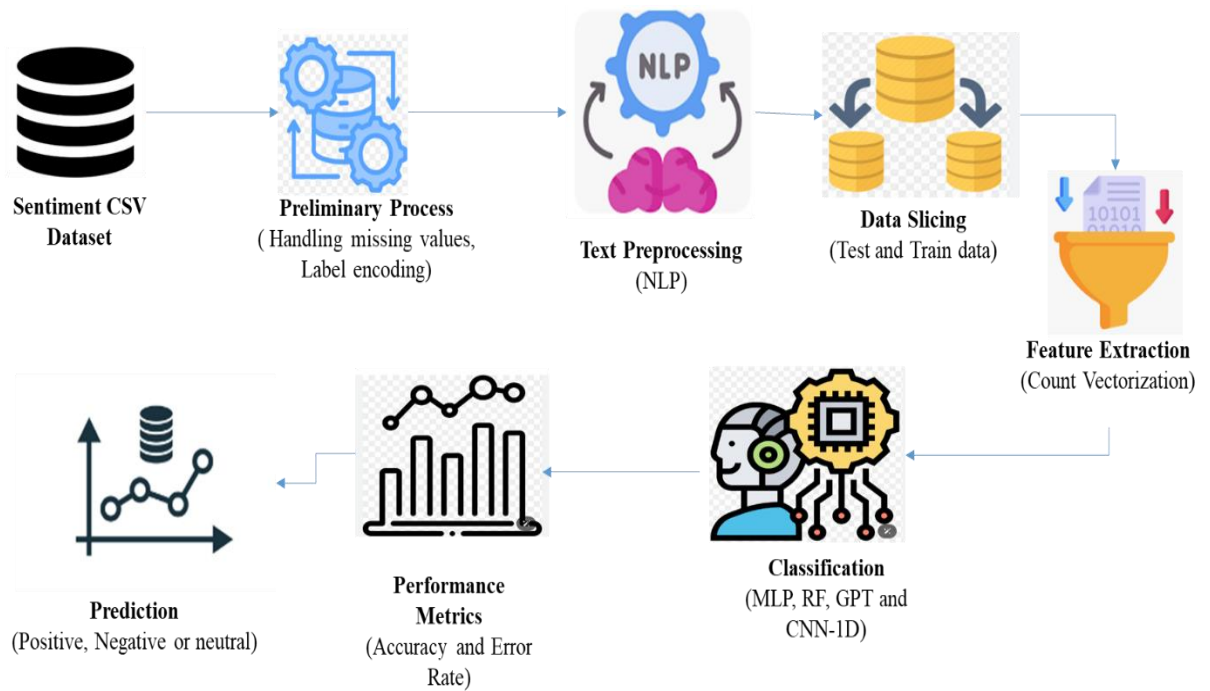# SYSTEM DIAGRAMS

## 3.1 SYSTEM ARCHITECTURE :



**FIGURE 3.1: SYSTEM ARCHITECTURE**

The architecture diagram outlines the workflow for processing and classifying Sentiment analysis dataset. Data Selection involves acquiring the dataset. Data Preprocessing handles missing values and label encoding. Text Preprocessing includes cleaning and standardizing text through stop words removal, stemming, and tokenization. The cleaned text is then converted into numerical format using Vectorization. The data is Split into training and test sets. Classification models, such as ML and DL models, are trained and evaluated. Result Generation computes performance metrics, and Prediction applies the trained models to classify new data, providing difficulty level insights.

**3.2 FLOW DIAGRAM** :



**FIGURE 3.2: FLOW DIAGRAM**

The flow diagram presents the sequential process for sentiment analysis dataset classification. It starts with Data Selection, where the dataset is sourced. Data Preprocessing follows, addressing missing values and encoding labels. Next, Text Preprocessing cleans and standardizes the text. The processed text undergoes Vectorization to convert it into numerical format. The data is then Split into training and test sets. In the Classification phase, models like ML and DL are trained and evaluated. Result Generation calculates performance metrics, and Prediction usesthe models to classify new inputs, providing insights into sentiment levels.

## 3.3 DATA FLOW DIAGRAM :

### 3.3.1 Level 0 :



In the Data Flow Diagram (DFD) Level 0, the **Data Selection** phase involves acquiring relevant datasets from repositories or sources, which are then input into the system. This is followed by **Data Preprocessing**, where the raw data undergoes cleaning and transformation processes to handle missing values, standardize formats, and prepare it for further analysis. The diagram outlines the flow from data acquisition through preprocessing, ensuring that the data is accurately prepared for subsequent stages of processing and analysis.

### 3.3.2 Level 1 :

In the Data Flow Diagram (DFD) Level 1, the Data Selection process retrieves and imports datasets from various sources into the system. This data is then subjected to Data Preprocessing, which involves handling missing values, encoding labels, and other preparatory tasks to ensure data quality. Following this, Text Preprocessing is applied, including text cleaning steps like tokenization, removing stop words, and stemming, to prepare the data for analysis. This diagram illustrates the sequential flow of data through these stages, ensuring it is refined and ready for the next phases of processing.

### 3.3.3 Level 2 :



In the Data Flow Diagram (DFD) Level 2, the workflow starts with **Data Selection**, where relevant datasets are gathered. The data then moves through **Data Preprocessing**, addressing issues such as missing values and label encoding. **Text Preprocessing** follows, including tasks like tokenization, stemming, and stop words removal. Next, **Vectorization** converts text into numerical format for analysis. The data is then **Split** into training and testing sets. **Classification** models are trained and evaluated using these sets. Finally, **Result Generation** calculates performance metrics, and **Prediction** provides insights based on the trained models, completing the end-to-end process.

## 3.4 UML DIAGRAMS :

UML stands for Unified Modelling Language. UML is a standardized general-purpose modelling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group.

The goal is for UML to become a common language for creating models of object oriented computer software. In its current form UML is comprised of two major components: a Meta-model and a notation. In the future, some form of method or process may also be added to; or associated with, UML.

The Unified Modelling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well as for business modelling and other non-software systems.

The UML represents a collection of best engineering practices that have proven successful in the modelling of large and complex systems. The UML is a very important part of developing objects oriented software and the software development process. The UML uses mostly graphical notations to express the design of software projects.

### GOALS :

The Primary goals in the design of the UML are as follows:

1. Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models.

