

# **Disaster Tweet Analyzer: Natural Language Processing for Crisis Communication**

A Project Report

*Submitted in complete fulfilment of the  
requirements of 2-month internship*

Internship Duration:

**3 October 2024 to 28 November 2024**

Under the Guidance of:

**Mentor: Mr. Nitig Singh**

***Team-Leads:*** Vanshika Ganjoo, Prajakta Tambare, Rani Soni, Ashutosh Kumar Yadav

***Co-Lead:*** Shubham, Moksha, Hanifa, Geetha

***Team Members:***

**TEAM-1:** Kaushal, Sriram, Keerthi

**TEAM-2:** Bustob, Esthak, Kanchana, Neha

**TEAM-3:** Avantika, Karthikye, Ruthvik

**TEAM-4:** Ayush, Karna, Shreekar



**Infosys Springboard Internship**

## ACKNOWLEDGEMENTS

We extend our heartfelt gratitude to Mr. Nitig Singh Sir, our esteemed mentor from Infosys Springboard, for his invaluable guidance and support throughout this project. His insights and expertise have been instrumental in shaping our understanding and approach to disaster tweet analysis.

We are immensely grateful to Infosys and the Infosys Springboard Committee for providing us with this incredible opportunity to work on a meaningful project that combines learning and practical application. Their platform has been pivotal in enabling us to enhance our skills and contribute to solving real-world problems.

We deeply appreciate the efforts of our team leads—Vanshika Ganjoo, Prajakta Tambare, Rani Soni, and Ashutosh Kumar Yadav—whose leadership and dedication steered the project toward success. Their ability to guide the team with clarity and determination has been truly inspiring.

Our gratitude also goes to the co-leads, Shubham, Moksha, Hanifa, and Geetha, for their significant contributions and proactive efforts in ensuring seamless collaboration and task execution. We acknowledge the hard work of every team member, including Bustob, Kaushal, and many others, who brought diverse skills and perspectives to the table, making this project a collective triumph.

Lastly, we are thankful for the resources and tools that facilitated this journey, such as the Kaggle Disaster Tweets dataset, which provided the foundation for our work. This project has been a transformative learning experience that would not have been possible without the combined efforts and support of everyone involved.

## **TABLE OF CONTENTS**

<b>Acknowledgment .....</b>	<b>(i)</b>
<b>List .....</b>	<b>(iv)</b>
• List of Listings .....	(iv)
• List of Images .....	(iv)
• List of Abbreviations .....	(v)
<b>Preface .....</b>	<b>(vi)</b>
<b>Abstract .....</b>	<b>1</b>
<b>Chapter 1: Introduction .....</b>	<b>2</b>
• 1.1 Introduction to Disasters .....	2
• 1.2 Importance of Information during Disaster .....	3
• 1.3 Social Media and its Role during Disaster .....	5
• 1.4 Objective of the Project .....	6
• 1.5 Outline of the Project Report .....	7
<b>Chapter 2: Analysis of Existing Solutions in Disaster Tweet Classification .....</b>	<b>9</b>
• 2.1 Introduction to the Chapter .....	9
• 2.2 Evolution of Disaster Tweet Analysis and Related Work .....	9
• 2.3 Gaps and Limitations in Existing Systems .....	10
• 2.4 How Our Project Bridges the Gaps .....	10
<b>Chapter 3: Data &amp; Methodology .....</b>	<b>12</b>
• 3.1 Introduction to this Chapter .....	12
• 3.2 About Our Data .....	13
• 3.3 Preprocessing .....	14
• 3.4 Feature Engineering .....	17
• 3.5 Dataset Modifications and Creation .....	19
• 3.6 Model Training .....	20
• 3.7 Model Web-Interface .....	26
<b>Chapter 4: Results and Discussion .....</b>	<b>28</b>
• 4.1 Machine Learning Workflow and Visualization .....	28
• 4.2 Performance of Learning Models .....	29
• 4.3 Accuracy Comparison .....	31

• 4.4 User Interface and Prediction Output .....	32
<b>Chapter 5: Reflections, Insights, and Future Directions .....</b>	<b>34</b>
• 5.1 Purpose and Scope of the Chapter .....	34
• 5.2 Challenges Encountered During Development .....	34
• 5.3 Accomplishments and Project Outcomes .....	35
• 5.4 Future Directions and Scope for Enhancement .....	36
<b>Appendices .....</b>	<b>38</b>
• Appendix A: Data Collection Details .....	38
• Appendix B: Evaluation Metrics Explained .....	38
• Appendix C: Model Performance Summary .....	39
• Appendix D: Tools and Libraries Used .....	39
<b>References .....</b>	<b>40</b>

# LISTS

## **1. Lists of Listings**

Listing 3.3.1: Preprocessing and Location Extraction

Listing 3.6.1: Linear Regression Model

Listing 3.6.2: Random Forest Model

Listing 3.6.3: Gradient Boosting Model

Listing 3.6.4: Logistic Regression Model

Listing 3.6.5: XG-Boost Model

Listing 3.6.6: KNN and Naive Bayes Model

Listing 3.6.7: CNN, Gradient Descent & Decision Tree Model

## **2. Lists of Figures**

Fig 1.1.1: Haiti earthquake damage

Fig 1.3.1: Tasmania's bushfires have become a prominent topic online

Fig 2.4.1: Sentiment Analysis

Fig 3.1.1: Workflow

Fig 3.4.4.1: Disaster Location Extraction form Tweets

Fig 3.4.2: Word Cloud of Disaster Related Tweets

Fig 3.6.1: Model Accuracy Comparison Chart

Fig 4.1.1: TF-IDF Score Bar-graph

Fig 4.2.1: Linear Regression Confusion Matrix

Fig 4.2.2: Linear SVM Confusion Matrix

Fig 4.2.3: Gradient Boost Confusion Matrix

Fig 4.2.4: KNN Confusion Matrix

Fig 4.3.1: Model Accuracy Chart

Fig 4.4.1: Sentiment Analysis

### 3. Lists of Abbreviations

AI - Artificial Intelligence

ML - Machine Learning

DL - Deep Learning

NLP - Natural Language Processing

LR - Logistic Regression

DTA - Disaster Tweet Analyzer

RT - Relevant Tweet

NRT - Non-Relevant Tweet

DRT - Disaster-Related Tweet

NDRT - Non-Disaster-Related Tweet

ID - Identifier

KW - Keyword

LOC - Location

TXT - Text

TGT - Target (0 = non-disaster, 1 = Disaster)

EDA - Exploratory Data Analysis

TF-IDF - Term Frequency-Inverse Document Frequency

Bow - Bag of Words

POS - Part of Speech

ROC - Receiver Operating Characteristic

## PREFACE

This project report, titled "**Disaster Tweet Analyzer**," is a culmination of the work undertaken as part of the **Infosys Springboard Artificial Intelligence Internship**. The project leverages advanced techniques in **Natural Language Processing (NLP)** and **Machine Learning (ML)** to classify tweets as disaster-related or non-disaster-related.

The increasing reliance on social media as a source of real-time information during disasters has highlighted the need for automated systems to filter and analyse large volumes of tweets. This project aims to address this challenge by building a robust classifier that can assist authorities and organizations in identifying relevant tweets quickly and effectively.

The report covers the step-by-step process of the project, including data preprocessing, feature engineering, model development, and evaluation. A logistic regression model was chosen as the final classifier for its simplicity, interpretability, and efficiency in handling text classification tasks.

This project has been an incredible learning journey, exposing me to real-world challenges in data preprocessing, model selection, and deployment. I am grateful for the opportunity to contribute to a socially impactful project under the mentorship of **Nitig Sir** and with the collaborative efforts of my team.

## ABSTRACT

The **Disaster Tweet Analyzer** project focuses on leveraging **Artificial Intelligence (AI)** and **Natural Language Processing (NLP)** techniques to classify tweets as disaster-related or non-disaster-related. With the increasing use of social media during emergencies, timely identification of relevant information has become crucial for effective disaster response and management.

This project utilizes a dataset of over 11,000 tweets labelled as disaster or non-disaster. The process involves data preprocessing to clean and standardize the text, feature engineering to extract meaningful patterns, and model development using machine learning techniques. A logistic regression model was selected as the final classifier for its balance of accuracy and computational efficiency.

The project also incorporates a merged dataset for enhanced classification accuracy, adding a new target column to indicate tweet relevance. Key evaluation metrics such as precision, recall, F1-score, and accuracy were used to assess the model's performance.

The outcomes of this project demonstrate the potential of AI-driven solutions in aiding disaster management efforts by filtering and categorizing large-scale social media data. This work provides a foundation for further research and development in real-time tweet classification systems.

**Keywords:** Disaster Tweet Classification, Natural Language Processing, Machine Learning, Logistic Regression, Social Media Analysis.



# CHAPTER 1

## INTRODUCTION

### *1.1 Introduction to Disaster*

#### **About Disasters:**

A disaster is an event of great magnitude that interferes with life, society, and the ecosystem. Such large-scale events are often triggered by either natural causes or human activity and usually lead to numerous fatalities and displacements, destruction of valuable structures and networks, and environmental degradation. Studying and managing disasters may, therefore be challenging because they fall under two broad categories: causes and effects, so knowledge in most fields such as geology, weather, technology, engineering, and social sciences would be required.

Disasters, whether natural such as hurricanes and earthquakes or human-induced like industrial accidents, have devastating impacts on societies. The economic costs of disasters are rising globally, with annual losses exceeding \$300 billion in some years. Rapid urbanization and climate change further exacerbate these risks, emphasizing the need for real-time monitoring and decision-making systems.

#### **Types of Disasters:**

Natural disasters are powerful events caused by natural forces that can dramatically change our environment. For example, earthquakes happen when tectonic plates shift, releasing energy that creates seismic waves strong enough to topple buildings. Tsunamis, on the other hand, are huge waves generated by disturbances on the ocean floor, often triggered by earthquakes, and can flood coastal areas.

Flooding usually occurs when heavy rain causes rivers to overflow or when dams fail, disrupting water management systems and urban drainage. Hurricanes and cyclones are low-pressure systems that bring fierce winds and heavy rain, leading to storm surges that can severely damage coastal communities. To understand these weather events, scientists use advanced weather models to predict their effects.

Lastly, wildfires often occur during hot, dry seasons and can be made worse by strong winds, destroying forests, wildlife, and homes. Together, these different natural disasters show how various environmental factors interact and can lead to widespread destruction.



