

Abstract

This dissertation presents a comprehensive exploration of Geo Spatial Data Analytics applied to Natural Disasters, integrating advanced Machine Learning techniques for predictive insights. Investigating diverse geospatial patterns, our study leverages longitude and latitude data to predict disaster types with precision. Through systematic data collection, preprocessing, and rigorous algorithmic training, this research contributes to enhanced disaster management strategies. The findings underscore the pivotal role of geospatial analysis and machine learning in bolstering disaster preparedness and response.

Contents

Declaration of original work	1
Acknowledgements	2
Abstract	3
List of Figures	7
List of Tables	9
1 Introduction	10
1.1 Background and Context	10
1.2 Research Problem and Research Questions	11
1.3 Significance and Justification	12
1.4 Scope and Limitations	13
1.4.1 Scope	13
1.4.2 Limitations	14
2 Literature Review	15
2.1 The Fusion of Geospatial Data and Disaster Risk Assessment	15
2.2 Machine Learning's Rising Role in Disaster Risk Management	16
2.3 Operationalizing Geospatial Data and Machine Learning in Disaster-Prone Regions	16
2.4 A Paradigm Shift in Disaster Risk Management	16
3 Methodology	18
3.1 Research Design and Approach	18

3.2	Data Collection and Preprocessing Overview	19
3.3	Machine Learning Approach	20
3.4	Code Implementation Environment	20
4	Data Collection and Preprocessing	22
4.1	Geocoded Disasters (GDIS) Dataset	22
4.2	Emergency Events Database (EM-DAT)	23
4.3	Data Merging: GDIS and EM-DAT	24
4.4	Conclusion	25
5	Machine Learning Models	26
5.1	Correlation Analysis of Geographic Features	26
5.2	Selection of Machine Learning Models	27
5.3	Final Model and Disaster Type Prediction	28
5.3.1	Architecting the Ultimate: Final Random Forest Classifier	29
5.3.2	From Coordinates to Insights: predict_disaster_type Function	30
5.3.3	A Glimpse into the Future: Applying the Function	30
5.4	Understanding the Random Forest Classifier	30
5.4.1	Ensemble Techniques	31
5.4.2	Decision Trees: Random Forest's Foundation	31
5.4.3	Gini Index	32
5.4.4	The Implementation of Decision Trees in the Random Forest Algorithm	36
6	Result	38
6.1	Predictive and Visual Analysis: predict_and_visualize_disaster()	38
6.2	Disaster Type Ranking: A Pragmatic Approach	40
6.3	Validation in the Real World: A Case for Model Efficacy	41
7	Visualization of Geo-Spatial Data	43
7.1	Global Disasters Overview	43

7.2	Temporal Patterns	44
7.3	Disaster Type Analysis	45
7.4	Affected Population by Country	46
7.5	Impact and Fatality Analysis	46
7.6	Geographic Distribution	47
7.7	Spatial Concentration	47
7.8	Earthquake Mapping	48
7.9	Temporal Evolution	48
7.10	Country-Specific Analysis: India and Indonesia	49
7.10.1	Indonesia	51
7.10.2	India	52
7.11	Regional Analysis: Top adm1 (Region/Province) for All Disasters	53
7.12	City Analysis: Top 10 Locations for all Disasters	54
8	Discussion	55
9	Conclusion	60
9.1	Summary of Findings	60
9.2	Conclusion Statement	61
9.3	Recommendations	61
9.4	Reflection and Self-Evaluation	62
Bibliography		63

List of Figures

4.1	df_gdis head	23
4.2	df_emdat head	23
4.3	df head	24
5.1	Disaster Ranking by Prediction probability Code	29
5.2	https://wiki.pathmind.com/decision-tree	32
5.3	Tree Diagrams	33
5.4	Tree Diagrams	34
5.5	https://www.section.io/engineering-education/introduction-to-random-forest-in-machine-learning/	36
6.1	Prediction and Visulization	39
7.1	Global Disaster Count by Country	44
7.2	Chart of Disasters over Time	45
7.3	Bar Chart of Disasters by Type	45
7.4	Visualisation of Total Affected Population by Country	46
7.5	Choropleth Map of Total Global Deaths	46
7.6	Geographic Distribution	47
7.7	3D Representation of the Geographic Concentration of a Specific Disaster Type	48
7.8	Earthquake Mapping	48
7.9	Animated Scattergeo Plot of Specific Disasters over Years	49
7.10	Animation of Specific Disasters over Years in Different Countries	49