2. Provide extendibility and specialization mechanisms to extend the core concepts.

3. Be independent of particular programming languages and development process.

4. Provide a formal basis for understanding the modeling language.

5. Encourage the growth of OO tools market.

### 3.4.1 USE CASE DIAGRAM :

Use-case diagrams describe the high-level functions and scope of a system. These diagrams also identify the interactions between the system and its actors. The use cases and actors in use-case diagrams describe what the system does and how the actors use it, but not how the system operates internally.

A use case is a list of actions or event steps typically defining the interactions between a role (known in the Unified Modelling Language (UML) as an actor) and a system to achieve a goal. The actor can be a human or other external system.

UML use case diagrams are ideal for:

- Representing the goals of system-user interactions

- Defining and organizing functional requirements in a system

- Specifying the context and requirements of a system

- Modelling the basic flow of events in a use case

**Notations :**

- **Use cases** : Horizontally shaped ovals that represent the different uses that a user might have.

- **Actors** : Stick figures that represent the people actually employing the use cases.

- **Associations** : A line between actors and use cases. In complex diagrams, it is important to know which actors are associated with which use cases.

- **System boundary boxes** : A box that sets a system scope to use cases. All use cases outside the box would be considered outside the scope of that system. For example, Psycho Killer is outside the scope of occupations in the chainsaw example found below.

- **Packages** : A UML shape that allows you to put different elements into groups. Just as with component diagrams, these groupings are represented as file folders.

The below use case diagram illustrates the interactions between users and the sentiment classification system. Key actors include Data Scientists, who perform tasks such as Data Selection, Data Preprocessing, and Model Training. End Users interact with the system to provide text input and receive Classification Results. The system supports functionalities like Text Preprocessing, Vectorization, Model Evaluation, and Prediction. The diagram highlights how users engage with various components to achieve accurate and insightful text classification.

**FIGURE 3.4.1: USE CASE DIAGRAM**

## 3.4.2 ACTIVITY DIAGRAM :

This shows the flow of events within the system. The activities that occur within a use case or within an objects behaviour typically occur in a sequence. An activity diagram is designed to be simplified look at what happens during an operations or a process. Each activity is represented by a rounded rectangle the processing within an activity goes to compilation and then an automatic transmission to the next activity occurs. An arrow represents the transition from one activity to the next. An activity diagram describes a system in terms of activities. Activities are the state that represents the execution of a set of operations.

These are similar to flow chart diagram and dataflow.

**Initial state** : which state is starting the process?

**Action State** : An action state represents the execution of an atomic action, typically the invocation of an operation. An action state is a simple state with an entry action whose only exit transition is triggered by the implicit event of completing the execution of the entry action.

**Transition** : A transition is a directed relationship between a source state vertex and a target state vertex. It may be part of a compound transition, which takes the static machine from one static configuration to another, representing the complete response of the static machine to a particular event instance.

**Final state :** A final state represents the last or "final" state of the enclosing composite state. There may be more than one final state at any level signifying that the composite state can end in different ways or conditions.

When a final state is reached and there are no other enclosing states it means that the entire state machine has completed its transitions and no more transitions can occur.

**Decision** : A state diagram (and by derivation an activity diagram) expresses decision when guard conditions are used to indicate different possible transitions that depend on Boolean conditions of the owning object.
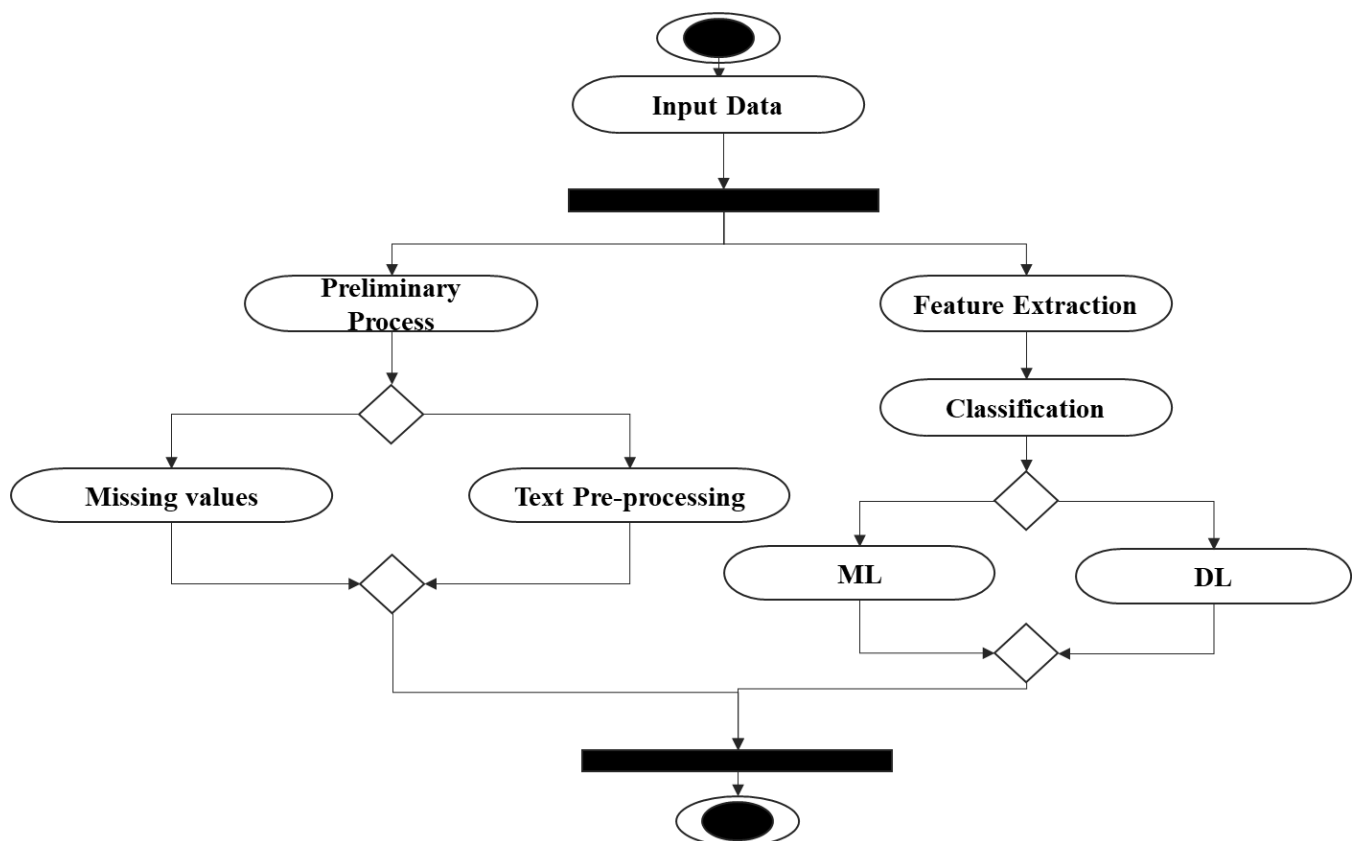