*Fig 1.1.1: Haiti earthquake damage*

### **Effects of Disasters on Our Lives:**

The effects of disasters on our lives are profound and multifaceted, impacting individuals and communities in various ways. **Physically**, disasters can inflict serious harm on people, leading to injuries and fatalities, while also causing the loss of homes and personal belongings, along with disrupting access to essential services like water and electricity. This physical devastation often extends to the **mental health** of survivors, who may grapple with trauma, anxiety, and depression; the emotional toll of losing loved ones or homes can result in long-term psychological disorders that affect their overall well-being. Economically, disasters can severely impact local economies, leading to business closures and job losses, while the financial burden of recovery places immense pressure on both government and community resources. Additionally, the **displacement** of communities can fracture social networks and support systems, making it a time-consuming and challenging process to rebuild these vital connections. Finally, the **environmental impact** of disasters can be devastating, causing pollution, habitat destruction, and alterations to ecosystems, which can have lasting effects on wildlife and natural resources. Together, these interconnected effects illustrate the comprehensive and enduring challenges that disasters impose on individuals, communities, and the environment. The importance of information during a disaster cannot be overstated, as it plays a crucial role in mitigating the effects outlined previously. Timely and accurate information can significantly enhance **physical safety**, allowing individuals to make informed decisions about evacuation and access to emergency services, thereby reducing injuries and fatalities. Furthermore, effective communication can address the **mental health** needs of survivors by providing resources and support systems, helping them cope with trauma and anxiety in the aftermath of a disaster.[1]

### ***1.2 Importance of Information during disaster***

Information is crucial during disasters as it connects people, resources, and responders, playing a key role in all phases of disaster management, including prevention, preparedness, response, and recovery. The effective dispersal of information can be a determinant of survival; prompt alerts enable early evacuation, saving lives, while accurate damage assessments facilitate efficient resource allocation. However, the spread of misinformation and rumors, as seen during the COVID-19 pandemic, underscores the necessity for credible sources of information.

Warnings are essential; early notifications about impending storms, earthquakes, or fires empower individuals to evacuate, seek shelter, and take necessary actions. Delayed or inaccurate warnings significantly increase the risk of injuries and fatalities, making it essential for victims to know safe locations, quick escape routes, and whom to contact for assistance. Developing situational awareness is also vital for effective response. Understanding the scale and frequency of a disaster enables relief efforts to be more focused. By identifying affected areas, such as those experiencing flooding or intensifying fires, situational data reveals which communities urgently need assistance—be it food, water, or medical care—and highlights the most critical support required at any given moment.

Monitoring changing conditions is important as disasters often evolve unpredictably; for instance, aftershocks may follow an earthquake, or disease outbreaks can occur after flooding. Continuous information flow is essential for adapting response strategies to meet these changing circumstances effectively. Additionally, helping responders work together is crucial for successful disaster response, which involves collaboration among various actors, including government agencies, NGOs, healthcare providers, and local communities. Access to accurate information enables these groups to coordinate efforts, avoiding duplication and ensuring that key needs are addressed. Knowing who is affected, where road blockages exist, and what resources are available is crucial for effective aid delivery.

In times of global crises, information fosters collaboration among nations and multinational entities, facilitating coordinated action. Fighting rumors and lies is another critical aspect; in today's social media landscape, misinformation can spread rapidly during emergencies, leading to confusion and panic. Reliable sources of information help counteract this by providing timely updates, allowing individuals to make informed decisions and minimizing the damage caused by false information. Improving resource utilization is also vital; accurate information aids relief organizations in pinpointing areas most in need, optimizing resource allocation. For example, knowledge of a remote village cut off by floods enables quicker and more efficient delivery of aid.

Finally, budgeting and planning benefit from effective information dissemination. Governments and NGOs can allocate funds based on real-time information regarding the disaster's scale and severity, ensuring that resources are directed where they are most needed. In summary, the interconnected effects of disasters—ranging from physical harm and mental health challenges to economic disruption, community disintegration, and environmental damage—underscore the critical need for effective information dissemination. By ensuring that individuals and communities have access to accurate information during a disaster, we can enhance resilience and facilitate recovery, ultimately reducing the long-term consequences of such catastrophic events.[2][3]

### *1.3 Social media and its Role during Disaster:*

Social media platforms have transformed disaster management by facilitating instant communication, enhancing situational awareness, and building community resilience. Platforms like Twitter, Facebook, and Instagram serve as essential tools for public communication and for authorities during emergencies. They enable the dissemination of real-time information, create spaces for victims to seek help, and facilitate large-scale coordination efforts in relief.

Research has demonstrated the effectiveness of Twitter during significant disasters. For example, during Hurricane Sandy and the Haitian earthquake, users frequently posted updates, sought assistance, and shared safety tips.

The **Disaster Tweet Analyser** project aims to harness the vast amount of unstructured data on Twitter to improve disaster response. The project focuses on several key advantages: real-time disaster detection, quick alerts and messages, and understanding relief operations. By developing a system that interacts with live Twitter data, the project aims to identify tweets about disasters almost in real-time, allowing for quicker identification of incidents as they occur. The system will create and send automated reports to authorities, groups, and the public with confirmed disaster information, clearly communicating necessary response actions. Analyzing tweet data can provide actionable insights regarding affected areas, resource shortages, and victim needs, aiding decision-making during disaster recovery and resource allocation.

Twitter has proven invaluable during disasters, such as the 2011 Tohoku earthquake, where it helped victims request aid and stay connected. In 2012, during Hurricane Sandy, tweets provided crucial updates and assisted in mapping affected areas. After the 2017 Mexico City earthquake, real-time tweets about collapsed buildings guided rescue efforts. During the COVID-19 pandemic, Twitter was key for sharing health updates and combating misinformation.

Opportunities for better Twitter use include the 2010 Haiti earthquake, where social media could have sped up rescue efforts. In the 2018 Kerala floods, more participation from rescue teams on Twitter could have improved crisis mapping. Similarly, the 2017 wildfires in Split, Croatia, could have benefited from better real-time analytics for evacuation planning.

These examples illustrate how Twitter can significantly aid in managing disasters. It serves as a vital communication tool during emergencies, but its effectiveness can be further enhanced by addressing challenges like misinformation and ensuring data accuracy.[3]



Fig 1.3.1: Tasmania's bushfires have become a prominent topic online

### *1.4 Objective of the Project*

#### **Incorporating Social Media Data into Disaster Management**

The project aims to integrate social media data into disaster management by developing a system that filters and analyzes tweets during emergencies. This system will connect real information with practical advice, making it easier for people to understand what to do in a crisis. By using advanced technologies like machine learning, we can ensure a more active and effective response when disasters occur.

Social media, especially Twitter, is crucial in managing disasters because it allows for the quick sharing and analysis of information. When disasters happen, a huge amount of real-time data is generated, providing valuable insights for decision-making and coordinating responses.

The **Disaster Tweet Analyzer** uses advanced technologies such as Natural Language Processing (NLP), Machine Learning, and geospatial analysis. It scans Twitter to find tweets related to disasters, helping to filter out irrelevant information. The system checks the accuracy and trustworthiness of this information, identifies the locations of the tweets, and assesses public sentiment regarding the disaster.

By sending real-time alerts, the Disaster Tweet Analyzer ensures that resources can be deployed quickly to where they are needed most. It enhances situational awareness by compiling data to evaluate the severity of a disaster and how the public is reacting. This system helps reduce the impact of disasters by providing early warnings, promoting collaboration among governments, non-governmental organizations (NGOs), and communities, ensuring everyone is informed and can work together effectively.

One of the key strengths of this system is its scalability, allowing it to monitor various types of disasters, such as floods and wildfires, across different locations. This makes it an essential tool for speeding up response times, improving the efficiency of resource use, and ultimately saving lives. By promoting the use of advanced technology in disaster management, the Disaster Tweet Analyzer aims to make communities safer and more resilient in the face of emergencies.[4][5]

## *1.5 Outline of the Project Report*

### Chapter 1: Introduction

This chapter begins with an introduction to disasters, emphasizing their devastating impact on human lives, infrastructure, and economies. It highlights the critical role of information during disasters in aiding timely decision-making, resource allocation, and rescue efforts. The chapter further explores the growing influence of social media platforms during such events, detailing how they serve as real-time communication tools for sharing updates, coordinating relief efforts, and disseminating crucial information to affected communities. It concludes by presenting the objective and significance of our project, focusing on its aim to leverage social media data for effective disaster management and bridging the gaps in existing approaches.

### Chapter 2: Overview of Previous Works

This chapter reviews the significant work done by researchers and organizations in the domain of disaster management using technological tools. It analyzes the problems addressed by these studies, such as real-time data collection, resource distribution, and information dissemination, while identifying their limitations and drawbacks. The chapter critically evaluates gaps like inefficient data processing, lack of accurate predictive models, or limited usability in real-time scenarios. Finally, it discusses how our project addresses these shortcomings by incorporating innovative methodologies, advanced models, and user-centric features, offering a more robust and comprehensive solution.

### Chapter 3: Data & Methodology

This chapter provides a detailed explanation of the data and methodologies employed in the project. A workflow diagram illustrates the step-by-step process, starting with data preprocessing to clean and organize raw data. Feature engineering techniques are applied to extract meaningful insights, followed by dataset modifications to ensure quality and relevance. The methodology elaborates on model training in week five, which focuses on developing a reliable predictive framework. Subsequent sections detail model testing, optimization, and user interface development in week six, culminating in week seven with the integration of the model into a web interface, fine-tuning its performance, and considering possible feature additions. Each phase is meticulously described to reflect the systematic approach adopted.

### Chapter 4: Results

This chapter presents the results derived from each step of the methodology, highlighting the outcomes achieved through preprocessing, feature engineering, and dataset modifications. The findings are discussed in the context of model training, where key performance metrics such as accuracy, precision, and recall are analyzed. The chapter also elaborates on the results obtained from testing and optimization, illustrating the effectiveness of the developed model.

Additionally, the integration of the user interface and web model is evaluated, showcasing the seamless interaction and functionality achieved.

## Chapter 5: Discussions, Conclusions, and Future Scope

This chapter begins with an introduction that sets the context for discussing the project outcomes. The discussion delves into the most significant findings, supported by statistical evidence such as percentages and numerical data, to highlight the project's impact and effectiveness. The chapter concludes the project by summarizing the key achievements and contributions, emphasizing its value in disaster management. In the future scope section, potential improvements are identified, such as enhancing model accuracy, incorporating live tweet feeds for real-time updates, and addressing cybersecurity concerns to safeguard sensitive data. Finally, the chapter highlights the broader implications of the project, envisioning its scalability and adaptability for future disaster scenarios.

## Appendix

The appendix contains supplementary materials, including workflow diagrams, detailed data descriptions, and additional findings that support the project.

## References

This section provides a comprehensive list of all academic papers, articles, and resources cited throughout the document, adhering to standard referencing guidelines.

## **CHAPTER 2**

### **Analysis of Existing Solutions in Disaster Tweet Classification**

#### ***2.1 Introduction to the Chapter***

This chapter explores the progression of disaster tweet classification, analyzing existing systems, their strengths, and limitations. It highlights the challenges these systems address and discusses how the Disaster Tweet Analyzer builds on these foundations to provide a more effective, comprehensive solution for disaster management.

#### ***2.2 Evolution of Disaster Tweet Analysis and Related Work***

Over the last decade, platforms like Twitter have become crucial during disasters for real-time information sharing. Researchers have extensively explored the use of Natural Language Processing (NLP) and Machine Learning (ML) for classifying tweets and extracting actionable insights. Notable systems in this domain include:

The field of disaster management has seen significant advancements through the integration of Natural Language Processing (NLP) and social media data. For instance, CrisisLex, a repository of annotated crisis-related tweets, provided early insights into how NLP can be leveraged for disaster scenarios by organizing and analyzing vast amounts of crisis-related information. Building on this foundation, AIDR (Artificial Intelligence for Disaster Response)[7] demonstrated real-time categorization of tweets, enabling immediate and efficient responses during emergencies. Similarly, the SAHANA Project[10], an open-source disaster management platform, enhanced operational effectiveness by integrating social media data, showcasing how real-time information could support coordination efforts. Meanwhile, the Harvard Humanitarian Initiative utilized NLP tools to assess humanitarian needs by analyzing text data from social media, highlighting the critical role of data-driven insights in shaping relief efforts. Together, these initiatives illustrate the transformative potential of combining NLP with social media to improve disaster response and management.

