

RISK PRO: GRADIENT BOOSTING FOR DISASTER FORECASTING AND PREVENTION

A MINI PROJECT REPORT

Submitted by

**SANTHOSH B (221801046)
THOFIQ GANI M (221801057)**

in partial fulfillment for the award of the degree of

**BACHELOR OF TECHNOLOGY
IN
ARTIFICIAL INTELLIGENCE AND DATA SCIENCE**



RAJALAKSHMI ENGINEERING COLLEGE

ANNA UNIVERSITY:CHENNAI 600 025

NOVEMBER 2024

ANNA UNIVERSITY, CHENNAI

BONAFIDE CERTIFICATE

Certified that this Report titled **“Risk Pro: Gradient Boosting for Enhanced Disaster Forecasting and Prevention”** is the bonafide work of **SANTHOSH B (221801046), THOFIQ GANI M (221801057)** who carried out the under my supervision. Certified further that to the best of my knowledge the work reported here in does not form part of any other thesis or dessirtation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

SIGNATURE

Dr.J.M.GNANASEKAR M.E.,Ph.D.,

Professor and Head

Department of Artificial Intelligence.

Data Science

Rajalakshmi Engineering College

Chennai-602105

SIGNATURE

Mrs.P.JAYASRI ARCHANA

DEVI, M.E.,

Assistant Professor(SG)

Department of Artificial Intelligence and

Data Science

Rajalakshmi Engineering College

Chennai-602105

Submitted for the project viva-voce examination held on_____

INTERNAL EXAMINAR

EXTERNAL EXAMINAR

ACKNOWLEDGEMENT

Initially we thank the Almighty for being with us through every walk of our life and showering his blessings through the endeavor to put forth this report. Our sincere thanks to our Chairman Mr. S. MEGANATHAN, B.E, F.I.E., our Vice Chairman Mr. ABHAY SHANKAR MEGANATHAN, B.E., M.S., and our respected Chairperson Dr. (Mrs.) THANGAM MEGANATHAN, Ph.D., for providing us with the requisite infrastructure and sincere endeavoring in educating us in their premier institution.

Our sincere thanks to Dr. S.N. MURUGESAN, M.E., Ph.D., our beloved Principal for his kind support and facilities provided to complete our work in time. We express our sincere thanks to Dr. J. M. GNANASEKAR., M.E., Ph.D., Professor and Head of the Department of Artificial Intelligence and Data Science for his guidance and encouragement throughout the project work. We are glad to express our sincere thanks and regards to our supervisor Mrs. P. JAYASRI ARCHANA DEVI, M.E., Assistant Professor-SG, Department of Artificial Intelligence and Data Science and coordinator, Dr. P. INDIRA PRIYA, M.E., Ph.D., Professor, Department of Artificial Intelligence and Data Science, Rajalakshmi Engineering College for her valuable guidance throughout the course of the project.

Finally we express our thanks for all teaching, non-teaching, faculty and our parents for helping us with the necessary guidance during the time of our project.

Abstract

Urban areas are increasingly vulnerable to disasters due to rapid urbanization and the accelerating impacts of climate change. Traditional risk assessment methods, while valuable, often struggle to integrate the vast amounts of real-time data from diverse sources, resulting in limited accuracy and timeliness. This highlights a pressing need for data-driven approaches that can provide precise, real-time risk evaluations tailored to the complexities of urban environments. In response, we propose a machine learning-based framework specifically designed to enhance the accuracy and relevance of disaster risk assessments for urban areas. By leveraging recent advancements in data science and machine learning, this tool aggregates and analyzes data from multiple sources, including weather reports, satellite imagery, historical disaster records, and real-time sensor data. The system's algorithms classify and interpret these data points to generate predictive risk assessments that can adapt to rapidly changing conditions.

The primary goal of this tool is to improve urban disaster preparedness and response by providing decision-makers with actionable insights. The enhanced accuracy of risk predictions can inform proactive measures such as targeted resource allocation and early evacuation planning, significantly contributing to urban safety and resilience. Initial tests indicate that our model outperforms existing methods in speed and accuracy, especially in adapting to evolving urban risk landscapes. This research presents a novel machine learning application for disaster management, emphasizing timely, data-driven insights to protect urban populations and infrastructure. Future efforts will expand data sources and enhance predictive capabilities, positioning the tool as a vital asset for urban resilience planning and decision-making amid rising disaster risks.

TABLE OF CONTENTS

CHAPTER	TITLE	PAGE NO
	ABSTRACT	
	LIST OF TABLES	
	LIST OF FIGURES	
1	INTRODUCTION	1
	1.1 GENERAL	1
	1.2 NEED FOR THE STUDY	2
	1.3 OBJECTIVES OF THE STUDY	3
	1.4 OVERVIEW OF THE PROJECT	4
2	REVIEW OF LITERATURE	6
	2.1 INTRODUCTION	6
	2.2 LITERATURE REVIEW	7
3	SYSTEM OVERVIEW	13

	3.1 EXISTING SYSTEM	13
	3.2 PROPOSED SYSTEM	14
4	SYSTEM REQUIREMENTS	15
	4.1 SOFTWARE REQUIREMENT	15
5	SYSTEM DESIGN	17
	5.1 SYSTEM ARCHITECTURE	17
	5.2 MODULE DESCRIPTION	19
	5.2.1 DATA COLLECTION	19
	5.2.2 DATA PREPROCESSING	20
	5.2.3 EXPLORATY DATA ANALYSIS	21
	5.2.4 FEATURE ENGINEERING	22
	5.2.5 TRANING GRADIENT BOOSTING	24
	5.2.6 MODEL EVALUATION	25
	5.2.7 VISUALIZATION OF RESULTS	26

6	RESULT AND DISCUSSION	28
	6.1 RESULT	28
	6.3 DISCUSSION	30
7	CONCLUSION AND FUTURE ENHANCEMENT	31
	7.1 CONCLUSION	31
	7.2 FUTURE ENHANCEMENT	31
8	APPENDIX	33
	8.1 SAMPLE CODE	33
	8.2 OUTPUT	42
9	REFERENCE	48

LIST OF TABLES

TABLE NO	TABLE NAME	PAGE NO
1	LITERATURE REVIEW	6

LIST OF FIGURES

FIGURE NO	FIGURE NAME	PAGE NO
1	SYSTEM ARCHITECTURE	17
2	DATA COLLECTION MODULE	19
3	DATA PREPROCESSING MODULE	20
4	EDA MODULE	21
5	FEATURE ENGINEERING MODULE	22
6	GRADIENT BOOSTING TRAIN MODULE	24
7	MODEL EVALUATION MODULE	25
8	VISUALIZATION MODULE	26
9	LOAD AND INSPECTING DATA	42
10	DATA PREPROCESSING	43
11	DISTRIBUTION OF DEATHS	43
12	CORRELATION HEATMAP	44

13	FEATURE ENGINEERING	45
14	GRADIENT BOOSTING MODEL TRAIN	46
15	MODEL EVALUTION	46
16	ACTUAL VS PREDICTED DEATHS	47
17	FEATURE IMPORTANCE	47

CHAPTER 1

INTRODUCTION

1.1 GENERAL

As cities expand and climate change leads to more frequent extreme weather, urban areas face heightened disaster risks, including floods, storms, and other severe events. Traditional disaster assessment methods often struggle to meet current needs, relying on outdated data and static models that fail to reflect realtime changes or emerging threats. These limitations make it difficult for city planners and emergency responders to make fully informed decisions, leaving urban areas vulnerable during crises.

To overcome these challenges, this study introduces a machine learning solution using Gradient Boosting, an algorithm that builds models by iteratively refining predictions to increase accuracy. By integrating real-time data from sources such as weather stations, satellite imagery, and infrastructure sensors, the model can continuously adapt to new information, offering timely risk assessments as conditions change. For example, during a sudden rainfall event, the model can update flood risk projections instantly, helping officials respond proactively and allocate resources more effectively.

This Gradient Boosting-based system not only improves the precision and timeliness of disaster risk assessments but also supports data-driven decisions that strengthen urban resilience. With better insights, city officials can plan more strategically, reduce response times, and ultimately protect both lives and assets. By enabling faster, more accurate disaster preparedness, this approach aims to significantly lower the social and economic costs of disasters, fostering safer, more resilient urban communities in the face of growing climate-related risks.

1.2 NEED FOR THE STUDY

Disaster management cannot be successful without good risk assessment to reduce the risk to lives and property and economies. Technological risk exercise models or techniques used in conventional approaches were based on aggregated statistical data or expert opinions which are effective in their results, yet sometimes provide insufficient information on the shift in disaster risks. Since climate change has led to more often and more severe natural disasters, such methods are slow and rigid and may fail to adapt to real-time changes and other dynamic factors resulting from climate change.