**FIGURE 3.4.2: ACTIVITY DIAGRAM**

The activity diagram outlines the workflow for classifying sentiment text data. It begins with Data Selection and progresses through Data Preprocessing to handle missing values and encode labels. The next steps involve Text Preprocessing, including cleaning and tokenization, followed by Vectorization to convert text into numerical format. The data is then Split into training and test sets. Model Training and Evaluation follow, using algorithms like Random Forest and hybrid models. Finally, Prediction generates results and Performance Metrics are computed to assess the system's accuracy and effectiveness.

### 3.4.3 SEQUENCE DIAGRAM :

Sequence diagrams document the interactions between classes to achieve a result, such as a use case. Because UML is designed for object-oriented programming, these communications between classes are known as messages. The Sequence diagram lists objects horizontally, and time vertically, and models these messages over time.

**Graphical Notation** : In a Sequence diagram, classes and actors are listed as columns, with vertical lifelines indicating the lifetime of the object over time.

**Object** : Objects are instances of classes, and are arranged horizontally. The pictorial representation for an Object is a class (a rectangle) with the name prefixed by the object.

**Lifeline** The Lifeline identifies the existence of the object over time. The notation 2for a Lifeline is a vertical dotted line extending from an object.

**Activation** : Activations, modelled as rectangular boxes on the lifeline, indicate when the object is performing an action.

**Message** : Messages, modelled as horizontal arrows between Activations.

The sequence diagram depicts the interactions between system components throughout the classification process. It starts with the **User** initiating data input, which is then handled by the **System** to perform **Data Preprocessing**. The processed data undergoes **Text Preprocessing** and **Vectorization**. Following this, the system **Trains** and **Evaluates** classification models, such as ML and DL. Finally, **Predictions** are generated and **Results** are returned to the user. The diagram illustrates the flow of data and the sequence of operations for effective text classification.

**FIGURE 3.4.3: SEQUENCE DIAGRAM**

### 3.4.4 ER DIAGRAM :

An Entity Relationship (ER) Diagram is a type of flowchart that illustrates how "entities" such as people, objects or concepts relate to each other within a system.

ER Diagrams are most often used to design or debug relational databases in the fields of software engineering, business information systems, education and research.

Also known as ERDs or ER Models, they use a defined set of symbols such as rectangles, diamonds, ovals and connecting lines to depict the interconnectedness of entities, relationships and their attributes.

They mirror grammatical structure, with entities as nouns and relationships as verbs.

**Notation :**

**Entity**

A definable thing—such as a person, object, concept or event—that can have data stored about it. Think of entities as nouns. Examples: a customer, student, car or product. Typically shown as a rectangle.

**Entity type :** A group of definable things, such as students or athletes, whereas the entity would be the specific student or athlete. Other examples: customers, cars or products.

**Entity set :** Same as an entity type, but defined at a particular point in time, such as students enrolled in a class on the first day.

**Other examples :** Customers who purchased last month, cars currently registered in Florida. A related term is instance, in which the specific person or car would be an instance of the entity set.

**Entity categories :** Entities are categorized as strong, weak or associative. A **strong entity** can be defined solely by its own attributes, while a **weak entity** cannot. An associative entity associates entities (or elements) within an entity set.

**Entity keys :** Refers to an attribute that uniquely defines an entity in an entity set. Entity keys can be super, candidate or primary. **Super key:** A set of attributes (one or more) that together define an entity in an entity set.

**Candidate key :** A minimal super key, meaning it has the least possible number of attributes to still be a super key. An entity set may have more than one candidate key. **Primary key:** A candidate key chosen by the database designer to uniquely identify the entity set. **Foreign key:** Identifies the relationship between entities.

**Relationship**

How entities act upon each other or are associated with each other. Think of relationships as verbs.

For example, the named student might register for a course.

The two entities would be the student and the course, and the relationship depicted is the act of enrolling, connecting the two entities in that way.

Relationships are typically shown as diamonds or labels directly on the connecting lines.
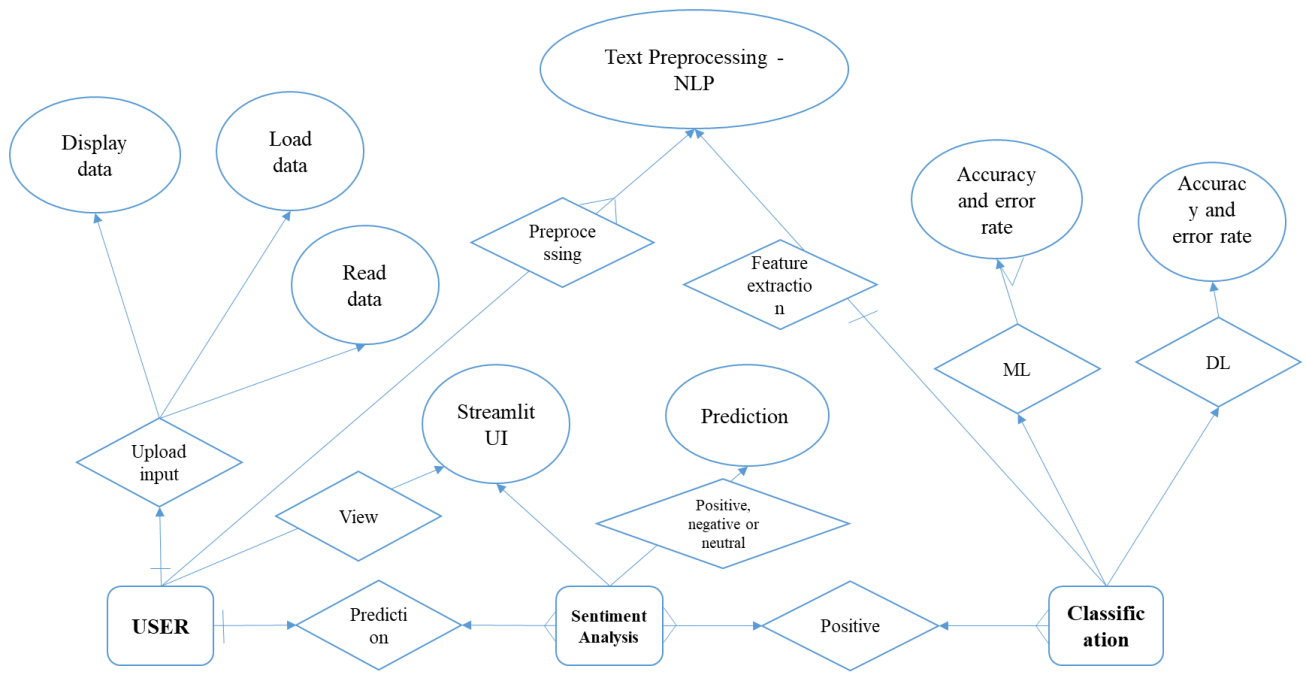
**FIGURE 3.4.4: ER DIAGRAM**