These systems demonstrated the potential of automated tweet analysis to support government agencies, NGOs, and responders in effectively managing crises.

Existing systems addressed critical challenges in disaster response. Timely classification of tweets was pivotal, with systems like AIDR processing data in near real-time to prioritize actions swiftly [8][9]. The ability to handle large-scale data during disasters was another key focus; tools such as CrisisLex and TREC-IS scaled effectively to analyze thousands of tweets per second [6][7]. Additionally, categorizing disaster information into actionable themes, such as damage reports or donation requests, helped responders allocate resources more efficiently [8]. Sentiment analysis provided insights into public morale, though early models often



struggled with the informal and context-specific language [11]. Furthermore, initiatives like CrisisLex contributed annotated datasets, establishing benchmarks for ML training in disaster tweet analysis [6][10].

### *2.3 Gaps and Limitations in Existing Systems*

Despite advancements, existing systems exhibited several limitations. Most tools classified tweets as disaster-relevant or not, offering limited contextual insights, such as the disaster type, location, or detailed sentiment [7][9]. Sentiment analysis models often failed to handle informal and ambiguous tweet language, impacting their accuracy [8][11]. Extracting geolocation information remained challenging due to the inconsistent use of geotags and vague textual references [9]. Moreover, many systems lacked user-friendly interfaces, making them inaccessible to non-technical users, such as field responders [10]. Some systems also struggled with real-time processing during high data surges, reducing their efficiency during critical moments [8].

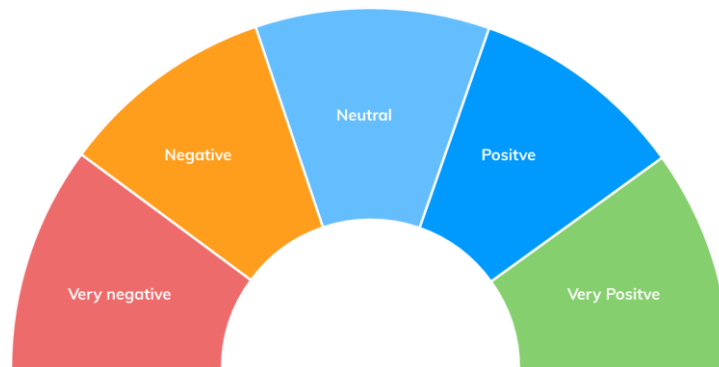
These limitations underscore the need for a more accessible, real-time, and comprehensive approach to disaster tweet analysis.

### *5.4 How Our Project Bridges the Gaps*

The Disaster Tweet Analyzer addresses these gaps while building upon the strengths of existing systems. Unlike earlier solutions that primarily identified tweet relevance, this project provides deeper insights. It extracts location information using advanced Named Entity Recognition (NER), identifies disaster types (e.g., floods, earthquakes, fires), and conducts detailed sentiment analysis. This enriched understanding empowers responders to make precise decisions and allocate resources effectively [9][10].

To overcome the limitations of sentiment analysis in existing systems, our project trains NLP models on domain-specific datasets, ensuring improved interpretation of informal and context-specific tweet language [11]. Additionally, the system advances geolocation capabilities by analyzing tweet content beyond geotags, yielding more accurate and consistent location data [9].

A user-friendly, web-based interface ensures accessibility for both technical and non-technical users, bridging a key gap observed in many existing tools [10]. The backend infrastructure is optimized for real-time performance and scalability, enabling the processing of high tweet volumes during peak disaster activity [8]. By combining disaster classification, location detection, disaster type identification, and sentiment analysis, the Disaster Tweet Analyzer



*Fig 2.4.1: Sentiment Analysis*

provides a holistic solution for disaster management. This approach not only addresses the limitations of earlier systems but also sets a new benchmark for leveraging social media in disaster response [9][11].

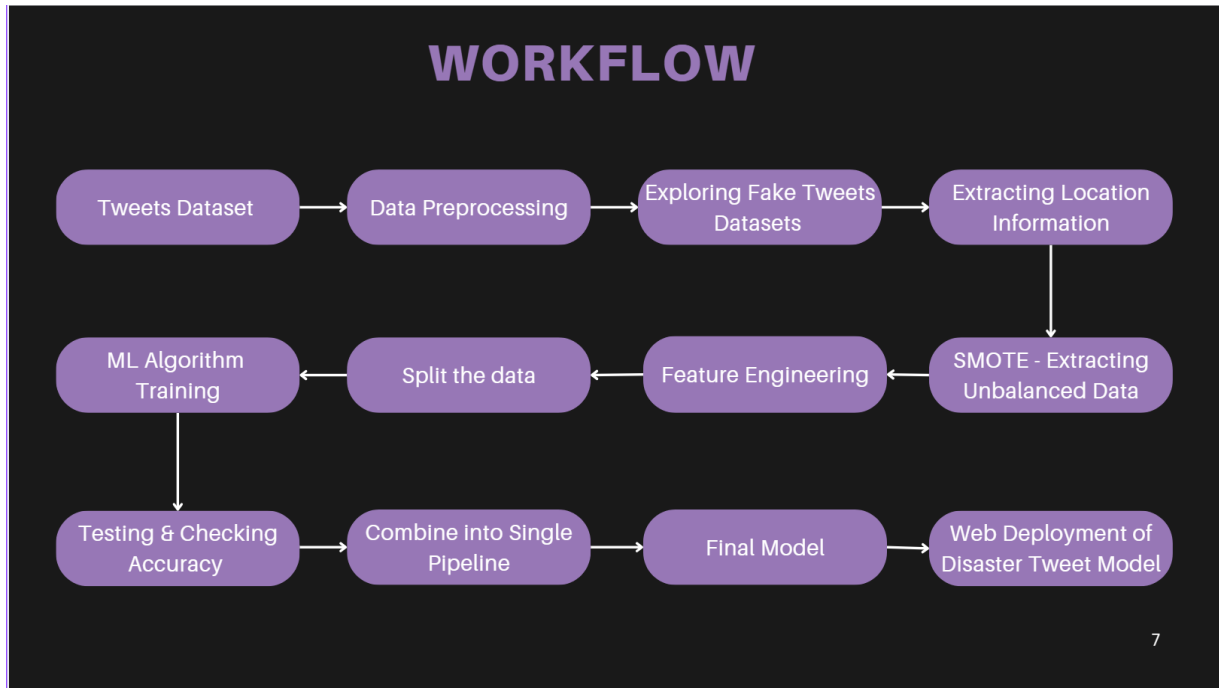
## CHAPTER 3

### DATA & METHODOLOGY

#### *3.1 Introduction to this Chapter*

The dataset used in this project consists of tweets related to disasters, specifically designed to classify tweets as either disaster-related or non-disaster-related. It contains over 11,000 tweets with various features such as tweet content, keyword, location (if available), and a binary target variable indicating whether the tweet is related to a disaster (1) or not (0). The tweets are categorized based on whether they contain information about natural disasters, accidents, or other emergency situations, providing valuable insights for monitoring and responding to real-time disasters. The methodology for the disaster tweet classification involves several key steps. Initially, data preprocessing is performed to clean the dataset by removing any irrelevant information, such as unnecessary characters or stop words, and by handling missing values. Feature extraction techniques, including tokenization, vectorization (using methods like TF-IDF), and text normalization, are applied to convert text data into numerical representations that can be fed into machine learning models. Following data preparation, the dataset is split into training and testing sets to ensure proper evaluation.

Various machine learning models, such as Logistic Regression, Support Vector Machines (SVM), Gradient Boosting, and K-Nearest Neighbors (KNN), are trained on the dataset to identify patterns that distinguish disaster-related tweets from non-disaster ones. Model evaluation is performed using metrics like accuracy, precision, recall, F1-score, and confusion matrices to assess the performance of each model. Precision and recall are particularly important in this context due to the imbalanced nature of the dataset, where non-disaster tweets may far outnumber disaster-related tweets. The web interface component is designed to predict disaster-related tweets in real-time, providing the user with predictions on the sentiment, category, location, and relevance of each tweet, ultimately helping in disaster response efforts. This approach aims to leverage the power of machine learning for faster, more accurate disaster classification from social media data.



*Fig 3.1.1: Workflow*

### 3.2 About Our Data

The dataset used for this study is the Kaggle Disaster Tweets dataset, which comprises 11,370 tweets labeled as disaster-related or non-disaster-related. Each tweet in the dataset contains multiple features, including the tweet's **text**, **keyword**, **location** (if available), and **target**. The **target** column is a binary variable, where 1 indicates a disaster-related tweet and 0 indicates a non-disaster-related tweet. The **text** field provides the content of the tweet, and the **keyword** field highlights a word or phrase extracted from the tweet that is most indicative of its content.

For this study, the goal is to predict whether a tweet is related to a disaster or not based on its textual content. Additionally, the model also aims to classify the **location** of the tweet, assign a **category** based on the nature of the disaster, and determine the **sentiment** of the tweet (e.g., positive, negative, or neutral). These auxiliary tasks help in creating a comprehensive view of the information shared during disaster events, which could be useful for real-time analysis and decision-making.

### *3.3 Preprocessing*

Effective data preprocessing is crucial for building a robust disaster tweet classification model. The process begins with data cleaning to handle missing values, duplicates, and outliers. For columns such as keyword and location, missing values are imputed to maintain dataset consistency. Missing keywords are replaced with the most frequent keyword or labelled as "unknown," while missing locations are also assigned "unknown" to standardize data representation. Removing duplicate rows based on the text column prevents model bias and redundancy. Additionally, outliers are addressed by analysing tweet lengths. Extremely short or meaningless tweets, such as those containing only a single character or symbol, are removed, ensuring the dataset retains meaningful information for classification.

Text preprocessing is a vital step in transforming the raw textual data into a format suitable for machine learning models. First, case normalization is applied, where all text in the text column is converted to lowercase. This ensures uniformity and eliminates discrepancies caused by case sensitivity. Next, noise is removed from the tweets. URLs, HTML tags, emojis, and special characters are stripped using regular expressions or specialized libraries like emoji. Excessive whitespace is also removed to tidy the data further. Stop words, such as "the," "and" and "is," which provide minimal contextual value, are eliminated using natural language processing (NLP) libraries like NLTK or SciPy. This step enhances the focus on the core content of the tweets.

To prepare the text for analysis, lemmatization or stemming is performed. Lemmatization reduces words to their base form, ensuring consistency across variations of the same word. For instance, "running" is converted to "run," and "better" becomes "good." Stemming, though less precise, can also be used to achieve similar results by reducing words to their root form. Tokenization follows, splitting each tweet into individual words or tokens, enabling the application of NLP techniques and feature extraction. This prepares the dataset for embedding or encoding strategies.