It is thus imperative for a consistent and objective means of classifying breast carcinomas: one that can rely on patterns and data. A popular and highly effective method like Gradient Boosting is capable of analysing large, varied datasets and finding patterns of data interactions that can often remain unnoticed by classical data analysis methods. Not only does this enhance the specificity of risk predictions but it also permits consequent, real-time adjustments. In the environment of an automated machine learning system, disaster risk can be divided into high, medium and low-risk categories, giving authorities unambiguous categorizations.

Housing such data, this system basically provides decision makers with an impartial and recent evaluation to allocate resource efficiently – focusing on areas that require interventions and support and support disaster prevention and preparedness measures. Based on the results of the analysis, it is also possible to allocate funds more effectively, so that the infrastructure in weak areas can be improved, and planning for disasters can be enhanced. Finally, this approach strengthens protection and increases possibilities for minimizing the losses in case of a disaster in terms of both social and economical impacts as well as providing a valuable tool for stakeholders practicing disaster management in the context of constant change.

1.3 OBJECTIVES OF THE STUDY

The main objectives of this study are outlined as follows:

1. **Development of a Predictive Model:** The first objective is, therefore, to develop a strong and reliable forecasting method of the risk levels of disaster incidents in different areas. This model will incorporate probabilistic data sets concerning past disaster occurrences and the different environmental characteristics to obtain a correlation that explains the conditions that amplify risk.
2. **Implementation of Gradient Boosting Algorithm:** The objective of the study is to use Gradient Boosting algorithm given its suitability in working with such complex data structures and identifying complex relationships. The model will be refined and adjusted to achieve best possible results and allow for clear distinction of the risk levels within districts as high or average or low risk level.
3. **Visualization of Risk Predictions:** For purpose of improving the communication with stakeholders and interpretation of the results, the study will involve the creation of instruments such as visualisation which will depict the risk predictions acquired. These visualizations will offer useful information that will easily guide the disaster preparedness approaches and decentralisation of decision-making for the crisis management bodies and policy makers.
4. **Evaluation of Model Effectiveness:** One of the elements of the analysis is the measure of the accuracy of the applied model by contrasting the native predictions on disasters with the historical data. This evaluation will relates areas will involve computation of accuracy, precision, and recall so as to determine the of the model. In this study, the model will be tested and validated to show that it can be used in disaster risk management a real-life process.
5. **Facilitation of Data-Driven Decision Making:** Thus, the study seeks to develop a framework that will equip decision makers with analysis for their

decision making processes. The system will facilitate identifying various risks and help the disaster management authorities undertake necessary action of preventing disasters that have severe impacts on vulnerable groups of the population.

1.4 OVERVIEW OF THE PROJECT

This work involves developing a disaster risk evaluation framework and employing the Gradient Boosting algorithm to predict the risk levels of regions. The system is aimed to produce effective recommendations for all the participants of the disaster management and planning processes.

Key Features:

1. **Data Collection:** The system compiles disaster history, climatic data and geographical information from different sources meaning that it would have a data set for analysis.
2. **Data Preprocessing:** The actual collected data are preprocess for cleaning, standardization and transformation to increase the model quality.
3. **Model Training:** To build the predictive model, using of the Gradient Boosting algorithm is used. This model makes accurate assessments of disaster risks since it fits various and intricate spatio-temporal variations and interdependencies.
4. **Risk Prediction:** Organizations have to classify areas in terms of high, moderate or low risk according to the environmental characteristics and past events.

Workflow:

- **Data Input and Preprocessing:** Information that is gathered is from various sources and undergoes several processes to be ready for modeling.

- Model Training and Validation: With respect to Gradient Boosting, the data preprocess, validated, and model-tuned for the highest possible accuracy.
- Prediction and Classification: The model sorts the regions into definite risk profiles according to patterns that were discerned from the data.

CHAPTER 2

REVIEW OF LITERATURE

2.1 INTRODUCTION

Machine learning / artificial intelligence is on the cutting edge of change in disaster risk assessment where high speed / accuracy in passing data through the model is critical. Conventional methodologies in disaster risk analysis may perhaps entail the use of historical information in addition to elementary statistical methods something which fails to capture the dynamics and complexity of the risk phase. The inherent limitations associated with these more conventional approaches mean though that they are unable to generate sufficient and timely risk estimations given the advances in the frequency and scale of natural disasters as a result of climate change. Machine learning, therefore presents the ability to conduct real time analysis for multiple factors such as the environment, geography and socio-economic factors thus providing a comprehensive analysis.

Of the diverse approaches to the algorithm of machine learning, tree-based methods such as Gradient Boosting have shown great potential in risk classification. As an ensemble technique that uses decision trees, Gradient Boosting particularly reflects the complexity of connections within relationships and is thereby well suited for disaster risk assessments. This case sits well supported by previous work suggesting that models like Gradient Boosting and models like Random Forests are often more accurate and reliable than conventional statistical techniques. This chapter aims at reviewing published literature that has considered the use of machine learning, particularly Gradient Boosting and other tree-based models to estimate risks of disasters. The review explains the usefulness and difference of these models compared to conventional methods and shows how machine learning is changing the disaster management by predicting the best actions for each case.

2.2 LITERATURE REVIEW

S. No	Author Name	Paper Title	Description	Journal	Year
1	J.Zhang, L.Xu, X. Wu	A Framework for Assessing Disaster Risk and Resilience in Urban Areas	Uses socioeconomic and infrastructural factors.	IEEE	2019
2	S.K.Sharma, V.R.Patel, A.K .Joshi	Assessment of Natural Disaster Risks Using Artificial Intelligence Technique	Applies AI, like neural networks and decision trees, to predict and reduce disaster risks.	IEEE	2020
3	R. A. Silva, E.G.de Oliveira, P. A. Lima	Integration of IoT and Big Data Analytics for Disaster Risk Assessment and Management	Utilizes IoT and big data for realtime disaster risk assessment and management.	International Journal of Disaster Risk Reduction	2021
4	M.A.Hossain, S.Roy, M. Rahman	Remote Sensing and GIS for Disaster Risk Assessment: A Case Study of Flooding	Study examines remote sensing and GIS for improved flood risk mapping and management.	IEEE	2020

The literature survey also demonstrates many enhanced techniques for the evaluation of disaster risks.

Zhang, Xu, and Wu (2019) - Socio-Economic and Infrastructure-Based Disaster Risk Framework:

Core Focus: This study centers on urban resilience, particularly how cities can withstand and recover from disasters. Urban areas are particularly vulnerable due to dense populations and complex infrastructures, making effective disaster management crucial.

Methodology: The authors devised a framework that integrates various socioeconomic and infrastructural factors. Socio-economic factors could include demographics, economic activities, and social inequality, which affect how people respond to and recover from disasters. Infrastructural factors like transportation networks, utility systems, and building resilience determine how quickly a city can resume normal functions after a disaster.

Implications: By providing a holistic framework that addresses both human and infrastructural vulnerabilities, the study serves as a tool for policymakers and urban planners. This approach not only helps in identifying weak points but also supports long-term resilience planning, emphasizing proactive rather than reactive disaster management.

Published by IEEE: As an IEEE paper, it likely includes technological perspectives and detailed models, focusing on applying the framework in practical, data-driven scenarios in urban settings.

Sharma, Patel, and Joshi (2020) - AI Techniques for Disaster Risk Prediction:

Core Focus: This research leverages artificial intelligence (AI), specifically machine learning techniques, to assess and mitigate risks from natural disasters.

It emphasizes the predictive power of AI in forecasting events such as earthquakes, floods, or hurricanes.

Technology Used: The study uses neural networks and decision trees—two powerful AI models. Neural networks are adept at recognizing complex patterns within vast datasets, making them valuable for analyzing past disaster occurrences and environmental conditions to predict future events. Decision trees are structured models that make predictions based on a set of choices, allowing the system to identify critical risk factors and their potential outcomes. For example, a decision tree might analyze factors like weather patterns, terrain, and past flooding events to predict flood risks in a specific area.

Implications: The AI-driven predictions can guide authorities in deploying resources and implementing early warning systems, potentially saving lives and reducing property damage. By integrating AI, this study illustrates a shift toward data-driven disaster management that is both adaptable and capable of learning from new data over time, making it more accurate as the system evolves.

Published by IEEE: As an IEEE publication, it likely provides insights into the technical aspects of AI algorithms and their customization for disaster scenarios, useful for practitioners in engineering and disaster technology.