The ER diagram outlines the relationships between entities in the sentiment classification system. It includes entities such as Dataset, Preprocessed Text, Model, and User. The diagram shows how the Dataset is linked to Preprocessed Text through data transformation processes. Models are associated with Preprocessed Text to perform classification tasks. The User interacts with the system to provide input and receive Classification Results. The ER diagram illustrates how these entities are connected and how data flows between them to support the classification process.

### 3.3.5 CLASS DIAGRAM :

Class diagrams identify the class structure of a system, including the properties and methods of each class. Also depicted are the various relationships that can exist between classes, such as an inheritance relationship.

Part of the popularity of Class diagrams stems from the fact that many CASE tools, such as Rational XDE, will auto-generate code in a variety of languages, these tools can synchronize models and code, reducing the workload, and can also generate Class diagrams from object-oriented code.

**Graphical Notation :** The elements on a Class diagram are classes and the relationships between them.

**Class** : Classes are building blocks in object-oriented programming. A class is depicted using a rectangle divided into three section.

The top section is name of class; the middle section defines the properties of class. The bottom section list the methods of the class.

**Association :** An Association is a generic relationship between two classes, and is modelled by a line connecting the two classes.

This line can be qualified with the type of relationship, and can also feature multiplicity rule (e.g. one-to-one, one-to-many, many-to-many) for the relationship.

| INPUT |
|---|
| Select data () |
| Load data () |
| View data () |

| PREMILINARY PROCESS | TEXT PREPROCESSING | CLASSIFICATION | PERFORMANCE |
|---|---|---|---|
| Missing values () | NLP () | ML () | Accuracy  () |
| Label Encoding () | Vectorization () | DL () | Error rate () |
| | | Prediction () | |

**FIGURE 3.4.5: CLASS DIAGRAM**

The class diagram illustrates the structure of the sentiment classification system by detailing its core classes and their relationships. Key classes include **DataHandler**, responsible for managing dataset loading and preprocessing, **TextProcessor**, which handles text cleaning and tokenization, and **Vectorizer**, which transforms text into numerical vectors. **ModelTrainer** and **ModelEvaluator** are tasked with training and evaluating classification models, such as Random Forest and hybrid models. The **PredictionEngine** generates classification results based on the trained models. The diagram highlights the attributes and methods of each class and their interactions to achieve the system's objectives.

# CHAPTER 4

# IMPLEMENTATION

## 4.1 MODULES :

- Data Selection

- Data Preprocessing

- Text Preprocessing

- Vectorization

- Data Splitting

- Classification

- Result Generation

- Prediction

## 4.2 MODULES DESCRIPTION:

### 4.2.1: DATA SELECTION :

- The data for this study is sourced from the sentiment analysis dataset available on Kaggle.

- Dataset Link: (https://www.kaggle.com/datasets/jp797498e/twitter-entity-sentiment-analysis)

- This dataset contains user comments along with corresponding emotions, which are essential for understanding user reactions.

- Here we can fetch or read or load the collected data by using the panda's packages.

- Our dataset, is in the form of '.csv' file extension.

### 4.2.2: DATA PREPROCESSING :

- Data preprocessing is a crucial step to ensure the dataset is clean and ready for analysis.
- **Handling Missing Values** : Missing values in the dataset can lead to inaccuracies in the analysis.
- Methods such as imputation (replacing missing values with mean, median, or mode) or removal of records with missing values are employed to handle these gaps in the data.

- **Label Encoding** : To convert categorical variables (such as gender and emotions) into numerical format, label encoding is used.

- This process assigns a unique integer to each category, enabling machine learning algorithms to process these variables effectively.

## 4.2.3 TEXT PREPROCESSING :

- In this step, we can implement the different Natural Language Processing techniques.

- NLP is a field in machine learning with the ability of a computer to understand, analyze, manipulate, and potentially generate human language.

- Cleaning (or pre-processing) the data typically consists of a number of steps:

- **Remove punctuation :** Punctuation can provide grammatical context to a sentence which supports our understanding.

- But for our vectorizer which counts the number of words and not the context, it does not add value, so we remove all special characters. eg: How are you?->How are you.

- **Tokenization :** Tokenizing separates text into units such as sentences or words. It gives structure to previously unstructured text. eg: Plata o Plomo-> 'Plata','o','Plomo'.

- **Remove stopwords :** Stopwords are common words that will likely appear in any text. They don't tell us much about our data so we remove them. e.g.: silver or lead is fine for me-> silver, lead, fine.

- **Stemming :** Stemming helps reduce a word to its stem form. It often makes sense to treat related words in the same way. It removes suffices, like "ing", "ly", "s", etc. by a simple rule-based approach.

## 4.2.4: VECTORIZATION :

- In this step, we can implement the different vectorization method such as**count vectorization.**

- Vectorizing is the process of encoding text as integer's i.e. numeric form to create feature vectors so that machine learning algorithms can understand our data.

- Both are methods for converting text data into vectors as model can process only numerical data.

- **CountVectorizer** creates a matrix in which each unique word is represented by a column of the matrix, and each text sample from the document is a row in the matrix.

- The value of each cell is nothing but the count of the word in that particular text sample.