7.11 Disasters in Indonesia	51
7.12 Spatial Visualisation of Different Types of Disasters in Indonesia	51
7.13 Animated Scattergeo Plot Illustrating Disaster Locations over Time in Indonesia	52
7.14 Interactive Scattergeo Plot of Spatial Distribution of Different Disasters	52
7.15 Animated Scattergeo Plot Depicting Disaster Locations in India over Time	53
7.16 Total Disasters: India Vs Indonasia	53
7.17 Graph of Top 10 Global Provinces and Taxes Scatterplot	54
7.18 Top 10 Locations	54

List of Tables

5.1	A 3x5 Table	33
5.2	A 3x6 Table	34

Chapter 1

Introduction

The dissertation combines geographical data, machine learning methodologies, and visualisation approaches to enhance disaster risk evaluation, planning, and administration. This dissertation offers explicit and informative visual depictions and analyses of worldwide disaster patterns, encompassing a thorough chronology and emphasising several incidences impacting India and Indonesia. The research incorporates an advanced machine learning model to predict disaster classes by leveraging geographical coordinates. These insights play a significant role in the allocation of resources, the advancement of urban infrastructure, and the enhancement of early warning systems. Furthermore, the inquiry offers diverse visual representations depicting disasters across different countries, facilitating a more comprehensive comprehension of the locations most vulnerable to worldwide calamities. Consequently, this fosters enhanced global cooperation and streamlined norms for decision-making.

1.1 Background and Context

The dissertation is situated within a framework characterised by an increased demand for a thorough understanding and proactive forecasting of natural calamities amidst a global environment that is progressively getting more volatile. The pressing nature of this requirement underscores the importance of conducting a systematic analysis of geospatial data on diverse categories of natural calamities observed worldwide. The

incorporation of geographical characteristics and dynamic environmental factors is crucial in the determination of disaster probability. Therefore, developing machine learning models integrating these factors is of utmost significance. The dissertation aligns with the prevailing tendency of integrating data science and machine learning frameworks into the evaluation, readiness, and management strategies for disaster risk. The research is particularly noteworthy due to its capacity to generate precise predictions of natural disasters, successfully conveyed through cartographic depictions and artistic images. Python modules, particularly Plotly, are of great importance in this undertaking as they enable the creation of interactive visualisations that aid in detecting spatial dispersion and occurrence rates of catastrophic occurrences.

1.2 Research Problem and Research Questions

This dissertation centres on the intricate issue of effectively forecasting and assessing the incidences and trends of disasters using geographical variables. Contemporary techniques necessitate the consideration of the dynamic repercussions stemming from climatic fluctuations, population movements, and technological progressions, all of which possess the potential to influence the probability of natural calamities.

The following research questions guide the study:

1. In what manner may machine learning methodologies facilitate the production of precise predictions for hazard categorisations by using geographical data, precisely latitude and longitude?
2. To maintain the accuracy of disaster forecasts, it is imperative to continually update predictive models with the most recent data and account for changing climatic circumstances.
3. How might using data visualisation techniques contribute to developing efficiently understandable representations of these forecasts, improving catastrophe preparedness and management efforts?

4. What are this attempt's potential impacts on assessing and managing disaster risks, urban planning initiatives, international aid partnerships, and policy formulation?
5. How can the outcomes of this project substantially contribute to the advancement of informed infrastructure initiatives, decision-making processes for resource allocation, and the communication of information pertaining to measures for disaster prevention?

1.3 Significance and Justification

The dissertation can significantly transform the assessment of disaster risks, the distribution of resources, and the management practises by employing machine learning approaches. This dissertation offers significant insights into proactive disaster preparedness and response strategies by correctly forecasting the types of disasters that may occur based on regional variables. The value of this predictive capability lies in its ability to enhance evacuation tactics, maximise resource allocation, and enable the development of timely early warning systems. As a result, it can save lives and minimise losses.

The motivation behind this project stems from the imperative to comprehensively understand and effectively respond to a climate exhibiting escalating levels of instability. The multifaceted and wide-ranging impacts of natural disasters, which present risks to human lives, infrastructure, and economies, underscore the importance of forecasting and preparedness measures. Furthermore, the legitimacy of the research stems from its utilisation of up-to-date data and its recognition of the environment's ever-changing nature. This facilitates the ability of researchers to generate precise and prompt predictions regarding disasters, hence enhancing the efficiency of initiatives related to disaster management.