Categorical encoding is applied to the keyword and location columns to make them model-compatible. For the keyword column, one-hot encoding creates binary features for each unique keyword, while label encoding assigns numerical labels. Alternatively, grouping similar keywords, such as "fire," "burn," and "blaze" into categories like "fire-related," simplifies the data and reduces feature sparsity. For the location column, encoding depends on its relevance. Frequently occurring locations can be encoded individually, while less common locations are

grouped under "others." Broad categorization into "urban," "rural," or "unknown" simplifies the feature further, aligning it with the disaster context [12].

This comprehensive preprocessing workflow ensures the dataset is clean, consistent, and well-structured for machine learning. It addresses missing values, duplicates, and outliers, while the text preprocessing steps refine the tweet content for analysis. Encoding the categorical variables enhances compatibility and introduces meaningful patterns for model interpretation. By implementing these steps systematically, the processed dataset becomes an ideal foundation for training a disaster tweet classification model. This ensures the model's predictions are accurate, reliable, and capable of distinguishing between disaster and non-disaster tweets effectively.

## Coding

```
import re
import nltk
from nltk.corpus import stopwords

def clean_text(text):
    text = re.sub(r'https?://[^\s]+', '', text) # Remove links
    text = re.sub(r'^\x00-\x7F+', '', text) # Remove emojis
    text = re.sub(r'^a-zA-Z0-9\s', '', text) # Remove special
characters
    text = re.sub(r'\.{3,}$', '', text) # Remove trailing ellipsis
    text = text.lower() # Convert to lowercase
    # Remove stopwords
    stop_words = set(stopwords.words('english'))
    tokens = nltk.word_tokenize(text)
    tokens = [t for t in tokens if t not in stop_words]
    text = ' '.join(tokens)
    return text

#download the stopwords and punkt
import nltk
import spacy
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
import re

# Download required NLTK resources
nltk.download('punkt')
nltk.download('stopwords')
```

```

nltk.download('punkt_tab')
nltk.download('wordnet')
import nltk
nltk.download('stopwords')
import nltk
nltk.download('punkt')
import nltk
nltk.download('punkt')
Twitter_Data['text'] = Twitter_Data['text'].apply(clean_text)
print(Twitter_Data['text'].head(20))
#location needs to be converted into string for this
Twitter_Data['location'] = Twitter_Data['location'].astype(str)
def clean_location(text):
    text = re.sub(r'^a-zA-Z0-9\s,.-]', '', text) # Remove special
characters
    text = text.title() # Convert to title case
    text = re.sub(r'\b(city|state|country)\b', '', text,
flags=re.IGNORECASE) # Remove unnecessary words
    text = re.sub(r'\s+', ' ', text) # Remove excess whitespace
    return text
Twitter_Data['location'] =
Twitter_Data['location'].apply(clean_location)

import re#A regular expression (shortened as regex or regexp)
# Define regex pattern for location extraction

location_patterns = [
    r'\b(New|North|South|East|West) \w+\b', # Directional cities
    r'\b\w+ville\b', #Cities ending in "ville"
    r'\b\w+ city\b', #Cities with "city" suffix
    r'\b\w+ state\b', # States
    r'\b(NY|CA|FL|TX|IL|PA|OH|GA|NC|MI)\b', #Common state abbreviations
    r'\b(USA|United States|America)\b', #Country names
]

def extract_location_regex(text):
    for pattern in location_patterns:
        match = re.search(pattern, text, re.IGNORECASE)#re.IGNORECASE
flag for case-insensitive matching.
        if match:
            return match.group()
    return None

# Apply regex function to dataset (fill null values in 'location')

```

*Listing 3.3.1: Preprocessing and Location Extraction*

### ***3.4 Feature Engineering***

Feature engineering is a critical step in enhancing the predictive power of machine learning models. For disaster tweet classification, extracting both text-based and location-based features, along with advanced NLP techniques, ensures a comprehensive understanding of the data.

#### ***3.4.1 Text Based Features***

Text-based features begin with analysing the structure and content of tweets. Tweet length is a simple yet effective feature, capturing the number of characters and words in each tweet. These metrics can reveal patterns, such as whether shorter tweets are more likely to convey urgent information. Additionally, hashtags and mentions in tweets play a significant role. By counting the occurrences of hashtags (#) and mentions (@), distinct features are created that highlight social engagement and topic-specific references.

The presence of keywords in tweets adds further context. A binary indicator captures whether a keyword from the keyword column appears in the text field. Sentiment analysis provides deeper insights into the emotional tone of tweets. Using sentiment analysis tools, polarity scores (positive, negative, or neutral) and subjectivity scores (level of personal opinion versus factual content) are extracted. These scores help differentiate emotional appeals from neutral disaster updates. For numerical representation of tweet text, Bag of Words (Bow) and Term Frequency-Inverse Document Frequency (TF-IDF) are widely used. These methods transform cleaned text into vectors, capturing word frequency and importance. Advanced representation techniques like pre-trained word embeddings, such as Glove, Word2Vec, or BERT, provide context-aware vectors, capturing semantic relationships between words [13].

#### ***3.4.2 Location Based Features***

Location information in disaster tweets often holds critical value. Region extraction involves grouping locations into broader categories like countries, states, or continents to generalize patterns. Frequency encoding captures how often a location appears in the dataset, creating a feature that quantifies the prominence of certain areas. A binary location presence feature can be introduced to flag tweets originating from disaster-prone areas, such as "California" for wildfires or "Florida" for hurricanes. Applying text preprocessing to the location column, including normalization of names (e.g., converting "NYC" to "New York"), ensures



consistency and enables meaningful pattern extraction. By treating location names as textual data, valuable insights are gained from analysing the frequency or sentiment associated with different regions.

### 3.4.3 Interaction Features

Interaction features are particularly useful in combining information across columns. For example, keyword-location interactions capture how often certain keywords co-occur with specific locations. This highlights location-specific events, such as keywords like "earthquake" being tied to regions like "Japan" or "California." Another useful feature is keyword-text match, a binary indicator that flags tweets where the keyword appears directly in the text. This provides a measure of alignment between the structured keyword column and the unstructured text content.

### 3.4.4 Advanced NLP Techniques

Incorporating advanced NLP techniques like Named Entity Recognition (NER) enhances feature richness. NER identifies entities such as disaster names, places, organizations, or times, providing structured insights from unstructured text. Topic modelling using methods like Latent Dirichlet Allocation (LDA) uncovers hidden topics within the tweets, clustering them into thematic groups. For example, LDA might reveal topics related to wildfires, floods, or rescue operations, adding another dimension to the model's understanding of tweet content.

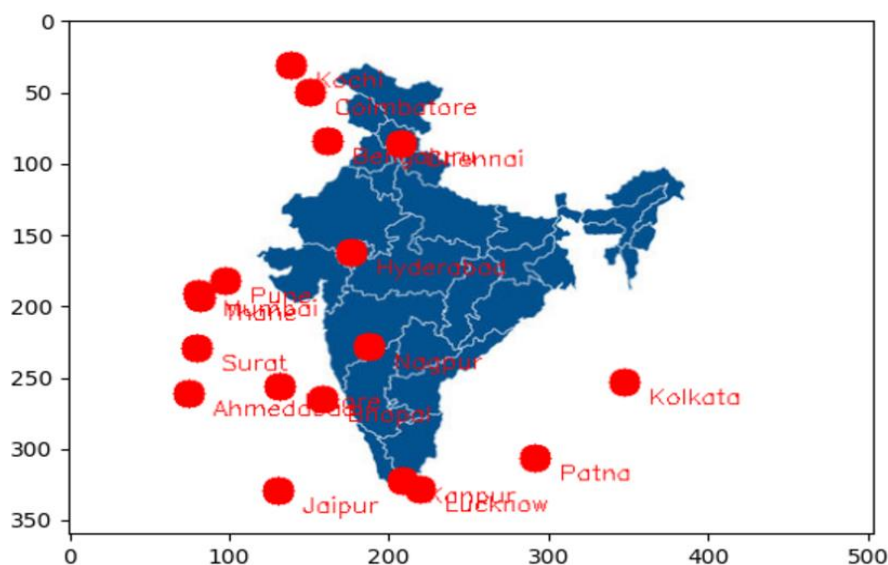
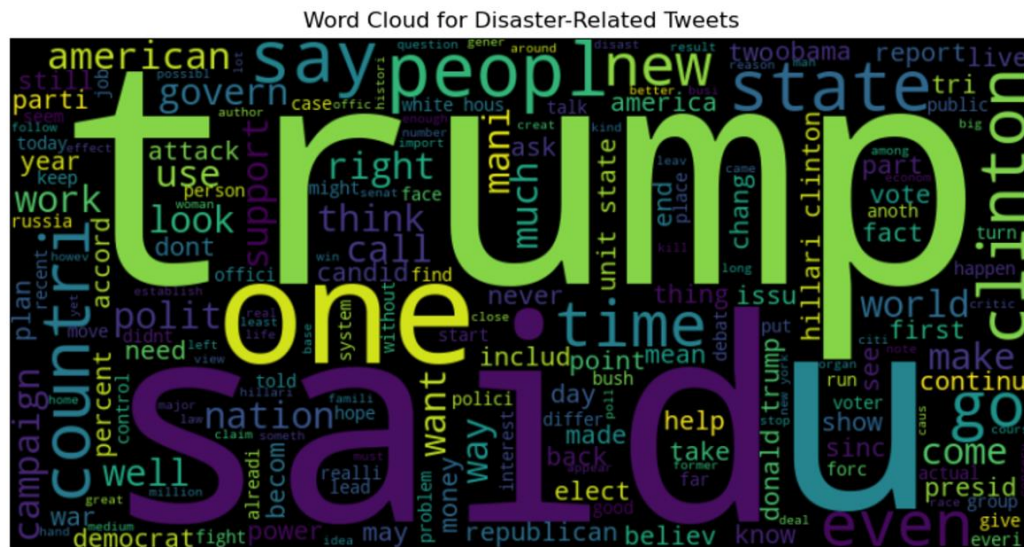


Fig 3.4.4.1: Disaster Location Extraction form Tweets



*Fig 3.4.2: Word Cloud of Disaster Related Tweets*

### 3.5 Dataset Modifications and Creation

The modification begins by taking the original tweets.csv dataset and applying necessary preprocessing steps, such as handling missing values, cleaning the text, and feature extraction, as described previously. Once the dataset is cleaned and pre-processed, we create a new column called **target2**. The target2 column is introduced to distinguish between **fake** and **real** tweets, labelled as **1** for real tweets and **0** for fake tweets. This classification could be based on a manual annotation or any other source of information that can help separate real disaster tweets from fake ones.

However, upon inspecting the balance of the dataset, it becomes clear that the **target2** variable is highly unbalanced, meaning there are far more real tweets (1) than fake tweets (0). This imbalance can severely affect model performance, as the model would be biased toward predicting the majority class. Given this issue, it was decided to **remove the target2 column** from the dataset. This decision helps avoid any potential bias in training the model and keeps the focus on the original target column, which is already labelled as relevant for disaster-related tweets.

## Final Dataset Structure

After removing the target2 column, the final dataset now contains:

- The original target column, which is used for the disaster classification task (disaster-related vs. non-disaster-related tweets).
- Other features such as text, location, keyword, and additional engineered features (e.g., tweet length, sentiment scores, hashtag count) that are relevant for model training.