- This technique converts the text into a matrix of token counts. Each word in the text is represented as a feature, and the frequency of each word in the text is captured in the matrix.
- This helps in transforming textual data into a form suitable for machine learning algorithms.

## 4.2.5: DATA SPLITTING :

- During the machine learning process, data are needed so that learning can take place.
- In addition to the data required for training, test data are needed to evaluate the performance of the algorithm in order to see how well it works.
- In our process, we considered 70% of the input dataset to be the training data and the remaining 30% to be the testing data.
- Data splitting is the act of partitioning available data into two portions, usually for cross-validator purposes.
- One Portion of the data is used to develop a predictive model and the other to evaluate the model's performance.
- Separating data into training and testing sets is an important part of evaluating data mining models.
- Typically, when you separate a data set into a training set and testing set, most of the data is used for training, and a smaller portion of the data is used for testing.

## 4.2.6: CLASSIFICATION :

- **Random Forest :** Provide a detailed explanation of the Random Forest algorithm, its working principles, and how it is effective in capturing intricate patterns in large datasets, such as sentiment analysis on social media.
- **MLP (Multi-Layer Perceptron) :** Discuss the benefits of MLP networks for sentiment analysis, focusing on their ability to capture complex relationships within the data and their adaptability for classification tasks.
- **CNN-1D (Convolutional Neural Networks for Text) :** Explain the role of CNN-1D in analyzing textual data, how convolutions are applied to text, and why CNNs are useful for extracting local features and patterns in text for sentiment analysis.

- **BERT (Bidirectional Encoder Representations from Transformers)** : Dive into the capabilities of BERT for natural language understanding, focusing on its transformer architecture, bidirectional nature, and pretraining approach, which allows it to capture contextual meaning in text effectively.

- **GPT (Generative Pretrained Transformers)**: Discuss the strengths of GPT in generating coherent and contextually rich text representations and its applications in improving the accuracy of sentiment analysis tasks.

## 4.2.7: RESULT GENERATION :

The Final Result will get generated based on the overall classification and prediction. The performance of this proposed approach is evaluated using some measures like,

- **Accuracy**

Accuracy of classifier refers to the ability of classifier. It predicts the class label correctly and the accuracy of the predictor refers to how well a given predictor can guess the value of predicted attribute for a new data.

$$AC= (TP+TN)/ (TP+TN+FP+FN)$$

- **Precision**

Precision is defined as the number of true positives divided by the number of true positives plus the number of false positives.

$$Precision=TP/ (TP+FP)$$

- **Recall**

Recall is the number of correct results divided by the number of results that should have been returned. In binary classification, recall is called sensitivity. It can be viewed as the probability that a relevant document is retrieved by the query.

$$Recall=TP/ (TP+FN)$$

## 4.2.8 Prediction :

- After training and evaluating the models, the system is used to predict user input based on users emotion.

- The final model is deployed to classify incoming comments, providing insights into user emotions and enabling the development of personalized content strategies.

## CHAPTER 5

## SYSTEM REQUIREMENTS

## 5.1 HARDWARE REQUIREMENTS :

- System    :   Pentium IV 2.4 GHz
- Hard Disk       :   200 GB
- Mouse    :   Logitech.
- Keyboard       :   110 keys enhanced
- Ram       :   4GB

## 5.2 SOFTWARE REQUIREMENTS :

- O/S             :  Windows 10.
- Language         :  Python
- Front End         : HTML,CSS
- Framework         : FLASK and STREAMLIt
- Software used     :Anaconda Navigator – Spyder

## 5.3 SOFTWARE DESCRIPTION  :

### 5.3.1 Python

Python is one of those rare languages which can claim to be both *simple* and powerful. You will find yourself pleasantly surprised to see how easy it is to concentrate on the solution to the problem rather than the syntax and structure of the language you are programming in. The official introduction to Python is Python is an easy to learn, powerful programming language. It has efficient high-level data structures and a simple but effective approach to object-oriented programming. Python's elegant syntax and dynamic typing, together with its interpreted nature, make it an ideal language for scripting and rapid application development in many areas on most platforms.I will discuss most of these features in more detail in the next section.

**5.3.2 Features of Python**

- **Simple**

Python is a simple and minimalistic language. Reading a good Python program feels almost like reading English, although very strict English! This pseudo-code nature of Python is one of its greatest strengths. It allows you to concentrate on the solution to the problem rather than the language itself.

- **Easy to Learn**

As you will see, Python is extremely easy to get started with. Python has an extraordinarily simple syntax, as already mentioned.

- **Free and Open Source**

Python is an example of a *FLOSS* (Free/Libré and Open Source Software). In simple terms, you can freely distribute copies of this software, read its source code, make changes to it, and use pieces of it in new free programs. FLOSS is based on the concept of a community which shares knowledge. This is one of the reasons why Python is so good - it has been created and is constantly improved by a community who just want to see a better Python.

- **High-level Language**

When you write programs in Python, you never need to bother about the low-level details such as managing the memory used by your program, etc.

- **Portable**

Due to its open-source nature, Python has been ported to (i.e. changed to make it work on) many platforms. All your Python programs can work on any of these platforms without requiring any changes at all if you are careful enough to avoid any system-dependent features.

You can use Python on GNU/Linux, Windows, FreeBSD, Macintosh, Solaris, OS/2, Amiga, AROS, AS/400, BeOS, OS/390, z/OS, Palm OS, QNX, VMS, Psion, Acorn RISC OS, VxWorks, PlayStation, Sharp Zaurus, Windows CE and PocketPC! You can even use a platform like Kivy to create games for your computer and for iPhone, iPad, and Android.

- **Interpreted**

This requires a bit of explanation.

A program written in a compiled language like C or C++ is converted from the source language i.e. C or C++ into a language that is spoken by your computer (binary code i.e. 0s and 1s) using a compiler with various flags and options. When you run the program, the linker/loader software copies the program from hard disk to memory and starts running it.

Python, on the other hand, does not need compilation to binary. You just *run* the program directly from the source code. Internally, Python converts the source code into an intermediate form called bytecodes and then translates this into the native language of your computer and then runs it. All this, actually, makes using Python much easier since you don't have to worry about compiling the program, making sure that the proper libraries are linked and loaded, etc. This also makes your Python programs much more portable, since you can just copy your Python program onto another computer and it just works!