Furthermore, this dissertation holds substantial significance for the field of urban planning and the advancement of infrastructure. The research aids authorities in identifying places that are susceptible to natural catastrophes and enhances their understanding of

the accompanying hazards. This is achieved through the provision of invaluable insights that serve as a guide for decision-making processes. Consequently, these observations aid in the development and execution of infrastructure that is more resistant to adverse events, thereby successfully reducing the potential losses associated with disasters. Furthermore, integrating data visualisation tools facilitates the effective communication of intricate geospatial and disaster-related data, hence augmenting the practical applicability of the research. Visual representations facilitate the practical application of information for various stakeholders involved in disaster management, encompassing government agencies, non-government groups, and local populations. This enables the cultivation of well-informed decision-making and fosters collaborative endeavours to bolster disaster preparedness and response strategies.

1.4 Scope and Limitations

1.4.1 Scope

This dissertation encompasses various topics by incorporating machine learning techniques and data visualisation strategies into analysing a substantial geographical dataset. The primary aim of this system is to predict trends and vulnerabilities related to catastrophes by utilising geographic coordinates. This information can be highly beneficial in informing essential areas such as emergency management, urban planning, infrastructure development, and policy formulation. The research's global applicability is noteworthy since it utilises latitude and longitude data to make predictions that may be applied to diverse geographical places worldwide. Furthermore, the study aims to generate quantifiable societal impacts beyond technical factors. This initiative aims to enhance community awareness, foster active engagement in resilience-building efforts, and facilitate international collaboration in mitigating disaster risks. Furthermore, the study establishes the fundamental basis for forthcoming investigations and advancements in forecasting and administration of disaster risks.

1.4.2 Limitations

Despite its comprehensive breadth, the research is bound by many noteworthy limitations requiring acknowledgement. The sole dependence on geographic coordinates for disaster predictions represents a significant constraint. Although the rationale behind this approach may be to minimise potential inaccuracies in data, it may fail to consider other critical location-specific variables that could impact the incidence of disasters. Furthermore, the effectiveness of the research is contingent upon the calibre and promptness of the input data employed for modelling, hence creating intrinsic constraints on the accuracy and reliability of the data. Although the research considers dynamic factors such as climate change and population dynamics, it may need to consider additional regional or cultural elements that could impact disaster planning and response. The research's scope should specifically encompass these factors. As a result, the predictions may offer a limited representation of the diverse characteristics of disaster events in particular areas, underscoring the need for further investigation to comprehend these intricate processes better.

Chapter 2

Literature Review

This chapter thoroughly examines the complex interplay between catastrophe risk assessment and management and the transformative capabilities of geospatial data and machine learning techniques. As the frequency and impact of natural disasters escalate, the imperative for robust assessment and management of disaster risk grows ever more pronounced. This chapter delves into the seminal work that underscores the potential of geospatial data and machine learning to enhance our capabilities in disaster risk assessment and response.

2.1 The Fusion of Geospatial Data and Disaster Risk Assessment

The evaluation and mitigation of disaster risk have been brought to new horizons through the infusion of geospatial data. Researchers have harnessed the intrinsic spatial dimension to foster an encompassing comprehension of the distribution of vulnerabilities and available resources (Harilal, 2021). Geospatial data encompasses an array of sources including satellite imagery, GPS traces, mobile Call Detail Records (CDRs), and social media posts (Rezaei, 2021). This multitude of data sources empowers an expedited and effective disaster response, enriching our understanding of the disaster landscape.

2.2 Machine Learning’s Rising Role in Disaster Risk Management

The amalgamation of machine learning methodologies with disaster risk management is ushering in a paradigm shift. These techniques enable holistic evaluations of hazards, exposure, and vulnerabilities, transcending the boundaries of natural and manmade disasters (Ekeanyanwu et al., 2022). Researchers have ingeniously combined data mining strategies, econometric regression models, and input-output models to generate comprehensive insights (Gao et al., 2022). The integration of various data dimensions, from geographical conditions to historical disaster data, augments the precision of risk estimation, resource allocation, and task prioritization.

2.3 Operationalizing Geospatial Data and Machine Learning in Disaster-Prone Regions

The application of geospatial data and machine learning resonates strongly in disaster-prone regions, as exemplified by studies in Indonesia and India. A profound illustration is the utilization of deep neural networks to detect submerged infrastructure in flooded areas (Pu, 2017). Through innovative algorithms, the study identified submerged stop signs and estimated floodwater depth, aiding in targeted rescue and recovery efforts. These endeavors shed light on the potency of these techniques in fostering community preparedness and resilience in the face of adversity.