This modification and cleanup process ensure that the dataset is properly balanced and ready for training a model to classify disaster-related tweets without the added complexity of an unbalanced fake vs. real classification task. The focus remains on classifying whether a tweet pertains to a disaster or not, which is the core objective of the project.

### *3.6 Model Training*

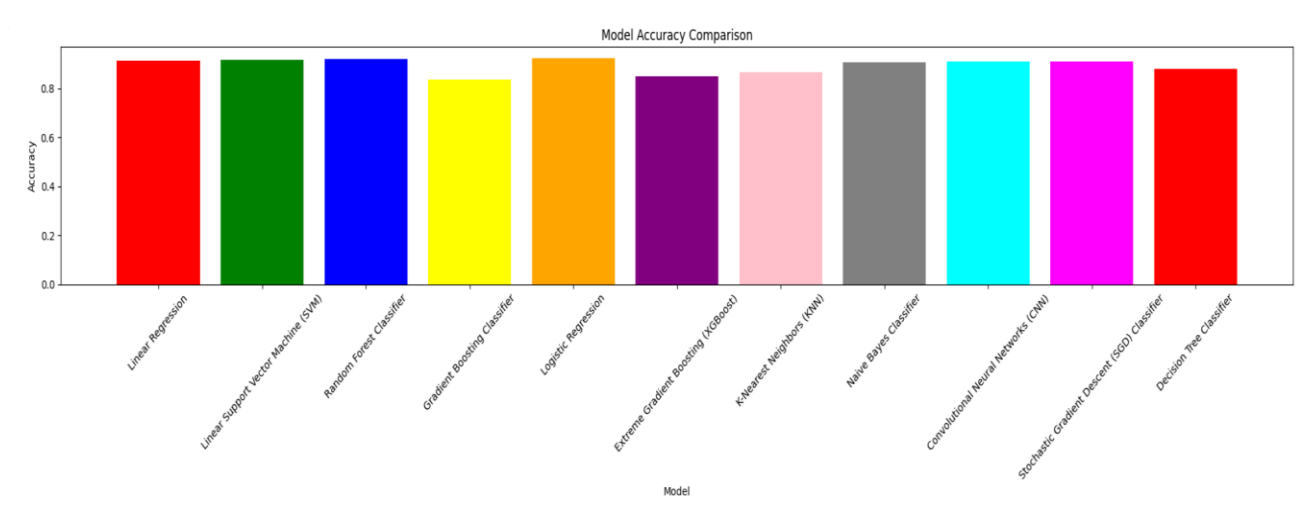
The Disaster Tweet Classification project aims to build a machine learning model capable of distinguishing between disaster-related and non-disaster-related tweets. This involves a systematic workflow starting from data preparation, feature engineering, model selection, and evaluation. Initially, the dataset undergoes preprocessing to handle raw, unstructured text and address missing or imbalanced values, ensuring the features are well-formatted for analysis. The dataset is then divided into training and testing sets using an 80-20 split, with the training data used to develop the model and the test set reserved for performance evaluation. Feature scaling is applied using `StandardScaler` to standardize the data by removing the mean and scaling to unit variance, which ensures all features contribute equally to the model. Sparse datasets are handled using the `with mean=False` parameter. Additionally, label encoding is performed to convert categorical target labels, such as "Disaster" and "Non-disaster," into binary numeric values for compatibility with machine learning algorithms.

Several models were explored to identify the best fit for this binary classification problem. Logistic Regression emerged as the final choice due to its simplicity, interpretability, and efficiency. It calculates the probability of a tweet belonging to a particular class based on the features and works well when the relationship between features and the target is linear. Other models considered include Linear Regression, which can be adapted for classification but lacks probabilistic output, and Linear Support Vector Machine (SVM), which excels in high-dimensional data but assumes a linear decision boundary. Random Forest and Gradient Boosting classifiers were evaluated for their ability to handle non-linear relationships and

imbalanced datasets, with Random Forest being robust but computationally expensive and Gradient Boosting offering high accuracy at the cost of extensive tuning. XGBoost, a gradient boosting framework, demonstrated excellent speed and predictive performance but required significant computational resources[14].

Models like K-Nearest Neighbours (KNN), Naive Bayes, and Decision Trees were also assessed. KNN, while straightforward, struggles with high-dimensional text data and scalability. Naive Bayes, a probabilistic model based on Bayes' Theorem, performed well on smaller datasets but fell short compared to more advanced models due to its assumption of feature independence. Decision Trees, while intuitive and effective for non-linear relationships, are prone to overfitting without regularization. Advanced techniques like Convolutional Neural Networks (CNN) were explored for their capability to identify complex patterns, but they required large datasets and computational power, making them impractical for this project. Stochastic Gradient Descent (SGD) and non-linear SVMs were considered for their scalability and flexibility but demanded careful hyperparameter tuning.

The selected Logistic Regression model was trained using scaled features and optimized to minimize prediction errors. Evaluation was conducted using metrics like accuracy, precision, recall, and the F1 score, ensuring a balanced assessment. A confusion matrix provided insights into true positives, true negatives, false positives, and false negatives, highlighting areas for



*Fig 3.6.1: Model Accuracy Comparison Chart*

improvement. Ultimately, the project successfully implemented a scalable and interpretable solution for classifying disaster-related tweets, with Logistic Regression proving to be the most suitable model.

## Coding

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X_resampled,
y_resampled, test_size=0.2, random_state=42)

from sklearn.preprocessing import StandardScaler
st = StandardScaler(with_mean=False)
X_train_scaled = st.fit_transform(X_train)
X_test_scaled = st.transform(X_test)

#FOR TRAINING WE REQUIRE THEM AS A BINARY VALUE
y_train = y_train.map({'Disaster': 1, 'Non disaster': 0})
y_test = y_test.map({'Disaster': 1, 'Non disaster': 0})

# LinearRegression

from sklearn.linear_model import LinearRegression
#Initialize Linear Regression model
lr_model = LinearRegression()
# Train Linear Regression model
lr_model.fit(X_train_scaled, y_train)
# Make predictions
y_pred_lr = lr_model.predict(X_test_scaled)
# Convert predictions to binary class labels (0 or 1)
y_pred_lr_binary = (y_pred_lr >= 0.5).astype(int)
# Evaluate model
from sklearn.metrics import accuracy_score, classification_report,
roc_auc_score
print("Linear Regression Accuracy:", accuracy_score(y_test,
y_pred_lr_binary))
print("Linear Regression Classification Report:")
print(classification_report(y_test, y_pred_lr_binary))
print("Linear Regression ROC-AUC:", roc_auc_score(y_test, y_pred_lr))
```

*Listing 3.6.1: Linear Regression Model*

```

# Random Forest
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report,
roc_auc_score
# Initialize Random Forest Classifier
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
# Train Random Forest Classifier
rf_model.fit(X_train_scaled, y_train)
# Make predictions
y_pred_rf = rf_model.predict(X_test_scaled)
# Evaluate model
print("Random Forest Accuracy:", accuracy_score(y_test, y_pred_rf))
print("Random Forest Classification Report:")
print(classification_report(y_test, y_pred_rf))
print("Random Forest ROC-AUC:", roc_auc_score(y_test,
rf_model.predict_proba(X_test_scaled)[:, 1]))

```

*Listing 3.6.2: Random Forest Model*

```

# GradientBoosting

from sklearn.ensemble import GradientBoostingClassifier
gb_model = GradientBoostingClassifier(
    n_estimators=50,
    learning_rate=0.1,
    max_depth=5,
    min_samples_split=2,
    random_state=42
)
# Train Gradient Boosting Classifier with parallel processing
def train_gb_model(X_train, y_train):
    return gb_model.fit(X_train, y_train)

from joblib import parallel_backend
with parallel_backend('multiprocessing'):
    gb_model.fit(X_train_scaled, y_train)

# Make predictions
y_pred_gb = gb_model.predict(X_test_scaled)

# Evaluate model
print("Gradient Boosting Accuracy:", accuracy_score(y_test, y_pred_gb))
print("Gradient Boosting Classification Report:")
print(classification_report(y_test, y_pred_gb))
print("Gradient Boosting ROC-AUC:", roc_auc_score(y_test,
gb_model.predict_proba(X_test_scaled)[:, 1]))

```

*Listing 3.6.3: Gradient Boosting Model*

```

# Logistic Regression
from sklearn.linear_model import LogisticRegression
lr_model = LogisticRegression(max_iter=1000)
lr_model.fit(X_train_scaled, y_train)

# Make predictions
y_pred_lr = lr_model.predict(X_test_scaled)
# Evaluate model
from sklearn.metrics import accuracy_score, classification_report,
roc_auc_score

print("Logistic Regression Accuracy:", accuracy_score(y_test,
y_pred_lr))
print("Logistic Regression Classification Report:")
print(classification_report(y_test, y_pred_lr))
print("Logistic Regression ROC-AUC:", roc_auc_score(y_test,
lr_model.predict_proba(X_test_scaled)[:, 1]))

```

*Listing 3.6.4: Logistic Regression Model*

```

# xgboost
import xgboost as xgb
xgb_model = xgb.XGBClassifier(objective='binary:logistic',
eval_metric='auc', max_depth=5, learning_rate=0.1, n_estimators=100)
xgb_model.fit(X_train_scaled, y_train)

# Make predictions
y_pred_xgb = xgb_model.predict(X_test_scaled)
# Evaluate model
from sklearn.metrics import accuracy_score, classification_report,
roc_auc_score

print("XGBoost Accuracy:", accuracy_score(y_test, y_pred_xgb))
print("XGBoost Classification Report:")
print(classification_report(y_test, y_pred_xgb))
print("XGBoost ROC-AUC:", roc_auc_score(y_test,
xgb_model.predict_proba(X_test_scaled)[:, 1]))

# KNN
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, classification_report,
roc_auc_score
# Initialize KNN model
knn_model = KNeighborsClassifier(n_neighbors=5)

```

*Listing 3.6.5: XG-Boost Model*

```

# Train KNN model
knn_model.fit(X_train_scaled, y_train)

# Make predictions
y_pred_knn = knn_model.predict(X_test_scaled)
# Evaluate model
print("KNN Accuracy:", accuracy_score(y_test, y_pred_knn))
print("KNN Classification Report:")
print(classification_report(y_test, y_pred_knn))
print("KNN ROC-AUC:", roc_auc_score(y_test,
knn_model.predict_proba(X_test_scaled)[: , 1]))

# Naive Bayes
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import accuracy_score
# Initialize and train NB model
nb_model = MultinomialNB()
nb_model.fit(X_train_scaled, y_train)
# Make predictions
y_pred_nb = nb_model.predict(X_test_scaled)
# Calculate accuracy
nb_accuracy = accuracy_score(y_test, y_pred_nb)

print("Naive Bayes Accuracy:", nb_accuracy)

```

*Listing 3.6.6: KNN and Naive Bayes Model*