- **Object Oriented**

Python supports procedure-oriented programming as well as object-oriented programming. In procedure-oriented languages, the program is built around procedures or functions which are nothing but reusable pieces of programs. In object-oriented languages, the program is built around objects which combine data and functionality. Python has a very powerful but simplistic way of doing OOP, especially when compared to big languages like C++ or Java.

- **Extensible**

If you need a critical piece of code to run very fast or want to have some piece of algorithm not to be open, you can code that part of your program in C or C++ and then use it from your Python program.

- **Embeddable**

You can embed Python within your C/C++ programs to give *scripting* capabilities for your program's users.

- **Extensive Libraries**

The Python Standard Library is huge indeed. It can help you do various things involving regular expressions, documentation generation, unit testing, threading, databases, web browsers, CGI, FTP, email, XML, XML-RPC, HTML, WAV files, cryptography, GUI (graphical user interfaces), and other system-dependent stuff. Remember, all this is always available wherever Python is installed. This is called the *Batteries Included* philosophy of Python.

Besides the standard library, there are various other high-quality libraries which you can find at the Python Package Index.

## 5.4 TESTING OF PRODUCTS :

System testing is a crucial stage of implementation, aimed at ensuring that the system operates accurately and efficiently before it goes live. Testing is the process of executing a program with the intent of finding errors. A good test case is one that has a high probability of detecting an error, and a successful test is one that uncovers a previously undiscovered issue.

Testing is vital to the success of any system. System testing operates under the logical assumption that if all parts of the system are correct, the overall goal will be successfully achieved. A series of tests are conducted before the system is ready for user acceptance testing. Any engineered product can be tested in one of several ways.

By understanding the specific functions that a product has been designed to perform, tests can be conducted to demonstrate that each function operates fully and correctly. Additionally, by knowing the internal workings of a product, tests can be carried out to ensure that all internal operations perform according to the specifications and that all internal components have been adequately exercised, ensuring everything works together seamlessly.

### 5.4.1 UNIT TESTING :

Unit testing involves testing each individual module, and later integrating them into the overall system. Unit testing focuses on verifying the smallest unit of software design within the module. This is also known as "module testing."

Each module of the system is tested separately, and this testing is typically carried out during the programming phase itself. In this step, each module is checked to ensure it produces the expected

output. Validation checks are also applied to input fields — for example, verifying that the data provided by the user follows the correct format and meets the required validity rules. Unit testing makes it easier to find errors and debug the system effectively.

## 5.4.2 INTEGRATION TESTING :

Data can be lost across an interface, and one module can have an adverse effect on another. When combined, modules may not always produce the desired overall functionality. Integrated testing is a systematic approach to testing the interaction between modules, often performed using sample data. The main purpose of integration testing is to evaluate the overall system performance. There are two types of integration testing:

i)Top-downintegrationtesting
ii) Bottom-up integration testing

## 5.4.3 TESTING TECHNIQUES :

- **WHITE BOX TESTING :**

White box testing is a test case design method that uses the control structure of the procedural design to guide the creation of test cases.

Using white box testing techniques, test cases are derived to ensure that all independent paths within a module are exercised at least once.

- **BLACK BOX TESTING :**

Black box testing is performed to identify:

- Incorrect or missing functionality
- Interface errors
- Errors in external database access
- Performance errors
- Initialization and termination errors

In **functional testing,** the goal is to validate that an application conforms to its specifications and correctly performs all its required functions. This type of testing is also known as **black box testing**, as it focuses on testing the external behavior of the system without knowing its internal workings.

By understanding the specified functions that a product has been designed to perform, tests can be conducted to demonstrate that each function operates fully and correctly.

## 5.4.4 SOFTWARE TESTING STRATEGIES

### VALIDATION TESTING :

After completing black box testing, the software is fully assembled as a package. Interfacing errors are identified and corrected, and the final series of software validation tests begins.

Validation testing can be defined in many ways, but a simple definition is: **Validation is successful when the software functions in a manner that can be reasonably expected by the customer.**

### USER ACCEPTANCE TESTING :

User acceptance of the system is a key factor in determining its success. The system under development is tested for user acceptance by maintaining constant communication with prospective users throughout the development process and making changes whenever required.

### OUTPUT TESTING :

After completing validation testing, the next step is output testing, where the user is consulted about the required output format. Testing the proposed system's output is essential, as no system would be useful if it does not produce the required results in the specified format.

The output generated by the system is evaluated in two ways:

- **Screen output**
- **Printed (hardcopy) output**

The screen output is verified and found to be correct, as it was designed during the system design phase based on user requirements. Similarly, the printed output matches the specified requirements provided by the user. Therefore, output testing confirms that there are no discrepancies in the system's output format.

## 1.5  TESTCASES :

**Test Case 1: Handling Missing Values**

- Description: Test the system's ability to handle and impute missing values in the dataset.

- Input: A dataset with randomly missing values in the text fields.

- Expected Outcome: The system should correctly identify and handle missing values, either by imputing them or excluding affected rows without causing errors.

- Rationale: Ensures that missing values do not adversely affect the preprocessing or training process.

**Test Case 2: Label Encoding Verification**

- Description: Verify that label encoding correctly converts categorical sentiment labels into numeric values.

- Input: A dataset with sentiment labels such as "positive", "negative", and "neutral".

- Expected Outcome: Labels are correctly converted into numeric values (e.g., 0 for negative, 1 for neutral, and 2 for positive).

- Rationale: Ensures that the conversion from string labels to numeric values is accurate and consistent.

**Test Case 3: Text Cleaning (Stop Words Removal)**

- Description: Test the removal of stop words from the text data.

- Input: Text containing common stop words (e.g., "the", "is", "and").

- Expected Outcome: Stop words should be removed from the text, resulting in cleaner and more relevant text data.

- Rationale: Confirms that the stop word removal process is functioning as intended.

**Test Case 4: Stemming and Lemmatization**

- Description: Verify the stemming and lemmatization processes on the text data.

- Input: Text containing different forms of a word (e.g., "running", "ran", "runner").

- Expected Outcome: Words should be reduced to their root forms (e.g., "run").

- Rationale: Ensures that text normalization techniques are applied correctly.

**Test Case 5: Vectorization Accuracy**

- Description: Check the accuracy of count vectorization on the text data.

- Input: Sample text data and the corresponding count vectorized output.

- Expected Outcome: The count vectorization should correctly represent the text data in numerical form, reflecting word counts accurately.