2.4 A Paradigm Shift in Disaster Risk Management

Scholarly discourse has consistently underscored the transformative potential of geospatial data and machine learning in revolutionizing disaster risk management. In the Philippines, sophisticated models leveraging geographical data and automatic learning algorithms predict areas susceptible to collapse (Anonymous, n.d.). The fusion of satellite imagery, digital elevation models, and advanced learning techniques leads to predictive precision that holds the promise of reshaping disaster risk assessment methodologies.

In summation, this chapter accentuates the vanguard role of geospatial data and machine learning in disaster risk assessment and management. By harnessing their synergistic potential, the research community propels disaster management strategies toward innovation, resilience, and adaptability in the face of an ever-changing natural landscape.

Chapter 3

Methodology

This chapter outlines the comprehensive methodology employed in the pursuit of geospatial data analytics and disaster-type prediction through machine learning. The synthesis of data analysis and predictive modeling forms the bedrock of this study's approach, fostering a holistic understanding of the complex interactions between geographical variables and natural disasters. A detailed exploration of the research design, data collection, preprocessing procedures, and machine learning techniques underscores the systematic manner in which insights were derived and predictive models were developed. This chapter serves as a roadmap that guides readers through the intricate processes, offering a clear perspective on the analytical framework that underpins the subsequent phases of our investigation.

3.1 Research Design and Approach

The research approach and methods employed in this study integrate data analytics and machine learning in a mutually beneficial manner, aiming to uncover valuable insights into geospatial data and natural disasters. The primary aim is to leverage computational methodologies to effectively harness the capabilities of identifying patterns, interconnections, and predictive models that significantly enhance the field of disaster-type forecasting. This research is situated at the convergence of geography, data science, and predictive modeling, aiming to generate practical knowledge that can be utilized to

improve disaster management and mitigation strategies.

The technique consists of discrete phases, each contributing differently to attaining the study's goals. The foundation of geospatial analysis lies in the collecting and pre-processing of data, encompassing the curation and refinement of raw geospatial data. After integrating machine learning methodologies, developing predictive models that can classify disaster types by utilizing geographical coordinates becomes feasible. It is essential to highlight that this approach uses a cyclical procedure typified by iterative assessment and enhancement to guarantee the strength and precision of the models.

The research design exhibits a solid commitment to maintaining methodological rigor by emphasising the significance of transparency, reproducibility, and integrating multiple data sources.

3.2 Data Collection and Preprocessing Overview

This work utilizes two comprehensive datasets, namely the Geocoded Disasters (GDIS) Dataset and the Emergency Events Database (EM-DAT) databases. The GDIS Dataset, which is an expansion of the EM-DAT and has been carefully curated by the Centre for Research on the Epidemiology of Catastrophes (CRED), comprises a comprehensive collection of 39,953 geocoded locations that are interconnected with records documenting 9,924 natural catastrophes that occurred between the years 1960 and 2018. The dataset includes various hazard categories: floods, cyclones, earthquakes, landslides, droughts, volcanic activities, and severe temperature events. It offers valuable insights for our analytical pursuits.

The GDIS Dataset provides spatial resolution principally at the first level of administration (state/province/region), supplemented by certain places at the third level of administration (district/commune/village) through a connection to the Global Administrative Areas database (GADM, 2018). The dataset, referred to as "df_gdis" in our

analytical framework, enables geographical data collection.

3.3 Machine Learning Approach

The heart of this study lies in the seamless fusion of geospatial data analytics and machine learning techniques to predict disaster types based on geographical coordinates. The machine learning approach adopted for this purpose entails a selection of well-established algorithms, each tailored to capture distinct patterns and relationships within the data.

The primary focus was on algorithmic diversity to ensure a comprehensive exploration of potential predictive models. Four key machine learning algorithms were chosen for evaluation: Logistic Regression, K-Nearest Neighbors Classifier, Decision Tree Classifier, and Random Forest Classifier. These algorithms were selected based on their suitability for classification tasks and their varied approaches to discerning underlying patterns.

The rationale behind algorithm selection drew from both their performance characteristics and their interpretability, a crucial aspect given the domain's high-stakes applications. The intent was to balance the predictive accuracy of the models with the ability to derive actionable insights from the results.