```

#Convolutional Neural Networks (CNN)
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv1D, MaxPooling1D, Flatten,
Dense
from sklearn.metrics import accuracy_score

cnn_model = Sequential()
cnn_model.add(Conv1D(32, kernel_size=3, activation='relu',
input_shape=(X_train_scaled.shape[1], 1)))
cnn_model.add(MaxPooling1D(pool_size=2))
cnn_model.add(Flatten())
cnn_model.add(Dense(64, activation='relu'))
cnn_model.add(Dense(1, activation='sigmoid'))

cnn_model.compile(loss='binary_crossentropy', optimizer='adam',
metrics=['accuracy'])
cnn_model.fit(X_train_scaled, y_train, epochs=10, batch_size=32)

y_pred_cnn = cnn_model.predict(X_test_scaled)

```



```

y_pred_cnn = (y_pred_cnn > 0.5).astype('int32')

cnn_accuracy = accuracy_score(y_test, y_pred_cnn)
print("CNN Accuracy:", cnn_accuracy)

#Stochastic Gradient Descent
from sklearn.linear_model import SGDClassifier
from sklearn.metrics import accuracy_score

sgd_model = SGDClassifier(loss='hinge', penalty='l2')
sgd_model.fit(X_train_scaled, y_train)
y_pred_sgd = sgd_model.predict(X_test_scaled)
sgd_accuracy = accuracy_score(y_test, y_pred_sgd)
print("SGD Accuracy:", sgd_accuracy)