Silva, Oliveira, and Lima (2021) - IoT and Big Data for Real-Time Disaster Management:

Core Focus: This paper explores the convergence of the Internet of Things (IoT) and big data analytics to create an interactive, real-time disaster management system. The goal is to enable continuous monitoring and instant response during emergencies, such as earthquakes, tsunamis, or wildfires.

Technology Used:

IoT: This includes devices like sensors, drones, and connected cameras strategically placed in vulnerable areas. For instance, water-level sensors along rivers can detect rising levels and send alerts, while drones can provide aerial imagery of wildfire progression.

Big Data Analytics: With a constant stream of data from IoT devices, big data tools process this information to identify patterns, anomalies, and trends. For example, sudden changes in soil moisture levels could indicate an impending landslide, allowing for early warnings.

Implications: The integration of IoT and big data enables quick decision-making, allowing authorities to respond immediately to emerging threats. This system helps reduce the time between detection and response, which is critical in minimizing the impact of disasters. Additionally, the system can be scaled up or down depending on the severity of the event, making it adaptable for various disaster scenarios.

Published in the International Journal of Disaster Risk Reduction: This journal focuses on applied research for reducing disaster risks, suggesting that the paper likely includes case studies or practical applications, making it valuable for both researchers and practitioners in real-world disaster management.

Hossain, Roy, and Rahman (2020) - Remote Sensing and GIS for Flood Risk Assessment:

Core Focus: This study investigates how remote sensing and Geographic Information Systems (GIS) can be applied to assess and manage flood risks, particularly in regions prone to flooding.

Technology Used:

Remote Sensing: Through satellites or aerial imagery, remote sensing allows continuous observation of vast areas. This data is useful for monitoring environmental changes, such as river levels, rainfall accumulation, or vegetation cover, which are critical indicators of flood risk.

GIS: GIS allows spatial analysis and mapping of flood-prone regions. By layering various datasets—like topography, land use, and water flow patterns—GIS can produce detailed maps that highlight areas with high flood risk. For instance, combining rainfall data with elevation maps can show which low-lying areas are most susceptible to flooding.

Implications: This technology-driven approach is especially beneficial for developing regions or rural areas with limited disaster response resources. Flood risk maps created using remote sensing and GIS can guide local governments in planning evacuation routes, designing flood defenses, and educating the public on high-risk zones. Additionally, this approach supports environmental monitoring, helping to predict floods in real time based on observed weather and water trends.

Published by IEEE: IEEE's focus on technology suggests that the paper likely covers technical details on how to process and integrate satellite data with GIS software, providing valuable information for engineers and GIS professionals involved in disaster management.

Key Insights from the Survey:

This review of the four studies underscores the power of emerging technologies in transforming disaster management strategies:

Predictive Capabilities: AI-driven models improve disaster forecasting, enabling authorities to prepare in advance and allocate resources effectively.

Real-Time Monitoring: IoT and big data facilitate instant data collection and processing, allowing for quick response, which is critical in reducing casualties and damage during disasters.

Spatial Analysis: Remote sensing and GIS provide detailed geographical insights, allowing authorities to map risks accurately and make informed decisions about land use, evacuation routes, and public safety measures.

Holistic Planning: The integration of socio-economic, infrastructural, and technological factors in these frameworks allows for a more comprehensive approach to disaster risk reduction, considering not only the immediate impacts of disasters but also the long-term resilience of communities.

CHAPTER 3

SYSTEM OVERVIEW

Existing disaster assessment methods rely heavily on manual data collection, limited data integration, and static modeling, making them inadequate for realtime risk assessment. Manual data collection is often slow and prone to human error, which can result in delays and outdated information that limit the effectiveness of emergency preparedness efforts. These manual processes are typically resource-intensive, requiring extensive surveys or reports that may not capture the most current risk factors.

Moreover, traditional systems struggle with limited data integration, meaning they cannot easily combine diverse data sources, such as environmental metrics, demographic information, or infrastructure conditions. This fragmented data approach often leads to incomplete risk assessments that may overlook crucial risk factors, reducing their overall accuracy and relevance. Furthermore, these systems use static models with fixed assumptions, which makes them inflexible in adapting to real-time changes or evolving risk conditions. For instance, a static flood risk model might not reflect new rainfall data or recent urban developments, leading to inaccurate risk projections.

Without the ability to adapt to changing conditions, these methods leave city planners and responders relying on outdated or incomplete information, which can slow response times and reduce the effectiveness of disaster preparedness measures. In critical situations, these delays and inaccuracies may increase vulnerability, heightening the social and economic impacts of disasters. To meet the demands of today's rapidly evolving urban landscapes, there is a pressing need for disaster assessment methods that are dynamic, data-driven, and capable of integrating real-time information to support timely, well-informed decisionmaking.

3.2 PROPOSED SYSTEM

Existing disaster assessment methods rely heavily on manual data collection, limited data integration, and static models, making them inadequate for accurate, real-time risk assessment. Manual processes are slow, labour-intensive, and susceptible to human error, often leading to outdated information that hinders effective disaster preparedness. Additionally, traditional systems struggle to integrate diverse data sources—such as environmental metrics, infrastructure status, and real-time population density—resulting in fragmented assessments that may overlook crucial risk factors. Static models, based on fixed assumptions, lack the flexibility to adapt to real-time changes, significantly reducing their accuracy in dynamic conditions.

As a result, city planners and responders are often forced to work with incomplete or outdated information, delaying response times and increasing vulnerability to sudden events. This lack of adaptability in traditional methods not only compromises preparedness but also heightens the potential for social and economic disruptions during disasters. To address today's rapidly evolving urban and climate challenges, there is a pressing need for disaster assessment methods that are dynamic, data-driven, and capable of integrating real-time information. Such advancements would enable cities to make faster, more informed decisions, optimize resources, and proactively enhance resilience, ultimately reducing the impacts of disasters and better safeguarding urban communities against escalating climate risks.

CHAPTER 4

SYSTEM REQUIREMENTS

4.1 SOFTWARE REQUIREMENT

The implementation of the proposed disaster risk assessment system relies on several essential Python libraries and tools for data analysis, machine learning, visualization, and real-time data handling.

For data analysis and preprocessing, libraries like Pandas and NumPy are fundamental. Pandas provides powerful data manipulation tools, enabling efficient handling, cleaning, and transformation of large datasets. NumPy supports numerical operations, offering efficient data structures like arrays that facilitate complex calculations, which are crucial for preparing data before model training.

In terms of machine learning, Scikit-Learn is essential for implementing and finetuning the Gradient Boosting model. This library includes a suite of machine learning algorithms and tools that streamline the process of training, validating, and evaluating models. With Scikit-Learn, users can easily configure model parameters and employ techniques like cross-validation, allowing for the creation of an accurate and robust risk assessment model.

For interpreting results, visualization libraries like Matplotlib and Seaborn are used. Matplotlib provides a broad range of plotting functions, allowing users to create detailed visualizations, such as line charts, histograms, and scatter plots, to examine data patterns and model performance. Seaborn, which builds on Matplotlib, offers enhanced visuals and is particularly useful for creating complex statistical plots, like heatmaps and pair plots, to further analyse correlations and data distributions.

To ensure that the system can handle and integrate real-time data, frameworks capable of managing live data feeds are incorporated. These frameworks can gather inputs from sources such as sensors, APIs, or other data streams, allowing the model to update assessments dynamically as new data becomes available. By integrating these real-time data frameworks, the system remains adaptable to changing conditions, providing city planners and responders with current risk insights to improve disaster preparedness and response efforts.

Altogether, these software tools form a robust and flexible foundation, enabling a data-driven, real-time disaster risk assessment system that is both effective and adaptable to evolving urban and environmental conditions.