Further detailing of these algorithms and their performance evaluation is provided in Chapter 5, where we delve into the specifics of each model's architecture, training, and testing procedures. Through this approach, we ensure a thorough exploration of algorithmic strengths and limitations in the context of disaster-type prediction.

3.4 Code Implementation Environment

The execution of data analysis and machine learning endeavors transpired within a robust coding environment, aligning with the principles of reproducibility and transparency. Python, a versatile and widely adopted programming language, served as the cornerstone of our codebase. The integration of Python enabled efficient data manipulation,

statistical analysis, and seamless integration of machine learning libraries.

The interactive nature of Jupyter Notebook further facilitated the research process, fostering an environment where code, visualizations, and narrative text converge in a coherent narrative. This integration enabled a step-by-step exploration of data pre-processing, feature engineering, algorithm implementation, and results interpretation. Each Jupyter Notebook encapsulated a specific phase of the research journey, enhancing clarity, documentation, and ease of collaboration.

The utilization of Python and Jupyter Notebook aligns with the broader open science ethos, ensuring that our methodologies and findings can be readily accessed, scrutinized, and reproduced by fellow researchers. Through this approach, the study embraces a commitment to robust research practices and accessibility.

In the forthcoming chapters, we shall introduce the code samples and methodologies delineated within the Jupyter Notebook framework. These lectures aim to offer exemplary instances of the analytical and machine learning methodologies employed in this research.

Chapter 4

Data Collection and Preprocessing

This chapter delves into the fundamental elements of our project, "Geospatial Data Analytics of Natural Disasters," focusing on the crucial stages of data gathering and preparation. To thoroughly examine natural disasters, it is imperative to have access to accurate and extensive datasets that encompass geographical information and specialized data on the tragedies. We have carefully selected and organized two significant datasets to achieve this objective, each serving a distinct purpose in providing valuable insights for our research attempt.

4.1 Geocoded Disasters (GDIS) Dataset

The primary dataset utilized in our research project, titled "Geospatial Data Analytics of Natural Disasters," is the Geocoded Disasters (GDIS) Dataset, version 1 (1960-2018). This dataset serves as a fundamental component, offering a comprehensive compilation of natural catastrophes that transpired globally within the designated timeframe. The GDIS Dataset is a significant resource for spatial analysis and exploration since it offers geocoded data for around 39,953 distinct sites closely linked to approximately 9,944 thoroughly recorded disasters.

Problems Encountered with the GDIS Dataset

Even with the GDIS Dataset's provision of a robust basis for geographical investigations,

	id	country	iso3	gwno	year	geo_id	geolocation	level	adm1	adm2	adm3	location	historical	hist_country	disastertype	disasterno	latitu
0	109	Albania	ALB	339.0	2009	346	Ana E Malit	3	Shkoder	Shkodres	Ana E Malit	Ana E Malit	0	NaN	flood	2009-0631	42.0209
1	109	Albania	ALB	339.0	2009	351	Bushat	3	Shkoder	Shkodres	Bushat	Bushat	0	NaN	flood	2009-0631	41.9592
2	175	Angola	AGO	540.0	2001	760	Onjiva	3	Cunene	Cuanhama	Onjiva	Onjiva	0	NaN	flood	2001-0146	-17.0934
3	187	Angola	AGO	540.0	2009	710	Evale	3	Cunene	Cuanhama	Evale	Evale	0	NaN	flood	2009-0092	-16.5315
4	187	Angola	AGO	540.0	2009	749	Mupa	3	Cunene	Cuvelai	Mupa	Mupa	0	NaN	flood	2009-0092	-16.2000

Figure 4.1: df_gdis head

it is imperative to acknowledge its inherent constraints. The dataset's significant reliance on geographical coordinates as the sole source of disaster-related information is a considerable concern. The GDIS Dataset's limited availability of comprehensive disaster information required a purposeful approach to augment its analytical depth. The implementation of this strategy action necessitated the integration of supplementary data to enhance and expand its informational scope.

4.2 Emergency Events Database (EM-DAT)

The Emergency Events Database (EM-DAT) significantly supplements the GDIS Dataset. Since the year 1900 up until the current time, EM-DAT has functioned as a comprehensive database containing more than 26,000 records documenting significant occurrences of large-scale catastrophes. The EM-DAT database offers extensive disaster information from reputable entities like UN agencies, research institutes, and non-governmental groups. The compilation encompasses data on the number of casualties, estimations of the extent of devastation, and analysis of the economic ramifications.