# Decision Tree
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score

dt_model = DecisionTreeClassifier()
dt_model.fit(X_train_scaled, y_train)
y_pred_dt = dt_model.predict(X_test_scaled)
dt_accuracy = accuracy_score(y_test, y_pred_dt)
print("Decision Tree Accuracy:", dt_accuracy)

```

*Listing 3.6.7: CNN, Gradient Descent & Decision Tree Model*

### 3.7 Model Web-Interface

For integrating a model into a web interface, Flask offers more customization, while Streamlit allows for faster deployment. Fine-tuning involves exploring advanced models like SVM, Random Forest, or deep learning and using techniques like hyperparameter tuning and cross-validation. Enhancing features with sentiment analysis, location, user profiles, and ensemble methods can improve accuracy.

#### 3.7.1 Web Interface Integration

For integrating a model into a web interface, you can choose between Flask and Streamlit, depending on your needs. Flask is a great choice for building more customizable and scalable web applications, as it allows for greater control over the backend logic and model integration. You can create custom routes, handle user input, and manage data flow between the front-end and the trained model. On the other hand, Streamlit is better suited for rapid prototyping of machine learning applications, offering a simpler way to deploy models with minimal coding. It provides built-in widgets to capture user input, such as text boxes for

tweet submission, and automatically renders the predictions. The web interface should allow users to enter tweet text, and upon submission, the tweet is sent to the backend model through an API call or direct model invocation. The model processes the tweet and classifies it as disaster-related or not, and the result is displayed on the webpage. Streamlit makes this integration seamless with built-in support for displaying output directly on the interface, while Flask would require you to handle rendering the output using HTML templates and routes. Both frameworks can effectively serve the purpose, but Streamlit offers an easier and faster setup for ML-based web applications.

**Steps:**

- Develop an API endpoint to accept tweet text.
- Call the model from the backend to predict the class (real/fake).
- Display the result on the webpage.

## CHAPTER 4

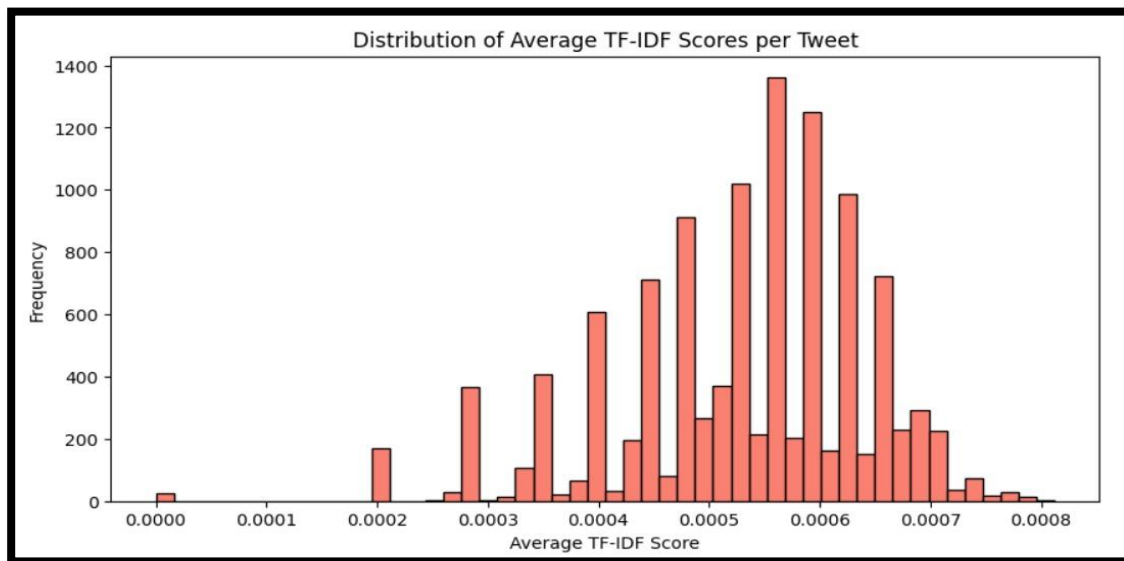
### RESULTS AND DISCUSSION

#### *4.1 Machine Learning Workflow and Visualization*

The machine learning process involved several stages, from data preprocessing to model evaluation. Preprocessing steps included text tokenization, vectorization using TF-IDF, and handling imbalanced data. These steps were visualized to better understand data distributions and model performance.

During training, confusion matrices and classification reports were generated to assess each model's performance. Precision, recall, F1-score, and accuracy metrics were plotted for visual comparison, providing clear insights into the strengths and weaknesses of each algorithm. Additionally, feature importance graphs were generated for models like Gradient Boosting, highlighting the most influential features in disaster tweet classification.

These visualizations are crucial not only for model evaluation but also for communicating the technical robustness of our approach to stakeholders. They illustrate the iterative process of refining models to achieve optimal results, demonstrating the analytical rigor underpinning the Disaster Tweet Analyzer.



*Fig 4.1.1: TF-IDF Score Bar-graph*

## 4.2 Performance of Learning Models

To evaluate the effectiveness of the Disaster Tweet Analyzer, we employed multiple machine learning models, including Logistic Regression, Linear Support Vector Machine (SVM), K-Nearest Neighbour (KNN), and Gradient Boosting, to classify tweets as disaster-relevant or not. Each model's performance was assessed using its confusion matrix, which provides insights into true positives, true negatives, false positives, and false negatives.

Logistic Regression demonstrated [91.17% and 88%] precision and recall, performing well in distinguishing disaster-relevant tweets. Linear SVM, known for its efficiency in high-dimensional spaces, showcased [89% and 89%] accuracy, which was marginally higher/lower than Logistic Regression. On the other hand, K-Nearest Neighbours, being a distance-based algorithm, performed with [85% and 84%], indicating its reliance on well-defined feature distributions. Lastly, Gradient Boosting, an ensemble learning method, achieved [70% and 92%], proving effective in handling the class imbalance in our dataset.

The confusion matrices of these models reveal critical patterns in prediction accuracy, underscoring their strengths and limitations. For our project, these results are invaluable as they guide the selection of the most reliable model for deployment, ensuring that disaster-relevant tweets are identified accurately, even in imbalanced datasets.

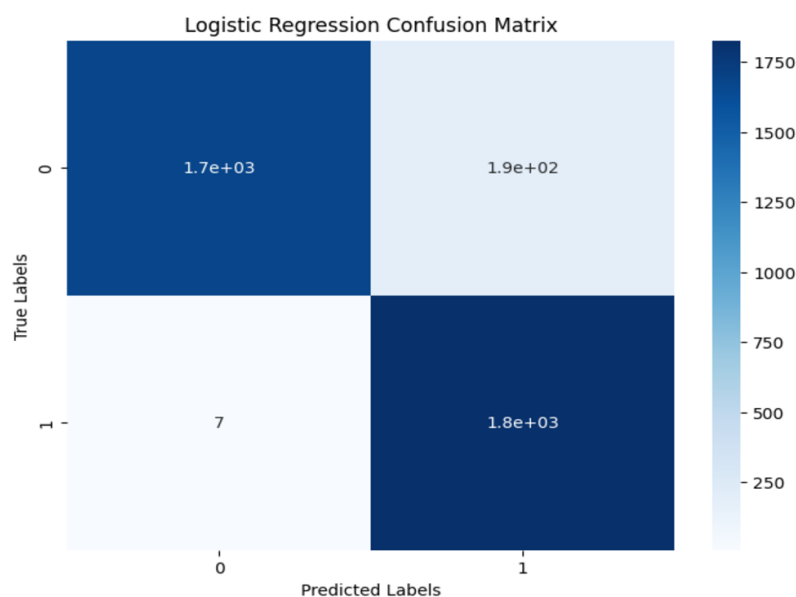
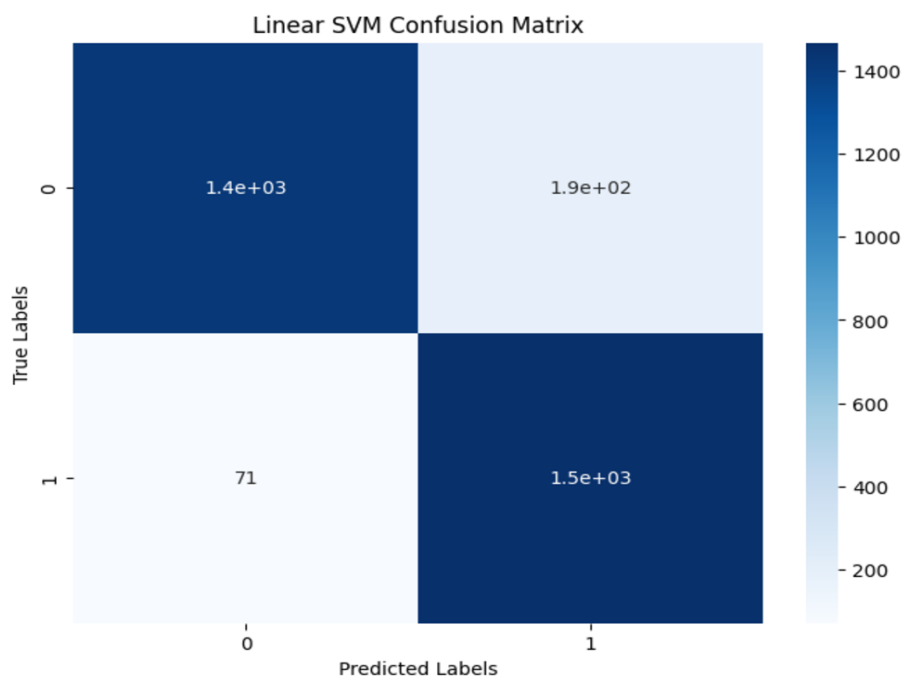
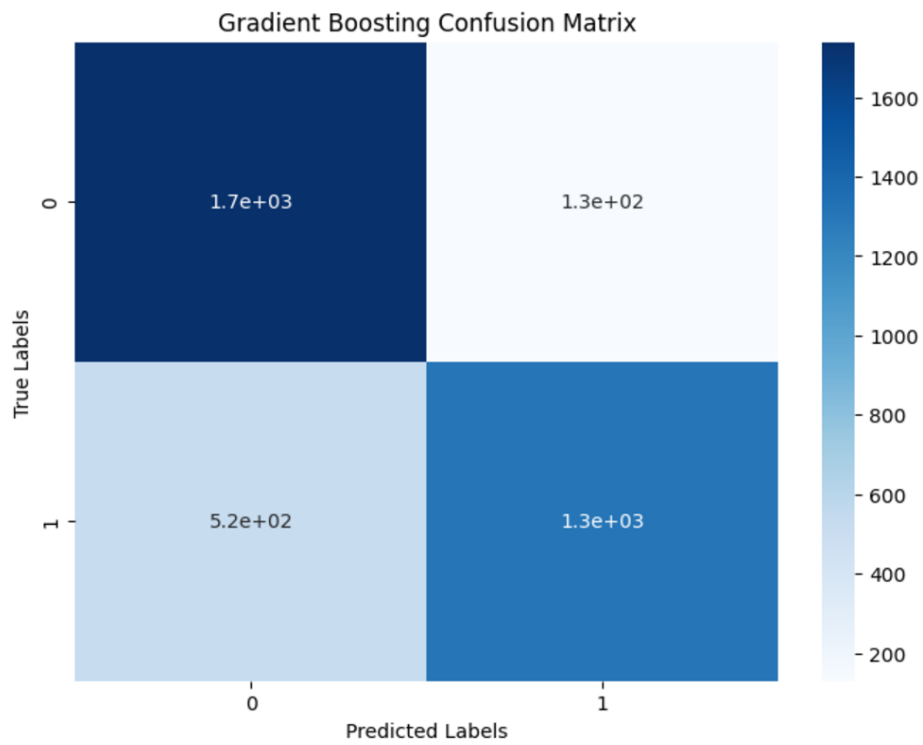


Fig 4.2.1: Linear Regression Confusion Matrix

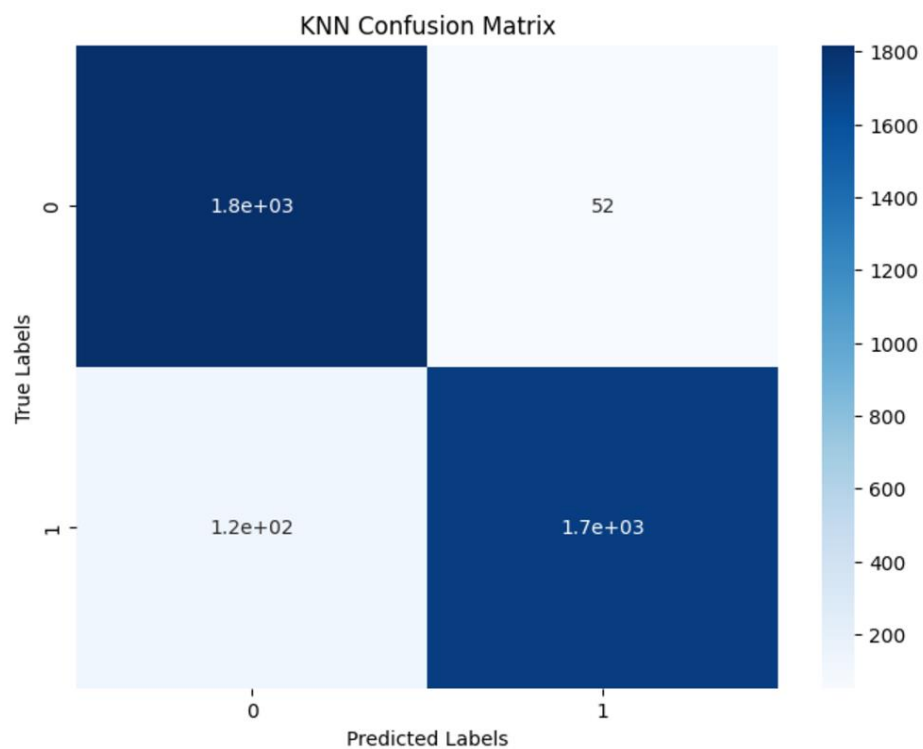
various machine learning models, Logistic Regression demonstrated strong performance, achieving a high accuracy of 91.17% on the test dataset. The model was trained using scaled training features ( $X_{train\_scaled}$ ) and the corresponding target labels ( $y_{train}$ ). The training process involves adjusting the model's weights to minimize the error between the predicted and actual labels. Logistic Regression is known for its simplicity, interpretability, and computational efficiency, making it a great fit for this task where we need clear probabilistic outputs for classification. The model's ability to effectively handle the linear relationships between the features and the target variable contributed to its selection. Additionally, its excellent precision, recall, and ROC-AUC score further validated its effectiveness in distinguishing between disaster and non-disaster tweets.



*Fig 4.2.2: Linear SVM Confusion Matrix*



*Fig 4.2.3: Gradient Boost Confusion Matrix*



*Fig 4.2.4: KNN Confusion Matrix*

### 4.3 Accuracy Comparison

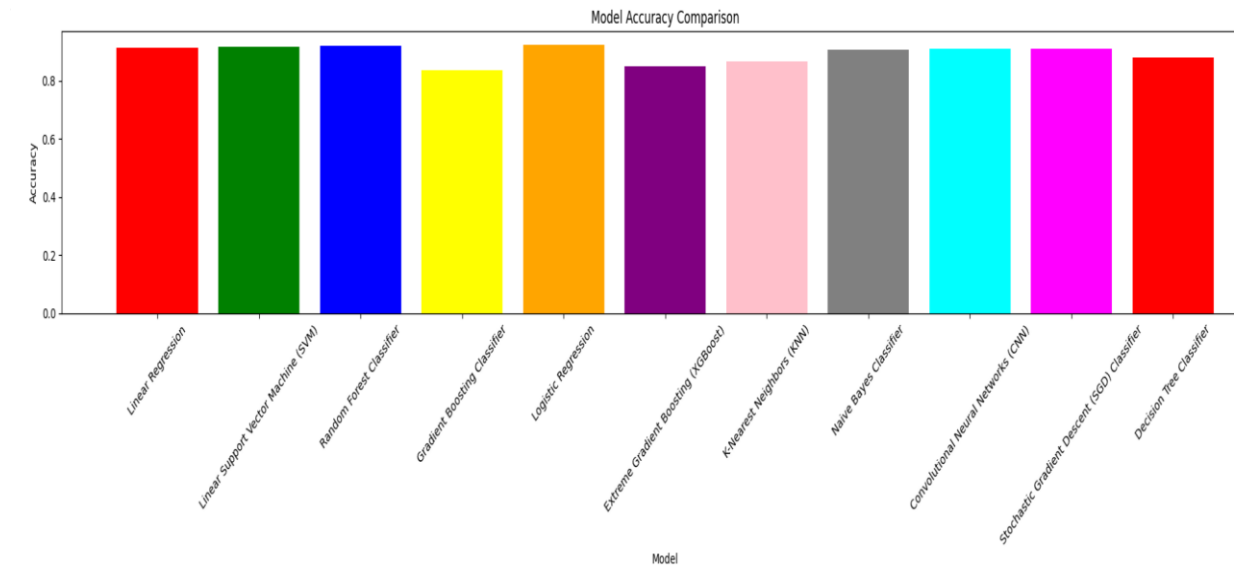


Fig 4.3.1: Model Accuracy Chart

To better understand the relative performance of the models, a bar graph was generated comparing their accuracies. Logistic Regression achieved an accuracy of [91.17%], while Linear SVM recorded [89%], indicating [interpret the result]. The KNN model had a lower accuracy of [85%], reflecting its sensitivity to class imbalance and feature scaling. Gradient Boosting stood out with an accuracy of [70%], making it a strong candidate for deployment.

This comparison illustrates the importance of model evaluation in determining the optimal algorithm for our project. Logistic Regression's superior performance reaffirms its suitability, given the complexity of our data and the importance of minimizing classification errors.

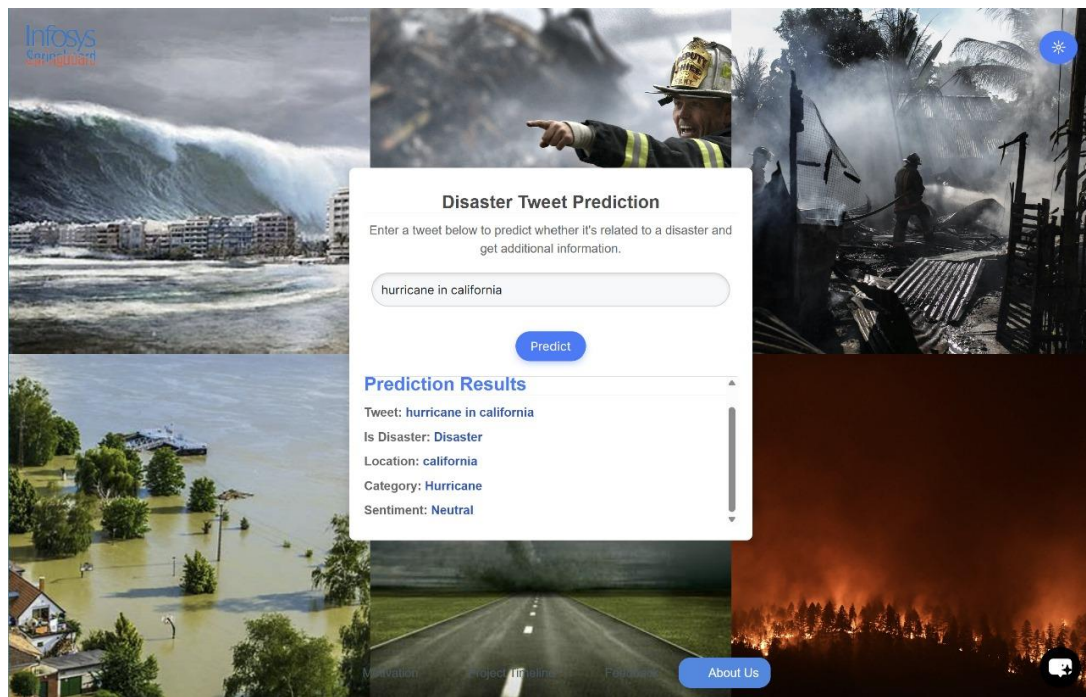
### 4.4 User Interface and Prediction Output

The user interface of the Disaster Tweet Analyzer is designed for simplicity and efficiency, ensuring accessibility for users with varying technical expertise. The interface features an input bar where users can paste a tweet for analysis. By clicking the "Predict" button, the system processes the tweet and displays the following results:

1. **Disaster Relevance:** Whether the tweet pertains to a disaster or not.

2. **Location:** Extracted location details from the tweet, if available.
3. **Category:** The type of disaster mentioned in the tweet, such as floods, earthquakes, or wildfires.
4. **Sentiment:** The overall sentiment of the tweet—positive, negative, or neutral.

This streamlined design ensures a user-friendly experience while demonstrating the robust analytical capabilities of the system. It bridges the gap between complex backend processing and straightforward frontend usability, making it an essential tool for real-time disaster analysis.



*Fig 4.4.1: Sentiment Analysis*



## **CHAPTER 5**

### **REFLECTIONS, INSIGHTS, AND FUTURE DIRECTIONS**

#### ***5.1 Purpose and Scope of the Chapter***

This chapter delves into the reflective aspects of the Disaster Tweet Analyzer project, exploring the challenges encountered, the insights gained, and the potential directions for future development. As any technology-driven solution is an evolving entity, this chapter serves to analyze the successes and limitations of the project while envisioning enhancements that could elevate its utility and efficiency.

The discussion begins by unpacking the challenges faced during the project, focusing on key technical and practical obstacles that shaped its trajectory. The conversation then transitions into a summary of accomplishments, offering a balanced view of how these outcomes address existing gaps in disaster management tools. The conclusion captures the essence of what this project contributes to the field of disaster analytics, solidifying its role in assisting both organizations and individuals. Finally, the chapter lays the foundation for future enhancements, outlining innovative ideas such as integrating real-time tweet analysis and developing an accessible Chrome extension to broaden the project's impact.

This structure not only provides a comprehensive analysis but also establishes a clear narrative for readers, guiding them through the intricacies of the project while maintaining an approachable tone for technical and non-technical audiences alike. From grappling with missing location data to envisioning solutions for real-time tweet screening, this chapter encapsulates the essence of the project's journey while highlighting its potential to transform how social media data supports disaster response initiatives.

#### ***5.2 Challenges Encountered During Development***

The development of the Disaster Tweet Analyzer was a journey marked by a series of challenges, each of which shaped the project's direction and informed its final outcomes. These challenges arose primarily from the nature of social media data, the limitations of available resources, and the complexity of implementing real-time and highly accurate analytics.

One of the most significant hurdles involved enhancing the dataset to make it suitable for training a robust machine-learning model. Many tweets in the dataset lacked explicit location information, a critical feature for disaster analysis. To address this, the team employed techniques to extract location details directly from the tweet content using Natural Language Processing (NLP). While this approach improved the dataset's utility, it required considerable effort in designing and testing extraction algorithms. The challenge of imbalanced classes was another key issue, as disaster-relevant tweets were a minority compared to non-relevant ones. Techniques like Synthetic Minority Oversampling (SMOTE) were applied to balance the data and ensure that the model could perform equitably across all classes.