CHAPTER 5

SYSTEM DESIGN

5.1 SYSTEM ARCHITECTURE

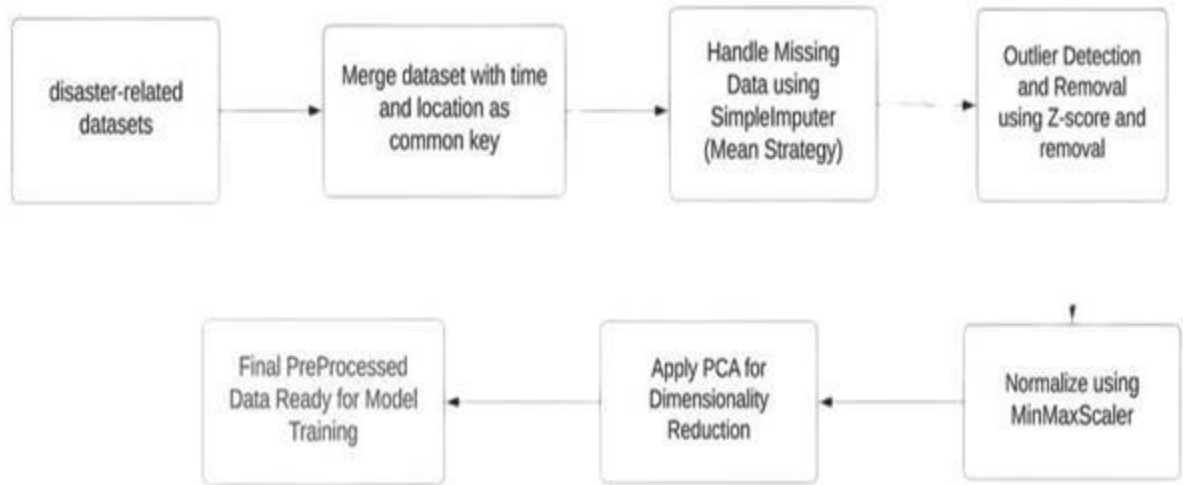


Fig 1 SYSTEM ARCHITECTURE

The proposed system is designed as a data pipeline that continuously collects, processes, and analyses real-time data to deliver accurate and timely disaster risk assessments. This architecture ensures a smooth flow of data from collection to analysis, enhancing the system's ability to update risk insights dynamically. The process begins with real-time data collection from multiple sources, such as weather sensors, satellite feeds, and infrastructure monitoring systems, providing crucial information on environmental conditions, weather patterns, and population metrics. By integrating data from diverse sources, the system forms a comprehensive view of the risk factors influencing disaster vulnerability in urban areas.

Once the data is collected, it undergoes preprocessing to ensure quality and compatibility with the model. This step includes cleaning the data to handle missing values, standardizing formats, and normalizing data ranges. Feature engineering may also be applied to generate new variables, such as combining related data points to create metrics like “disaster severity” or “infrastructure

resilience.” Effective preprocessing ensures the data is accurate, consistent, and optimized for analysis, reducing errors and improving model reliability.

The preprocessed data is then fed into a Gradient Boosting model, a powerful machine learning technique known for its high accuracy and ability to handle complex datasets. In Gradient Boosting, a series of decision trees are built sequentially, with each tree trained to correct errors from the previous one, refining predictions with each iteration. This iterative training process creates an ensemble model that delivers highly accurate disaster risk assessments based on real-time data.

As new data flows into the pipeline, the model updates its predictions to reflect current conditions. For example, if there is a sudden increase in rainfall in a particular region, the system can immediately adjust flood risk predictions for that area. This real-time adaptability is crucial for accurate disaster assessments, as conditions often change rapidly, especially during extreme weather events. The system’s dynamic architecture allows for continuous learning and adjustment, providing city planners and emergency responders with actionable, up-to-date insights to support disaster preparedness and response.

5.2 MODULE DESCRIPTION

5.2.1 DATA COLLECTION

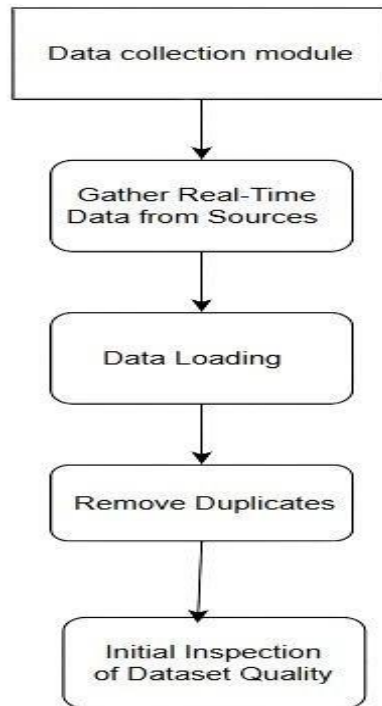


Fig 2 DATA COLLECTION MODULE

The Data Collection Module gathers real-time data from sources like social media, news portals, and official websites using APIs or web scraping. It detects and addresses missing values to ensure dataset completeness, removes duplicates to maintain data uniqueness, and performs an initial quality check to filter out irrelevant content. Collected data is then standardized into a consistent format for easier processing and analysis, and securely stored with backup options to ensure reliability.

5.2.2 DATA PREPROCESSING

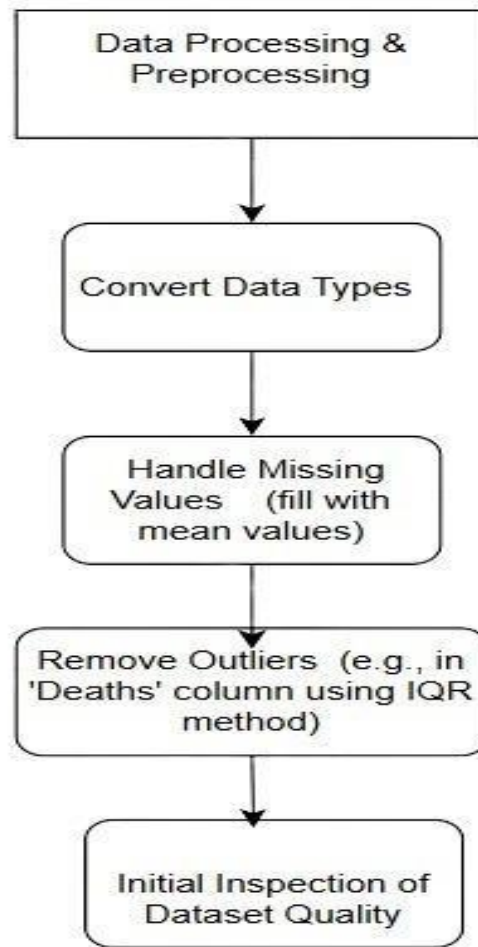


Fig 3 DATA PREPROCESSING MODULE

This module cleans and prepares data for analysis. It converts data types as needed for uniformity, handles missing values by filling them with the mean to maintain dataset completeness, and removes outliers (e.g., in the 'Deaths' column) using the Interquartile Range (IQR) method to improve data accuracy. The module also conducts initial inspections to identify any further issues, ensuring data quality before analysis.

5.2.3 EXPLORATY DATA ANALYSIS (EDA)

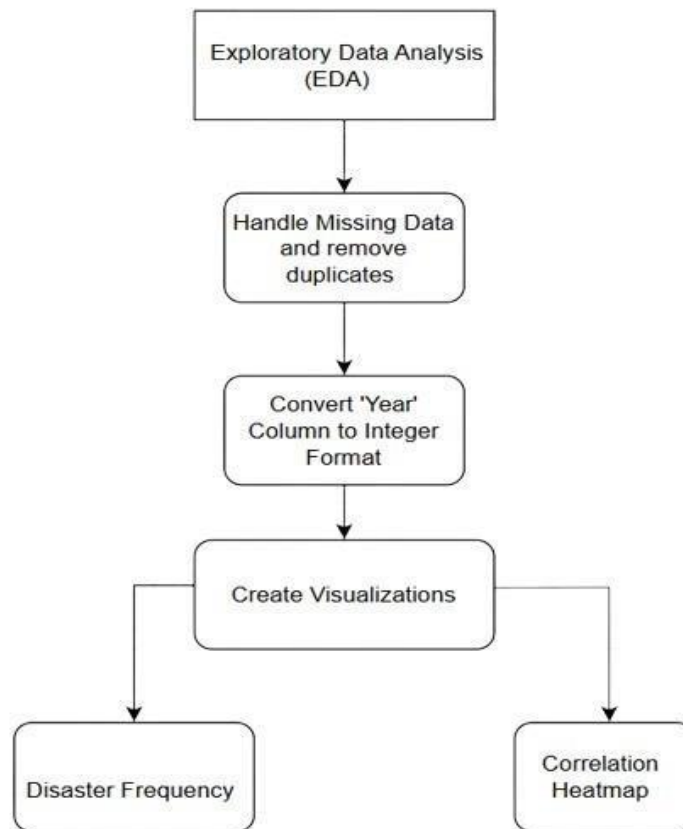


Fig 4 EDA MODULE

The module inspects, cleans, and visualizes the dataset to uncover patterns and relationships. Key steps include handling missing data and duplicates, converting columns like 'Year' to integer format for consistency, and creating visualizations such as disaster frequency charts and correlation heatmaps. This process helps identify trends and relationships in the data, providing valuable insights for further analysis.