Dis No	Year	Seq	Glide	Disaster Group	Disaster Subgroup	Disaster Type	Disaster Subtype	Disaster Subsubtype	Event Name	...	Reconstruction Costs, Adjusted ('000 US\$)	Insured Damages ('000 US\$)	Insured Damages, Adjusted ('000 US\$)	Total Damages in US Dollars	Total Damages, Adjusted ('000 US\$)	
0	1960-0013-CHL	1960	13	NaN	Natural	Geophysical	Earthquake	Tsunami	NaN	NaN	...	NaN	NaN	550000.0	5442439.0	1
1	1960-0026-AIA	1960	26	Nan	Natural	Meteorological	Storm	Tropical cyclone	NaN	Donna	...	NaN	NaN	35000.0	346337.0	1
2	1960-0025-ANT	1960	25	Nan	Natural	Meteorological	Storm	Tropical cyclone	NaN	Donna	...	NaN	NaN	NaN	NaN	NaN
3	1960-0024-ATG	1960	24	Nan	Natural	Meteorological	Storm	Tropical cyclone	NaN	Donna	...	NaN	NaN	NaN	NaN	NaN
4	1960-0030-BGD	1960	30	Nan	Natural	Meteorological	Storm	Tropical cyclone	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN

Figure 4.2: df_emdat head

Problems Encountered with the EM-DAT Dataset

Despite its comprehensive nature, the EM-DAT Dataset presents certain limitations that impact its suitability for our spatial analysis goals. One notable constraint is the absence of geocoded information for individual disasters. This absence hampers the ability to directly correlate disasters with specific geographic locations, a critical aspect of our research focus.

4.3 Data Merging: GDIS and EM-DAT

In pursuit of a comprehensive dataset that marries geographic insights with disaster-specific details, a merging process was undertaken. By strategically employing the "disasterno" variable, we orchestrated the integration of GDIS and EM-DAT datasets. The resulting amalgamation combines the strengths of both datasets, encompassing latitude, longitude, and vital disaster particulars.

	Dis No	Disaster Group	Disaster Subgroup	Disaster Subtype	Region	Continent	Dis Mag Value	Dis Mag Scale	Total Deaths	No Injured	...	adm1	adm2	adm3	location	historical	hist_c
0	1960-0013	Natural	Geophysical	Tsunami	South America	Americas	10.0	Richter	6000.0	3000.0	...	Los Lagos	Llanquihue	Puerto Montt	Puerto Montt	0	
1	1960-0013	Natural	Geophysical	Tsunami	South America	Americas	10.0	Richter	6000.0	3000.0	...	Los Rios	Valdivia	NaN	Valdivia	0	
2	1960-0030	Natural	Meteorological	Tropical cyclone	Southern Asia	Asia	NaN	Kph	3000.0	NaN	...	Chittagong	Noakhali	NaN	Noakhali	1	P
3	1960-0031	Natural	Meteorological	Tropical cyclone	Southern Asia	Asia	NaN	Kph	5149.0	NaN	...	Chittagong	NaN	NaN	Chittagong area	1	P
4	1960-0033	Natural	Geophysical	Ground movement	Northern Africa	Africa	6.0	Richter	57.0	NaN	...	M'Sila	NaN	NaN	Beni Illmane	1	

Figure 4.3: df head

Problems Encountered with the Merged Dataset

While the merged dataset showcases improved completeness, it is essential to acknowledge certain limitations. The merged dataset constitutes a subset of the original GDIS data, limiting the scope of analysis. Furthermore, the geographic coverage of the merged dataset is constrained relative to the original GDIS dataset and EM-DAT, necessitating a nuanced approach to interpretation.

4.4 Conclusion

Why I Chose the Geocoded Disasters (GDIS) Dataset

The selection of the GDIS Dataset as a priority is based on its alignment with the project's core purpose, which is to examine spatial phenomena. The primary focus of the GDIS Dataset lies in the inclusion of geographical coordinates, which facilitates accurate mapping and correlation between occurrences of disasters and distinct geographical locations. The EM-DAT Dataset, while its comprehensive coverage of disasters, could be more applicable for spatial analysis due to the absence of geocoded data. After conducting a thorough assessment of the available datasets, the GDIS Dataset was identified as the most suitable basis for fulfilling the geographic study goals of our project.

This structured outline comprehensively addresses the Geocoded Disasters (GDIS) Dataset, the Emergency Events Database (EM-DAT), their integration, and the rationale behind dataset selection. It emphasizes the insights gained and challenges encountered during the data collection and preprocessing phases of your Geo-Spatial Data Analytics of Natural Disasters project.