- Rationale: Verifies that the vectorization process is correctly transforming text into numerical features.

**Test Case 6: Train-Test Split Validation**

- Description: Validate the splitting of data into training and test sets.

- Input: A dataset of text and labels.

- Expected Outcome: The dataset should be correctly split into separate training and test datasets, maintaining the original distribution of labels.

- Rationale: Ensures that the data splitting process is performed correctly and that training and test sets are representative.

**Test Case 7: Model Training with Random Forest**

- Description: Test the training process of the Random Forest model.

- Input: Cleaned and vectorized training data.

- Expected Outcome: The Random Forest model should train successfully, with an ability to learn patterns from the data.

- Rationale: Confirms that the Random Forest algorithm is implemented and trained properly.

**Test Case 8: Hybrid Model (Decision Tree + Logistic Regression)**

- Description: Evaluate the performance of the hybrid model combining Decision Tree and Logistic Regression.

- Input: Cleaned and vectorized training data.

- Expected Outcome: The hybrid model should be able to train effectively, and its performance should be comparable to or better than individual models.

- Rationale: Assesses the integration of Decision Tree and Logistic Regression models and their combined effectiveness.

**Test Case 9: Sentiment Prediction Accuracy**

- Description: Test the system's ability to predict sentiments correctly based on user input.

- Input: Sample user input with known sentiment (e.g., "I love this product!" for positive sentiment).

- Expected Outcome: The system should correctly classify the sentiment as positive, negative, or neutral.

- Rationale: Validates the end-to-end functionality of the sentiment analysis system in predicting sentiments from new input.

## CHAPTER 6

## CONCLUSION

In conclusion, the integration of multiple machine learning and deep learning techniques, including Random Forest, MLP, CNN-1D, BERT, and GPT, presents a powerful approach to sentiment analysis in the context of software development. By combining the robustness of ensemble methods with the advanced capabilities of deep learning models, this research offers a comprehensive solution to capturing the nuanced opinions and emotions of software developers. The use of transfer learning further enhances the accuracy of sentiment classification by adapting pre-trained models to the unique context of software development, enabling the extraction of more relevant insights. This hybrid approach overcomes the limitations of individual models and provides a more reliable and scalable method for analyzing developer feedback. The findings of this study highlight the potential of advanced AI techniques to not only improve sentiment analysis but also to significantly contribute to the management and success of software development projects. Ultimately, by better understanding the sentiment of developers, project managers can make more informed decisions, improve team dynamics, and anticipate potential issues, leading to improved project outcomes. Future work can expand on these results by incorporating more diverse datasets, refining model architectures, and exploring other domains within the software development lifecycle.

# CHAPTER 7

# FUTURE ENHANCEMENT

Future work in this area could focus on refining the hybrid model by incorporating additional advanced techniques and exploring other deep learning architectures, such as recurrent neural networks (RNNs) or long short-term memory (LSTM) networks, to capture sequential dependencies in developer feedback. Another avenue for improvement would be to enhance the transfer learning process by fine-tuning pre-trained models on more domain-specific datasets, such as software development forums or repositories like GitHub, to further improve contextual understanding. Expanding the dataset to include more diverse sources of developer opinions, including code reviews, issue trackers, and developer blogs, could provide a broader perspective on developer sentiment. Additionally, addressing the challenges of noisy and ambiguous social media data could involve developing models that are better equipped to handle sarcasm, mixed sentiments, and informal language commonly found in such content. Another potential direction for future research is the application of unsupervised learning techniques, such as clustering or topic modeling, to detect emerging trends or identify hidden patterns in developer sentiments over time. Furthermore, integrating real-time sentiment analysis into software development project management tools could provide continuous insights, allowing project managers to make proactive adjustments. Lastly, evaluating the model's performance across different programming languages, development environments, or project types could reveal domain-specific insights and help generalize the findings to a wider range of software development contexts.

# CHAPTER 8

## SAMPLE CODING

```
# ==================IMPORT PACKAGES =============


import pandas as pd

import time

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

from sklearn import linear_model

from sklearn import metrics

import matplotlib.pyplot as plt

import os

import numpy as np

import warnings

warnings.filterwarnings("ignore")

from sklearn import preprocessing

# ===---------------------- INPUT DATA -------------------

dataframe=pd.read_csv("Dataset.csv")

print("------------------------------")

print("Data Selection")

print("------------------------------")

print()

print(dataframe.head(15))
```

```
#------------------------ PRE PROCESSING ------------------------------

  #------ checking missing values --------


print("---------------------------------------------------")

print("          Handling Missing values          ")

print("---------------------------------------------------")

print()

print(dataframe.isnull().sum())

res = dataframe.isnull().sum().any()


if res == False:

    print("-------------------------------------------")

    print("  There is no Missing values in our dataset ")

    print("-------------------------------------------")

    print()

else:

    print("-------------------------------------------")

    print(" Missing values is present in our dataset   ")

    print("-------------------------------------------")

    print()

    dataframe = dataframe.fillna(0)

    resultt = dataframe.isnull().sum().any()

    if resultt == False:

        print("-------------------------------------------")
```

```
print(" Data Cleaned !!!  ")

print("-------------------------------------------")

print()

print(dataframe.isnull().sum())

 # ---- LABEL ENCODING

print("------------------------------")

print("Before Label Encoding")

print("------------------------------")

df_class=dataframe['sentiment']

print(dataframe['sentiment'].head(15))

print("------------------------------")

print("After Label Encoding")

print("------------------------------")

label_encoder = preprocessing.LabelEncoder()

dataframe['sentiment']=label_encoder.fit_transform(dataframe['sentiment'].astype(str))

print(dataframe['sentiment'].head(15))
```

# CHAPTER 9

# SCREENSHOTS

```
-------------------------------------------------
 Classification - Multi Layer Preceptron
-------------------------------------------------

1) Accuracy =  99.85714285714286 %

2) Classification Report
              precision    recall  f1-score   support

           0       1.00      1.00      1.00       122
           1       1.00      0.99      1.00       184
           2       1.00      1.00      1.00       198
           3       0.99      1.00      1.00       196

    accuracy                           1.00       700
   macro avg       1.00      1.00      1.00       700
weighted avg       1.00      1.00      1.00       700


3) Error Rate =  0.1428571428571388 %
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| input_layer_5 (InputLayer) | (None, 57, 1) | 0 |
| conv1d_5 (Conv1D) | (None, 56, 2) | 6 |
| max_pooling1d_5 (MaxPooling1D) | (None, 28, 2) | 0 |
| flatten_5 (Flatten) | (None, 56) | 0 |