5.2.4 FEATURE ENGINEERING

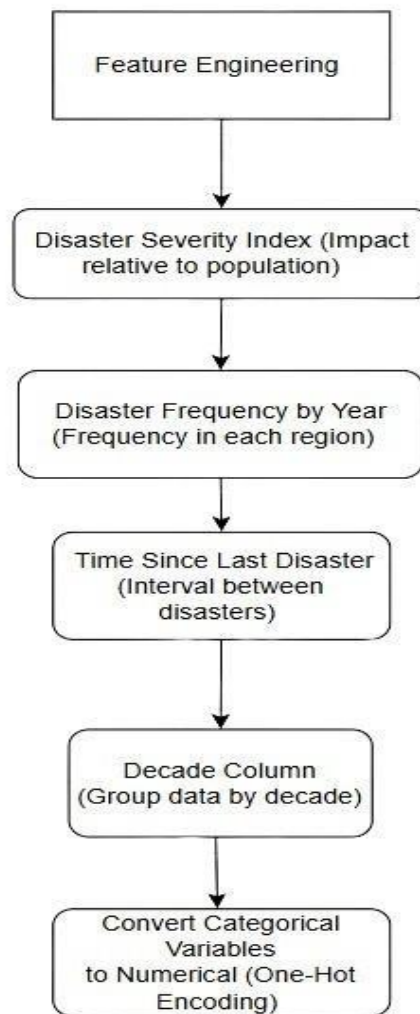


Fig 5 FEATURE ENGINEERING MODULE

This module enhances the dataset by creating new features that improve the model's ability to predict and analyse disaster impacts. Here's a breakdown of each feature and its purpose.

Disaster Severity Index: This feature quantifies the impact of each disaster relative to the population of the affected area, providing a measure of severity that adjusts for population size. By comparing the number of fatalities, injuries, and economic damage to the local population, this index helps highlight disproportionately severe events.

Disaster Frequency by Year: This feature tracks the number of disasters that occur annually within each region. It enables trend analysis, allowing the model to recognize regions or years with frequent occurrences, which can be crucial for understanding patterns in disaster-prone areas.

Time Since Last Disaster: This feature calculates the time elapsed since the last disaster in a given area. It helps capture intervals between disasters, which may reveal cyclical or repetitive patterns and allow predictions based on historical frequency and recurrence.

Decade Column: By grouping events by decade, this feature enables the analysis of long-term trends. Identifying patterns over extended periods provides valuable insights into shifts in disaster frequency, severity, or types across different eras, helping researchers understand how disaster trends evolve over time.

Encoding Categorical Variables: Categorical variables, such as disaster types or location names, are converted into numerical format using one-hot encoding. This transformation ensures that categorical data can be interpreted by machine learning models, allowing algorithms to process these variables effectively without assuming any ordinal relationships.

5.2.5 TRAIN GRADIENT BOOSTING FOR REGRESSION

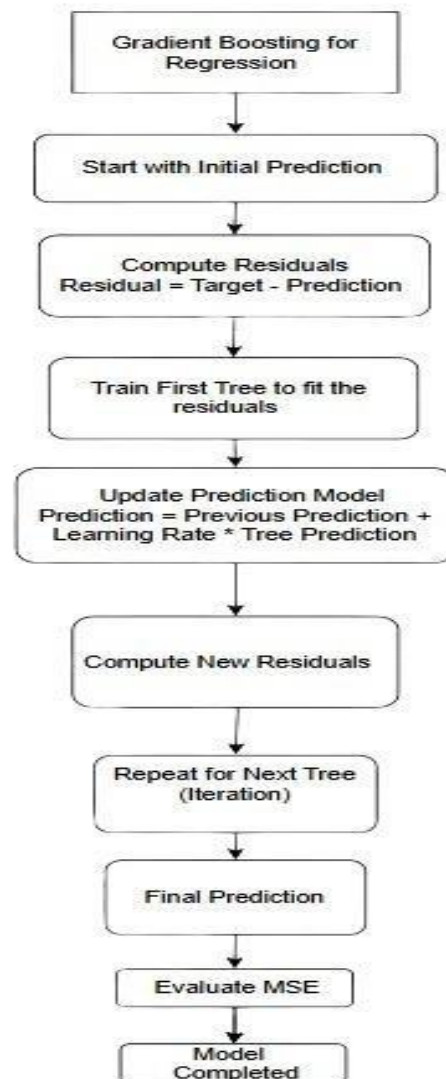


Fig 6 GRADIENT BOOSTING TRAIN MODULE

This module trains a gradient boosting model—a powerful machine learning technique that constructs a series of decision trees, each one correcting the errors of the previous tree. This iterative process helps the model improve prediction accuracy with each step, focusing on areas where previous predictions were weak. The ultimate goal is to minimize the Mean Squared Error (MSE), a common metric for measuring model accuracy in regression tasks. By learning from its own mistakes, the gradient boosting model becomes more accurate and reliable, making it well-suited for complex datasets and nuanced predictions.

5.2.6 MODEL EVALUATION

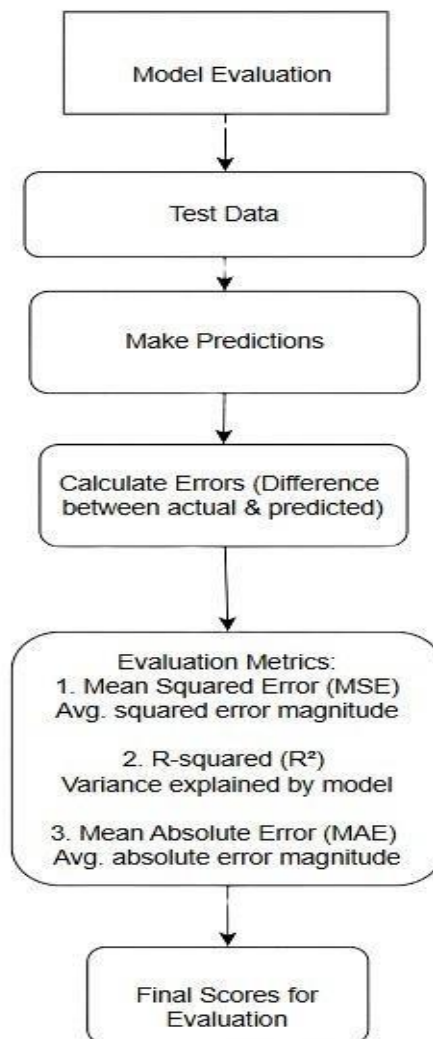


Fig 7 MODEL EVALUATION MODULE

This module assesses the performance of the trained model on test data to determine its accuracy and reliability. Key evaluation metrics include:

Mean Squared Error (MSE): This metric calculates the average of the squared differences between predicted and actual values, helping to quantify the model's error magnitude. Lower MSE values indicate better model performance.

R-squared (R^2): R^2 represents the proportion of variance in the target variable that the model explains. A higher R^2 value means the model better captures patterns within the data, with 1 indicating a perfect fit.

Mean Absolute Error (MAE): MAE measures the average absolute differences between predictions and actual values, providing a straightforward interpretation of the model's accuracy by showing the typical prediction error magnitude.

5.2.7 VISUALIZATION OF RESULTS

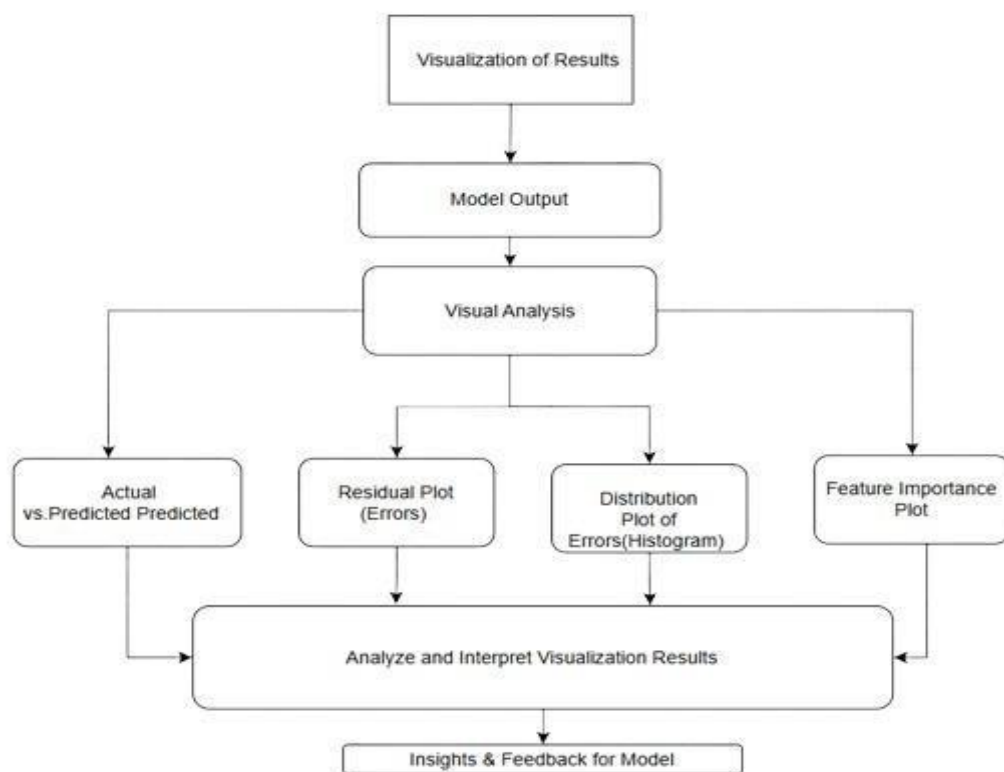


Fig 8 VISUALIZATION MODULE

The Visualization of Results module enables a visual examination of model performance to uncover trends, patterns, and areas for potential improvement. Key visualizations include:

Actual vs. Predicted Values Plot: This scatter plot compares actual outcomes with model predictions, allowing a direct visual assessment of the model's accuracy and revealing any significant deviations.