Chapter 5

Machine Learning Models

The machine learning aspect of this endeavor is to utilize geographic coordinates to forecast different categories of calamities. Machine learning methodologies primarily aim to discern associations between geographic areas and the incidence of specific mishaps. This chapter utilizes state-of-the-art algorithms to offer in-depth disaster insights that surpass a superficial examination. Integrating data-driven predictive models with spatial information presents an opportunity to improve catastrophe management tactics and facilitate informed decision-making. This chapter will provide an overview of the technique utilized, selecting appropriate machine learning models, and the persuasive implications of employing data-driven insights for disaster prediction. This study aims to enhance our knowledge of the correlation between geographical characteristics and the incidence of disasters, thereby bolstering our capabilities in this field.

5.1 Correlation Analysis of Geographic Features

Within this particular section, an extensive examination of the association between geographic parameters and the dependent variable "disaster type" is undertaken. This study aims to determine the level of correlation between each attribute and the occurrence of different types of natural catastrophes. The correlation values generated in this study reveal the potential impact of various geographical factors on the event of diverse disasters.

The research demonstrates that some attributes, namely "location," "id," and "iso3," exhibit relatively strong correlation coefficients with the dependent variable. Nevertheless, these attributes may only sometimes hold the utmost importance when forecasting the various categories of natural calamities. The term "location" is deemed unsuitable for accurate predictions due to the potential lack of representation of all conceivable locations within the dataset. In disaster prediction, the identifiers "id" and "iso3" possess limited direct interpretability.

To tackle this matter, attention was redirected towards factors strongly linked to geographical attributes, precisely "latitude" and "longitude." A moderate association was seen between these two qualities and the dependent variable. In contrast to the term "location," the attributes of "latitude" and "longitude" are continuous numerical variables that offer a more precise representation of a site's geographical coordinates. Given the high degree of specificity exhibited, these individuals possess exceptional qualifications for enhancing the machine learning model and making predictions regarding various types of disasters through geographic data.

5.2 Selection of Machine Learning Models

The careful selection of a suitable machine learning model is paramount to the overall success of the undertaking. This section explores various methodologies for predicting the occurrence of different sorts of disasters using geographical coordinates. The correctness of each model is evaluated to assess its ability to capture the inherent patterns present within the dataset effectively.

The study encompassed assessments of Logistic Regression, K-Nearest Neighbors Classifier, Decision Tree Classifier, and two variants of Random Forest Classifier. The training and evaluation of each model involved utilizing the "latitude" and "longitude" features within the dataset to make predictions on the corresponding "disastertype."

The performance of the models was quantified using accuracy as the evaluation metric. The accuracy scores obtained for each model were as follows:

- **Logistic Regression Model:** Accuracy - 0.4374921787010387
- **K-Nearest Neighbors Classifier Model:** Accuracy - 0.5816543611563009
- **Decision Tree Classifier Model:** Accuracy - 0.6119384307345764
- **Random Forest Classifier Model:** Accuracy - 0.6221999749718433 (Highest Accuracy)
- **Random Forest Classifier (cross_val_score):** Average Accuracy - 0.40680

The Random Forest Classifier had the highest level of accuracy across all the models examined, achieving a value of 0.6221999749718433. The model effectively discerned the complex relationships between geographic coordinates and disaster categories, showcasing its aptness for this objective.

The Random Forest Classifier emerged as the most effective model for this task based on the obtained accuracy scores. Its robustness in capturing intricate patterns in the geographic data makes it a reliable choice for predicting disaster types based on latitude and longitude coordinates. The following sections delve into the mathematical foundations and operation of the Random Forest Classifier, shedding light on its predictive capabilities.