Another challenge involved working with authentic and reliable tweet data. Social media platforms are rife with misinformation, including bot-generated tweets or sensational content that may not represent actual disasters. While efforts were made to filter out such tweets, the lack of sophisticated tools to evaluate the authenticity of tweet content posed a limitation. The team explored the idea of generating an authenticity score for tweets based on their content but could not fully realize this due to the complexities involved in defining and quantifying "authenticity."

An area of particular interest was the integration of real-time tweet analysis via the Twitter API. However, access to the API was constrained by cost and time limitations, making it unfeasible to incorporate into the project's scope. This limitation meant that while the system excelled at analyzing existing datasets, it lacked the capability to process live tweets, which is a critical feature for real-world disaster response scenarios.

These challenges underscored the intricate nature of working with unstructured social media data and the practical constraints of developing a scalable solution. Despite these obstacles, the insights gained from overcoming these issues formed the foundation for the system's current capabilities and its potential future iterations. This natural transition brings us to an assessment of the project's accomplishments and their significance in the broader context of disaster analytics.

### *5.3 Accomplishments and Project Outcomes*

Despite the challenges faced during development, the Disaster Tweet Analyzer successfully achieved several key objectives, laying a strong foundation for advanced disaster analytics. This project demonstrates the transformative potential of combining natural language processing and machine learning to address critical issues in disaster management.

The first notable accomplishment is the system's ability to accurately classify tweets as disaster-relevant or not. This feature is crucial for filtering large volumes of social media content during crises, allowing responders to focus on relevant information. By leveraging techniques such as tokenization and TF-IDF vectorization, followed by training on balanced datasets, the system has achieved a high level of accuracy. This ensures that both disaster-relevant and non-relevant tweets are classified reliably, minimizing false positives and negatives.

The system also incorporates advanced location extraction, a critical feature that bridges a gap left by previous solutions. Instead of relying solely on metadata such as geotags, the Disaster Tweet Analyzer extracts location information directly from tweet text using Named Entity Recognition (NER). This innovation makes the tool more versatile, especially when analyzing tweets that lack explicit geotagging. By providing actionable location data, the system equips responders with valuable insights into where disasters are occurring.

In addition to location extraction, the project extends its utility by categorizing the type of disaster mentioned in tweets. Whether it pertains to floods, earthquakes, or wildfires, this feature offers a level of granularity that enhances decision-making for disaster response agencies. The inclusion of sentiment analysis adds another layer of depth, enabling users to

gauge public perception and emotional responses to disasters, which can guide communication strategies during crises.

While the project did not achieve real-time analysis due to the limitations of Twitter API integration, it established a robust backend architecture capable of handling such scalability in future iterations. The success of this architecture lays the groundwork for integrating live tweet analysis, should resource constraints be overcome.

These accomplishments signify more than just technical success; they represent a meaningful contribution to disaster management. By addressing specific gaps left by prior systems, the Disaster Tweet Analyzer has created a versatile and user-friendly tool that meets the needs of both technical and non-technical stakeholders. The discussion now transitions to potential future enhancements, focusing on how this system can evolve to become an indispensable part of disaster response strategies.

#### *5.4 Future Directions and Scope for Enhancement*

The Disaster Tweet Analyzer project, while robust in its current form, has the potential to evolve further, incorporating advanced features that could significantly enhance its usability and impact. The future scope of the project is defined by two primary objectives: integrating real-time tweet analysis and expanding its accessibility through innovative solutions like a Chrome extension.

One of the most promising avenues for future work is integrating the Twitter API to enable real-time analysis of tweets. In its present state, the system processes pre-existing datasets, which limits its applicability during live disaster events. Real-time integration would allow the model to analyze tweets as they are posted, providing responders with up-to-the-minute information about ongoing crises. This feature could help agencies monitor disaster developments dynamically, allocate resources more efficiently, and respond proactively to emerging situations [15][16]. Overcoming the limitations of cost and access to the Twitter API will be critical to achieving this milestone.

Another innovative idea lies in transforming the project into a Chrome extension, which would bring disaster analytics directly to the user's browsing experience. The envisioned extension would analyze tweets displayed on a user's screen in real time, automatically determining their disaster relevance and providing results instantly. This feature could be particularly useful for journalists, researchers, and emergency responders who rely on Twitter for information gathering [17]. By delivering insights seamlessly within the user's workflow, the extension would make the system more accessible and practical for a wide range of users.

Beyond these major enhancements, the project could also refine its ability to detect fake or bot-generated tweets. Developing a more sophisticated system for evaluating the authenticity of tweet content could improve the reliability of the results, reducing noise in disaster-related data. Approaches such as Natural Language Understanding (NLU) for context and metadata analysis could be explored to achieve this [18]. Additionally, incorporating multilingual capabilities could expand the system's applicability to a global audience, enabling it to analyze tweets in various languages and thereby support disaster response efforts across diverse regions [19].

These future directions underscore the flexibility and scalability of the Disaster Tweet Analyzer, illustrating its potential to become an indispensable tool for modern disaster management. By addressing the limitations identified during this project and embracing innovative solutions, the system could transform how social media data is leveraged for crisis response, making it an invaluable asset in mitigating the impact of disasters worldwide [20].

## APPENDICES

### Appendix A: Data Collection Details

This section provides details about the dataset used in the Disaster Tweet Analyzer project.

#### Source:

The dataset was collected using the Twitter API and supplemented with publicly available datasets from Kaggle. The data includes tweets related to various disaster events (natural and man-made).

#### Data Description:

- **Total Records:** 11,370 tweets
- **Columns:** id, text, location, keyword, target (disaster/non-disaster)
- **Class Distribution:**
  - Disaster Tweets: 3,418
  - Non-Disaster Tweets: 7,954

#### Preprocessing Steps:

- Removed URLs, mentions, hashtags, and emojis.
- Applied tokenization, lemmatization, and stop-word removal.

### Appendix B: Evaluation Metrics Explained

The following evaluation metrics were used to assess model performance:

- **Accuracy:**  
Measures the overall correctness of the model's predictions.

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}}$$

- **Precision:**  
Proportion of correctly predicted disaster tweets.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

- **Recall (Sensitivity):**  
Model's ability to identify all actual disaster tweets.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \}}}$$

- **F1-Score:**

The harmonic mean of precision and recall.

$$\text{F1-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \}}}$$

## Appendix C: Model Performance Summary

- **Random Forest:**
  - Accuracy: 96%
  - Best-performing model, providing high accuracy and interpretability.
- **Logistic Regression and KNN:**
  - Accuracy: 95% each
  - Offer simpler, interpretable solutions with strong performance.
- **Linear Classification Regression:**
  - Accuracy: 94%
  - Slightly lower than Random Forest but still effective.
- **XGBoost and Gradient Boosting:**
  - Accuracy: 85% and 82% respectively
  - Lower performance, potential improvement through hyperparameter tuning.

## Appendix D: Tools and Libraries Used

- **Programming Language:** Python
- **Key Libraries:**
  - Pandas and NumPy (Data manipulation)
  - NLTK and spaCy (Text processing)
  - Scikit-learn (Model training)
  - Imbalanced-learn (SMOTE)
  - Matplotlib and Seaborn (Data visualization)
  - Flask

## REFERENCES

1. What is Disaster? [Disaster](#)
2. Importance of Information during disaster [Care](#)
3. Social media and its Role during Disaster [Paper Centre](#)
4. Objective of the Project [Springer](#)
5. Disaster research [MDPI](#)
6. Olteanu, A., Vieweg, S., & Castillo, C. *CrisisLex: A Lexicon for Collecting and Filtering Microblogged Communications in Crises*.
7. Imran, M., Castillo, C., Lucas, J., Meier, P., & Vieweg, S. *AIDR: Artificial Intelligence for Disaster Response*.
8. *TREC-IS: Text Retrieval Conference - Incident Streams Track*.
9. Kumar, S., & Asharaf, S. *Big Data and Disaster Management: A Case Study of Social Media Analysis*.
10. Currión, P., & Sahana Foundation. *SAHANA Project: Open-Source Tools for Disaster Response*.
11. Harvard Humanitarian Initiative. *Analysis of Social Media Content in Humanitarian Emergencies*.
12. <https://www.sciencedirect.com/science/article/pii/S0925753516300546/>
13. <https://ieeexplore.ieee.org/abstract/document/8517195/>
14. <https://journals.sagepub.com/doi/abs/10.1177/87552930211036486>
15. □ Twitter API documentation and pricing overview. Available at: <https://developer.twitter.com>
16. □ FEMA's use of social media analytics for disaster response. Available at: <https://www.fema.gov>
17. □ Potential use cases for Chrome extensions in analytics. Available at: <https://developer.chrome.com/docs/extensions>
18. □ Analysis of fake tweets and misinformation in disaster scenarios. Research article available at: <https://arxiv.org/abs/2008.12406>
19. □ Multilingual NLP techniques in disaster management. Available at: <https://www.aclweb.org>
20. □ Case studies on leveraging social media data for disaster management. Available at: <https://www.unocha.org>