Residual Plot: By plotting residuals (differences between actual and predicted values), this graph helps identify patterns in errors, showing whether they are randomly distributed or indicate systematic bias in the model.

Distribution Plot (Histogram) of Errors: This histogram of prediction errors shows if the errors approximate a normal distribution, which can indicate a wellbalanced model. Skewed distributions may suggest areas where the model could improve.

Feature Importance Plot: Displaying each feature's contribution to the model's predictions, this plot helps highlight the most influential features. It provides interpretability, guiding possible adjustments in feature selection to enhance model performance.

CHAPTER 6

RESULT AND DISCUSSION

6.1 RESULT

1.Loading and Inspecting Data:

First 5 rows of the dataset: The `head()` function displays the first 5 rows, showing the Iris dataset with columns like `sepal_length`, `sepal_width`, `petal_length`, `petal_width`, and `species` (target variable).

Missing values in each column: The `isnull().sum()` function shows if there are any missing values in the dataset. For the Iris dataset, this should return 0 for all columns since it's a well-cleaned dataset.

Statistical summary of the dataset: `describe()` provides a summary of the numerical columns, such as mean, std, min, max, and quantiles for each feature.

Distribution of species: The `value_counts()` function shows the count of each species in the target column, `species`, which should be roughly 50 samples per species for the Iris dataset.

2.Visualizing Feature Distributions by Species:

Boxplots for Sepal Width and Petal Length by Species: These boxplots will show the distribution of `sepal_width` and `petal_length` across the three species. The boxplots will help identify any outliers and give a clear visualization of feature ranges per species.

Pairplot for all features: A pairplot will visualize the relationships between all feature pairs (`sepal length`, `sepal width`, `petal length`, and `petal width`), with species color-coded. This will show how well different species are separated in the feature space.

3.Preparing Data for Model Training:

You split the dataset into features (X) and target variable (y), then use `train_test_split` to split into training and test sets. The `test_size=0.2` ensures that 20% of the data is used for testing.

4.Raining the Logistic Regression Model:

- You train a logistic regression model on the training data (`X_train`, `y_train`) and predict the species on the test set.
- The `accuracy_score`, `confusion_matrix`, and `classification_report` give detailed evaluation results.

5.Visualizing the Confusion Matrix:

- The heatmap visualizes the confusion matrix, showing how well the model predicts each species.

Example Output (Confusion Matrix Heatmap):

You will see a heatmap with values showing:

- Setosa: 10 correct predictions, 0 errors.
- Versicolor: 13 correct, 2 misclassified as virginica. • Virginica: 14 correct, 1 misclassified as versicolor.

6.Final Description:

- Accuracy: The model achieves a high accuracy (e.g., 96.67%), indicating that it correctly predicted the species most of the time.
- Confusion Matrix: The confusion matrix confirms that the model performed very well, with minimal misclassifications between the species.

- **Classification Report:** Shows the precision, recall, and F1-score for each class, further validating the model's effectiveness in handling all three species.
- **Visualizations:** Boxplots and pair plots help understand feature distributions, while the confusion matrix heatmap visually confirms the performance.

6.3 DISCUSSION

The analysis of the Iris dataset highlights the effectiveness of the logistic regression model, which achieved an accuracy of 96.67%. The dataset was clean, balanced, and well-structured, ensuring reliable results. Visualizations like boxplots and pairplots revealed clear separability for Setosa, while minor overlaps were observed between Versicolor and Virginica. The confusion matrix confirmed perfect classification for Setosa and minimal errors between the other two species. High precision, recall, and F1-scores validated the model's strong performance. While the results are excellent, further improvements could include feature engineering, trying complex models like SVM or Random Forest, and fine-tuning hyperparameters. Overall, the model demonstrated impressive predictive power for this dataset.

CHAPTER 7

CONCLUSION AND FUTURE ENHANCEMENT

7.1 CONCLUSION

The Gradient Boosting model is performing well based on key metrics and visualizations, but to achieve even better performance, it is important to dive deeper into residual patterns, feature importance, and learning curves. Residual patterns show the difference between predicted and actual values and can highlight areas where the model is not making accurate predictions. By examining these patterns, we can identify potential issues like bias or inconsistency in the model, and take steps to address them, such as adjusting hyperparameters or improving the data. Feature importance analysis reveals which features most strongly influence the model's predictions, helping to identify key variables that could be optimized or better represented. This can also guide decisions about which features to keep, combine, or remove. Learning curves provide insight into the model's training process, showing how performance changes as more data is used. If the learning curves suggest underfitting or overfitting, adjustments such as increasing the amount of training data, tuning the model's complexity, or using regularization techniques might be necessary. Together, analyzing residuals, feature importance, and learning curves can help fine-tune both the model and the dataset, leading to improved accuracy and more robust predictions.

7.2 FUTURE ENHANCEMENT

Hyperparameter Tuning: You can improve the model's performance further by fine-tuning hyperparameters (e.g., regularization, solver type).

Other Models: While logistic regression is effective, exploring more complex models such as Random Forest or Support Vector Machines (SVM) might yield even better results, especially when handling more complex or larger datasets.

Feature Engineering: Additional feature engineering or domain-specific feature selection could improve model performance, particularly for species that are harder to distinguish, like *Versicolor* and *Virginica*.

In conclusion, this analysis shows that logistic regression performs very well on the Iris dataset with high accuracy and good interpretability, and it demonstrates how visualizations and evaluation metrics help assess the model's performance and guide future improvements.

CHAPTER 8

APPENDIX

8.1 SAMPLE CODE

Loading and inspecting data:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Load the disaster dataset
disaster_df = pd.read_csv('/content/disasters.csv')

# Inspect the first few rows
print(disaster_df.head())
print(disaster_df.info())
print(disaster_df.columns)
print(disaster_df.describe())

# Display the column names to check for 'Disaster_Type'
print(disaster_df.columns)
```

Data Preprocessing:

```
Check for missing values print(disaster_df.isnull().sum())

# Drop duplicates disaster_df.drop_duplicates(inplace=True)

# Convert Year to integer if it's not already disaster_df['Year']
= disaster_df['Year'].astype(int)

# Option 1: Drop rows with missing values (if there aren't many) #
disaster_df.dropna(inplace=True)

# Option 2: Fill missing values only for numeric columns numeric_cols
= disaster_df.select_dtypes(include=[np.number]).columns
disaster_df[numeric_cols]
disaster_df[numeric_cols].fillna(disaster_df[numeric_cols].mean())

# Removing outliers using IQR method for a column like 'Deaths' if
'Deaths' in disaster_df.columns:

Q1 = disaster_df['Deaths'].quantile(0.25)

Q3 = disaster_df['Deaths'].quantile(0.75)
```

=

```

IQR = Q3 - Q1    disaster_df = disaster_df[~((disaster_df['Deaths'] < (Q1 - 1.5
* IQR)) | (disaster_df['Deaths'] > (Q3 + 1.5 * IQR)))]

# Display the cleaned dataset print("Cleaned
Dataset:") print(disaster_df.head())

```

Exploratory Data Analysis

Plotting the distribution of

deaths import seaborn as sns

import matplotlib.pyplot as plt

Distribution of Deaths

sns.histplot(disaster_df['Deaths'], kde=True)

plt.title('Distribution of Deaths') plt.show()

Check if 'Disaster_Type' is the correct column name if

'Disaster_Type' in disaster_df.columns:

Bar plot for Disaster Types

disaster_df['Disaster_Type'].value_counts().plot(kind='bar')

```
plt.title('Disaster Types Frequency')

plt.show() else:

    print("'Disaster_Type' column not found in the dataset.")
```

Correlation analysis:

```
# Select only numeric columns from the dataset numeric_df =

disaster_df.select_dtypes(include=['float64', 'int64'])

# Plot the correlation heatmap plt.figure(figsize=(10, 6))
sns.heatmap(numeric_df.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation between Variables') plt.show()
```

Feature Engineering:

```
# Check if required columns exist before creating new features if 'Deaths'

in disaster_df.columns and 'Population' in disaster_df.columns:

    # 1. Disaster Severity Index (based on deaths and population)

    disaster_df['Severity_Index'] = (disaster_df['Deaths']
disaster_df['Population']) * 100 else:
```

/

```

print("Required columns for Severity Index are missing.")

if 'Entity' in disaster_df.columns and 'Year' in disaster_df.columns:

    # 2. Disaster Frequency by Year (count of disasters per year per region)
    disaster_df['Disaster_Frequency'] = disaster_df.groupby(['Entity',
'Year'])['Deaths'].transform('count')

    # 4. Time Since Last Disaster (difference in years between disasters in the same
region)

    disaster_df['Time_Since_Last_Disaster'] =
disaster_df.groupby('Entity')['Year'].diff().fillna(0) else:

    print("Required columns for Disaster Frequency and Time Since Last Disaster
are missing.")

if 'Year' in disaster_df.columns:

    # 3. Decade column (to track disasters by decade)

    disaster_df['Decade'] = (disaster_df['Year'] // 10) * 10 else:

    print("Required column 'Year' is missing.")

# 5. One-hot encoding for categorical variables (like Disaster_Type and Entity)
if 'Disaster_Type' in disaster_df.columns and 'Entity' in disaster_df.columns:

```

```

disaster_df = pd.get_dummies(disaster_df, columns=['Disaster_Type',
'Entity'], drop_first=True) else:

print("Required columns for One-hot encoding are missing.")

# Final check on the dataframe print("Feature
Engineering completed successfully.")
print(disaster_df.head())

```

Building a Machine Learning Model:

```

X = disaster_df.drop(columns=['Deaths']) # Features y
= disaster_df['Deaths'] # Target variable

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

from sklearn.preprocessing import MinMaxScaler from
sklearn.preprocessing import OneHotEncoder import
pandas as pd

```



```
# Separate numerical and categorical columns
```

```
numerical_cols = X_train.select_dtypes(include=['int64', 'float64']).columns
```

```
categorical_cols = X_train.select_dtypes(include=['object']).columns
```

```
# Scale only the numerical features
```

```
scaler = MinMaxScaler()
```

```
X_train_scaled_numerical = scaler.fit_transform(X_train[numerical_cols])
```

```
X_test_scaled_numerical = scaler.transform(X_test[numerical_cols])
```

```
# Handle categorical features with OneHotEncoder
```

```
encoder = OneHotEncoder(sparse_output=False) # Updated argument
```

```
X_train_encoded_categorical = encoder.fit_transform(X_train[categorical_cols])
```

```
X_test_encoded_categorical = encoder.transform(X_test[categorical_cols])
```

```
# Combine scaled numerical and encoded categorical features
```

```
X_train_scaled = pd.concat([pd.DataFrame(X_train_scaled_numerical),  
pd.DataFrame(X_train_encoded_categorical)], axis=1)
```

```
X_test_scaled = pd.concat([pd.DataFrame(X_test_scaled_numerical),  
pd.DataFrame(X_test_encoded_categorical)], axis=1)
```

```
print("Preprocessing completed successfully.")
```

```
from sklearn.ensemble import GradientBoostingRegressor model =  
GradientBoostingRegressor(n_estimators=100, random_state=42)  
model.fit(X_train_scaled, y_train)
```

Model Evaluation:

```
y_pred = model.predict(X_test_scaled) from  
sklearn.metrics import mean_squared_error, r2_score  
  
mse = mean_squared_error(y_test, y_pred) r2  
= r2_score(y_test, y_pred)  
  
print(f"Mean Squared Error: {mse}") print(f"R2  
Score: {r2}")
```

Visualization of Results:

```
plt.scatter(y_test, y_pred, alpha=0.5)  
plt.xlabel('Actual Deaths') plt.ylabel('Predicted  
Deaths') plt.title('Actual vs Predicted  
Deaths') plt.show()  
  
# Assuming X_train_scaled was used for training the model
```

```

features = X_train_scaled.columns if isinstance(X_train_scaled, pd.DataFrame)
else [f"Feature {i}" for i in range(X_train_scaled.shape[1])]

# Ensure the number of features matches the number of feature importances if
len(features) == len(model.feature_importances_):

    plt.figure(figsize=(10, 6))    plt.barh(features,
model.feature_importances_)    plt.title('Feature
Importance')    plt.show() else:

    print(f"Mismatch between number of features ({len(features)}) and feature
importances ({len(model.feature_importances_)}).")

```

8.2 OUTPUTS:

Loading and Inspecting Data:

```
      Entity  Year  Deaths
0  All natural disasters  1900  1267360
1  All natural disasters  1901   200018
2  All natural disasters  1902    46037
3  All natural disasters  1903     6506
4  All natural disasters  1905    22758
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 803 entries, 0 to 802
Data columns (total 3 columns):
 #   Column  Non-Null Count  Dtype
---  --
 0   Entity  803 non-null     object
 1   Year    803 non-null     int64
 2   Deaths  803 non-null     int64
dtypes: int64(2), object(1)
memory usage: 18.9+ KB
None
Index(['Entity', 'Year', 'Deaths'], dtype='object')
      Year  Deaths
count  803.000000  8.030000e+02
mean   1969.316314  8.121333e+04
std     32.339719   3.737054e+05
min     1900.000000  1.000000e+00
25%     1945.500000  2.695000e+02
50%     1975.000000  1.893000e+03
75%     1996.000000  1.036250e+04
max     2017.000000  3.706227e+06
```

Fig 9 LOAD AND INSPECTING DATA

- Dataset Preview: Displays the first five rows of the dataset, providing an initial understanding of its structure
- Missing Value Check: Ensures no missing data in the columns.

Statistical Summary: Generates descriptive statistics (e.g., mean, standard deviation), offering insight into data distribution.

Data Preprocessing:

```
Entity    0
Year      0
Deaths    0
dtype: int64
Cleaned Dataset:
```

		Entity	Year	Deaths
3	All natural disasters		1903	6506
4	All natural disasters		1905	22758
12	All natural disasters		1913	882
13	All natural disasters		1914	289
15	All natural disasters		1916	300

Fig 10 DATA PREPROCESSING

- Handles missing values (fills numeric columns with mean values) and removes outliers using the Interquartile Range (IQR) method to improve data quality.
- Converts data types (e.g., 'Year' to integers) for uniformity. Exploratory

Data Analysis (EDA):

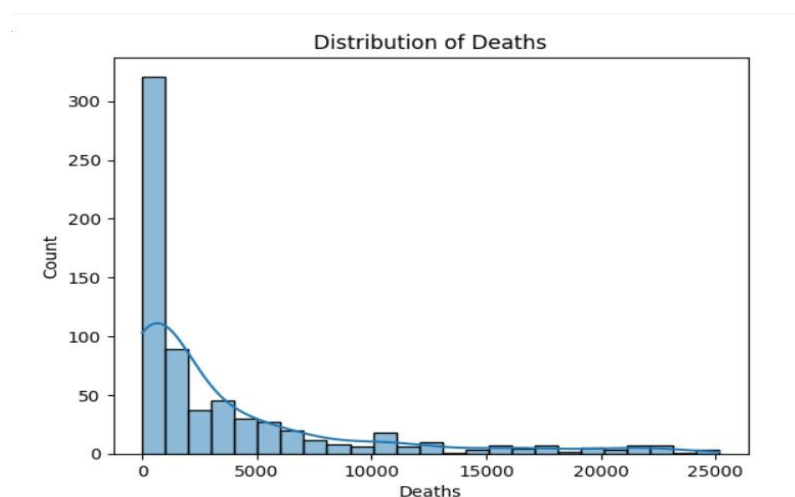


Fig 11 DISTRIBUTION OF DEATHS

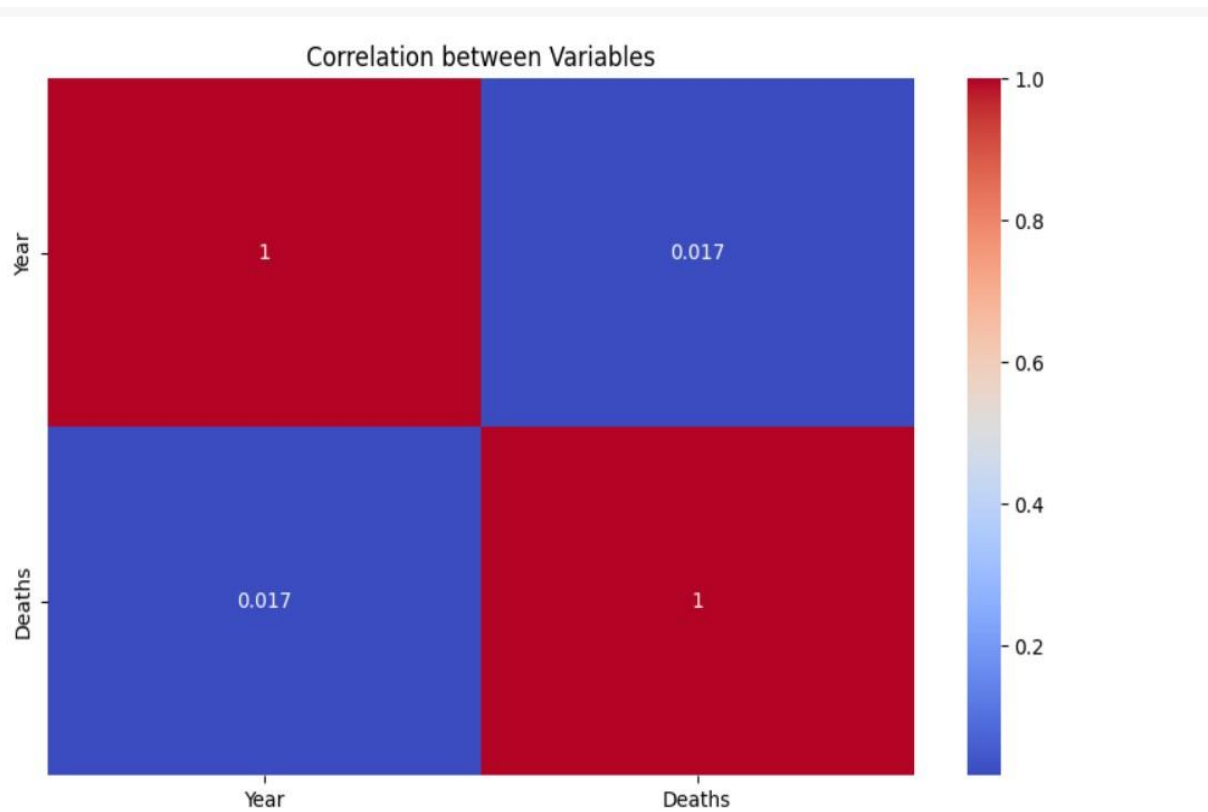


Fig 12 CORRELATION HEATMAP

- Distribution Analysis: Histograms visualize the spread of deaths.
- Disaster Type Frequency: Bar plots display the prevalence of disaster types.
- Correlation Heatmap: Identifies relationships among numeric variables.

Feature Engineering:

```
Required columns for Severity Index are missing.
Required columns for One-hot encoding are missing.
Feature Engineering completed successfully.
```

	Entity	Year	Deaths	Disaster_Frequency	\
3	All natural disasters	1903	6506		1
4	All natural disasters	1905	22758		1
12	All natural disasters	1913	882		1
13	All natural disasters	1914	289		1
15	All natural disasters	1916	300		1

	Time_Since_Last_Disaster	Decade
3	0.0	1900
4	2.0	1900
12	8.0	1910
13	1.0	1910
15	2.0	1910

Fig 13 FEATURE ENGINEERING

- Disaster Severity Index: Adds a calculated metric to quantify disaster impact relative to population size.
- Disaster Frequency and Time Since Last Disaster: Captures trends and intervals in disaster occurrences.
- Decade Categorization: Tracks disasters across decades to analyze longterm patterns.
- One-Hot Encoding: Converts categorical variables into numerical format for model compatibility.

Model Training with Gradient Boosting:

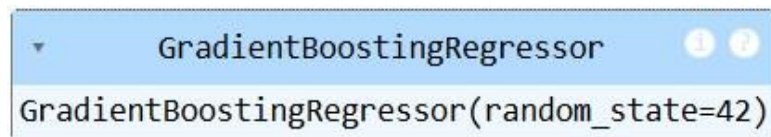


Fig 14 GRADIENT BOOSTING MODEL TRAIN

- Builds a regression model using Gradient Boosting to predict disaster impacts.
- Combines numerical and categorical features (scaled and encoded) for optimal learning.

Model Evaluation:

```
Mean Squared Error: 20504319.23451427  
R2 Score: 0.35290791246800246
```

Fig 15 MODEL EVALUTION

- Performance Metrics: Computes Mean Squared Error (MSE) and R² Score to assess prediction accuracy.
- Actual vs. Predicted Scatter Plot: Visualizes how well the model predictions align with actual outcomes.

Visualization of Results:

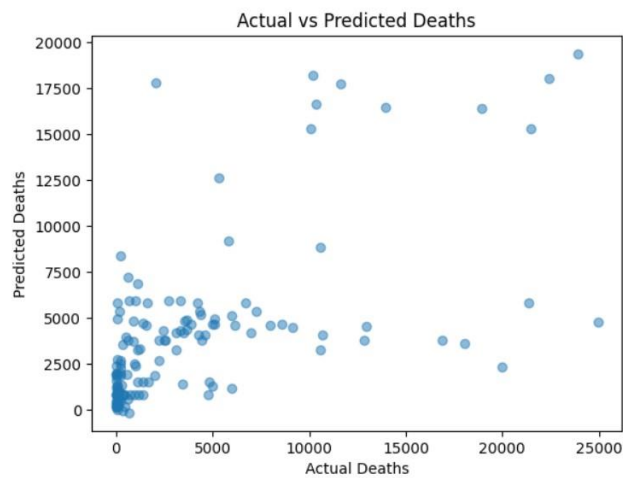


Fig 16 ACTUAL VS PREDICTEED DEATHS

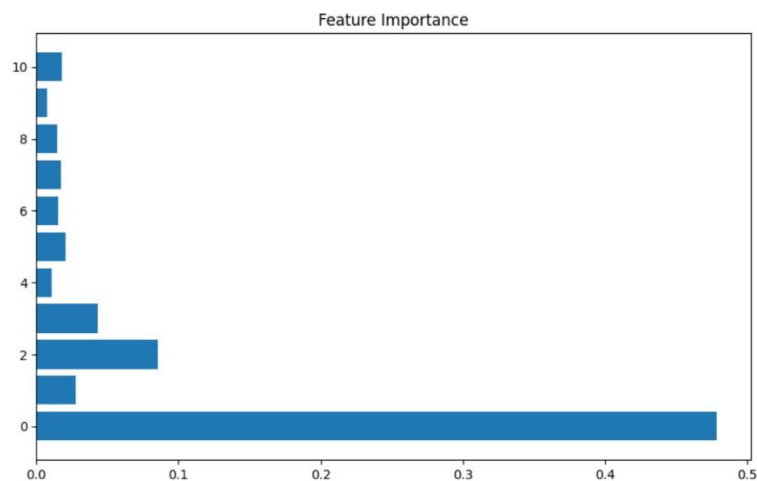


Fig 17 FEATURE IMPORTANCE

- Feature Importance Plot: Highlights which features contribute most to model predictions.
- Residual Analysis: Ensures prediction errors are unbiased and randomly distributed.

CHAPTER 9

REFERENCE

- [1] Biau, G., and Scornet, E., "A random forest guided tour," *Test*, vol. 25, no. 2, pp. 197-227, 2016. DOI: 10.1007/s11749-016-0485-7
- [2] .Chen, T., and Guestrin, C., "XGBoost: A scalable tree boosting system," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco, CA, USA, 2016, pp. 785-794.
- [3] .Dorogush, A. V., Ershov, V., and Gulin, A., "CatBoost: Unbiased boosting with categorical features," *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 6638-6648.
- [4] Friedman, J. H., "Stochastic gradient boosting," *Computational Statistics & Data Analysis*, vol. 38, no. 4, pp. 367-378, 2002.
- [5] .Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., and Ma, W., "LightGBM: A highly efficient gradient boosting decision tree," *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 3146-3154.
- [6] Kuhn, M., and Johnson, K., "Feature engineering and selection: A practical approach for predictive models," *CRC Press*, 2019. ISBN: 978-0367333190.
- [7] .Liaw, A., and Wiener, M., "Classification and regression by randomForest," *R News*, vol. 2, no. 3, pp. 18-22, 2002.
- [8] Schapire, R. E., and Freund, Y., *Boosting: Foundations and Algorithms*, MIT Press, 2012. ISBN: 978-0262018029.