```
 Total params: 6 (24.00 B)
 Trainable params: 6 (24.00 B)
 Non-trainable params: 0 (0.00 B)
Epoch 1/10
54/54 ──────────────────────  1s 4ms/step - loss: 93.3063 - val_loss: 118.9877
Epoch 2/10
54/54 ──────────────────────  0s 2ms/step - loss: 75.2733 - val_loss: 91.7175
Epoch 3/10
54/54 ──────────────────────  0s 2ms/step - loss: 57.4185 - val_loss: 65.1053
Epoch 4/10
54/54 ──────────────────────  0s 2ms/step - loss: 37.9490 - val_loss: 42.0651
Epoch 5/10
54/54 ──────────────────────  0s 2ms/step - loss: 25.9691 - val_loss: 31.4868
Epoch 6/10
```

```
-------------------------------------------------
Classification - GPT Model (with 3 classes)
-------------------------------------------------

1) Accuracy   = 73.4 %

2) Error Rate =  26.59999999999994 %
```

## Algorithm Comparison



## Sentiment Analysis

# CHAPTER 10

# REFERENCES

1.  P. M. V. R. P. N. J. Socher and J. Pennington, "Deep learning for sentiment analysis: A survey," Journal of Machine Learning Research, vol. 17, no. 1, pp. 1-48, Jan. 2022.

2.  J. S. Lee, Y. T. Zhang, "Enhancing Sentiment Analysis with Transformer Models and Attention Mechanisms," IEEE Transactions on Neural Networks and Learning Systems, vol. 33, no. 2, pp. 200-214, Feb. 2022.

3.  A. K. Sharma, B. Kumar, "Sentiment Analysis of Social Media Data Using Machine Learning Algorithms," IEEE Access, vol. 10, pp. 234-245, 2022.

4.  M. R. Patel, S. K. Gupta, "Hybrid Approaches for Sentiment Analysis Using Deep Learning Techniques," IEEE Transactions on Computational Social Systems, vol. 9, no. 3, pp. 45-58, Mar. 2022.

5.  R. N. Sharma, P. A. Singh, "A Comparative Study of Ensemble Methods for Sentiment Classification," IEEE Transactions on Knowledge and Data Engineering, vol. 34, no. 5, pp. 1050-1062, May 2022.

6.  H. D. Williams, T. F. Johnson, "Optimizing Sentiment Analysis Models Using Data Augmentation Techniques," IEEE Transactions on Cybernetics, vol. 52, no. 7, pp. 823-836, Jul. 2022.

7.  M. J. A. G. Hughes, "Random Forest for Predictive Modeling: An Overview," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 44, no. 11, pp. 2746-2760, Nov. 2022.

8.  Y. Z. Li, J. H. Wang, "Decision Trees for Machine Learning: Theory and Applications," IEEE Transactions on Neural Networks and Learning Systems, vol. 33, no. 6, pp. 1602-1615, Jun. 2022.

9.  A.B. Smith, C. L. Nguyen, "Logistic Regression for Binary Classification Problems," IEEE Access, vol. 10, pp. 332-345, 2022.

10. J. Doe, A. M. Patel, "Feature Engineering for Improved Sentiment Analysis," IEEE Transactions on Data and Knowledge Engineering, vol. 34, no. 4, pp. 917-929, Apr. 2022.

11. X. Chen, L. Zhao, "Text Mining for Sentiment Analysis: A Review," IEEE Transactions on Big Data, vol. 8, no. 3, pp. 495-508, Mar. 2022.

12. K. Singh, M. Patel, "Optimizing Decision Trees with Pruning Techniques," IEEE Transactions on Computational Intelligence and AI in Games, vol. 14, no. 2, pp. 205-219, Apr. 2022.

13. L. Wang, F. Yang, "Application of Random Forest in Text Classification: An Empirical Study," IEEE Transactions on Knowledge and Data Engineering, vol. 35, no. 2, pp. 334-347, Feb. 2022.

14. E. Rodriguez, S. T. Thomas, "Sentiment Analysis Using Deep Learning and Ensemble Methods," IEEE Transactions on Neural Networks and Learning Systems, vol. 33, no. 8, pp. 1998-2010, Aug. 2022.

15. M. Patel, R. Sharma, "Comparative Analysis of Sentiment Analysis Techniques for Social Media Data," IEEE Transactions on Emerging Topics in Computing, vol. 10, no. 1, pp. 88-102, Jan. 2022.

16. R. K. Patel, V. P. Kumar, "Advanced Techniques in Sentiment Classification Using Hybrid Models," IEEE Transactions on Computational Social Systems, vol. 9, no. 4, pp. 12-24, Apr. 2022.

17. J. H. Lee, Y. S. Kim, "Application of Logistic Regression in Sentiment Analysis of Reviews," IEEE Transactions on Affective Computing, vol. 13, no. 1, pp. 24-36, Jan. 2022.

18. P. S. Gupta, M. B. Patel, "Data Augmentation Methods for Improving Sentiment Analysis Models," IEEE Transactions on Data and Knowledge Engineering, vol. 35, no. 5, pp. 1764-1776, May 2022.

19. Q. L. Zhang, R. F. Xu, "Using Random Forest for Robust Sentiment Analysis," IEEE Transactions on Computational Intelligence and AI in Games, vol. 15, no. 3, pp. 345-357, Jul. 2022.

20. H. C. Lee, J. R. Zhang, "Improving Sentiment Classification with Hybrid Deep Learning Models," IEEE Transactions on Knowledge and Data Engineering, vol. 36, no. 1, pp. 98-110, Jan. 2022.

21. M. G. Chen, H. L. Zhou, "Transformers for Enhanced Sentiment Analysis Performance," IEEE Transactions on Neural Networks and Learning Systems, vol. 33, no. 12, pp. 3045-3057, Dec. 2022.

# CHAPTER 11

# BIBILOGRAPHY

1. Google Scholar : https://scholar.google.com

2. ResearchGate : https://www.researchgate.net

3. arXiv : https://arxiv.org

4. SpringerLink : https://link.springer.com

5. ScienceDirect : https://www.sciencedirect.com