5.3 Final Model and Disaster Type Prediction

The culmination of our journey through machine learning models leads us to the pinnacle—the embodiment of predictive prowess: the final Random Forest Classifier. In this section, we delve into the intricate architecture of this model and its transformative capacity to predict disaster types based on geographic coordinates.

```

# Create the Random Forest classifier
model = RandomForestClassifier(n_estimators=100, random_state=42)

# Fit the model to the data
model.fit(X_encoded, y)

def predict_and_visualize_disaster(latitude, longitude):
    # Function to predict the disaster type based on latitude and longitude
    def predict_disaster_type(latitude, longitude):
        # Create a dataframe with the input latitude and longitude
        input_data = pd.DataFrame({'latitude': [latitude], 'longitude': [longitude]})

        # Encode the input data
        input_encoded = pd.get_dummies(input_data)

        # Make the prediction
        prediction_proba = model.predict_proba(input_encoded)[0]
        disaster_types = model.classes_

        # Get the ranking based on probability
        ranking = sorted(zip(disaster_types, prediction_proba), key=lambda x: x[1], reverse=True)

        return ranking

    # Get the ranking of predicted disaster types
    predicted_ranking = predict_disaster_type(latitude, longitude)

    # Print the ranking
    for rank, (disaster_type, probability) in enumerate(predicted_ranking, 1):
        print(f'Rank {rank}: {disaster_type} (Probability: {probability})')

    # Function to zoom to the specified location on an interactive map
    def zoom_to_location(latitude, longitude):

        # Get the top 3 predicted disaster types and their probabilities
        predicted_ranking_2 = predict_disaster_type(latitude, longitude)[:2]

        # Create a color map based on the top 3 disaster types
        color_map = {
            'flood': 'skyblue',
            'earthquake': 'brown',
            'storm': 'gray',
            'extreme temperature': 'red',
            'landslide': 'darkbrown',
            'volcanic activity': 'orange',
            'drought': 'yellow',
            'mass movement (dry)': 'lightbrown'
        }

        # Create a Folium map centered at the specified location
        map_location = folium.Map(location=[latitude, longitude], zoom_start=8)

        # Create a HeatMap-Like overlay covering the district of the specified location with more random points
        heat_data = [
            [random.uniform(latitude - 0.1, latitude + 0.1), random.uniform(longitude - 0.1, longitude + 0.1)]
            for _ in range(500)
        ]
        HeatMap(
            heat_data,
            radius=10,
            gradients=[0.03: color_map[predicted_ranking_2[1][0]], 0.05: color_map[predicted_ranking_2[0][0]]]
        ).add_to(map_location)

        # Display the map
        display(map_location)

        # Call the function to zoom to the specified location
        zoom_to_location(latitude, longitude)

    # Example usage
    input_latitude = 28.644800
    input_longitude = 77.216721

    # Call the combined function to predict and visualize disaster for the specified location
    predict_and_visualize_disaster(input_latitude, input_longitude)

```

Figure 5.1: Disaster Ranking by Prediction probability Code

5.3.1 Architecting the Ultimate: Final Random Forest Classifier

With unwavering dedication, our research culminates in the establishment of the final Random Forest Classifier—an ensemble of decision trees harnessed to predict disaster types. The finesse of this model goes beyond mere complexity; it thrives on the aggregation of individual tree predictions, culminating in a harmonious orchestration that excels in accuracy and reliability. Anchored by a careful selection of estimators—100 in our case—the final Random Forest Classifier stands poised to unveil insights from geographical coordinates.

5.3.2 From Coordinates to Insights: predict_disaster_type Function

The pinnacle of our predictive journey is marked by the innovative predict_disaster_type function. This function encapsulates the essence of our endeavors, translating raw coordinates into actionable intelligence. As the function springs into action, it initiates a symphony of operations: encoding the input data, predicting disaster probabilities, and ranking these types based on likelihood. With every invocation of this function, a microcosm of analysis unfolds, bridging the geographical divide to offer invaluable insights into impending disaster scenarios.

5.3.3 A Glimpse into the Future: Applying the Function

Envision the scenario—a specific latitude and longitude beckoning for predictions. We set forth, feeding these coordinates into our predict_disaster_type function. The outcome transcends data points; it's a prognosis, a roadmap into potential disasters. The function's output isn't confined to mere numbers; it's a ranked hierarchy of disaster types, each accompanied by a probability score. This glimpse into the future arms stakeholders with actionable intelligence, enabling swift response, targeted resource allocation, and preemptive disaster management.

In this pivotal section, we cross the threshold from theory to tangible outcomes. The final model, fortified by meticulous training, emerges as a sentinel against uncertainty, while the predict_disaster_type function stands as a beacon, illuminating the uncharted territory of disaster prediction. Together, they encapsulate the culmination of our research—a journey that transforms geographical coordinates into a tapestry of insights, ready to reshape disaster risk assessment and management on a global scale.

5.4 Understanding the Random Forest Classifier

Random forests are a popular supervised Machine Learning methodology extensively employed in regression and classification tasks. Notably, they often exhibit exceptional performance even without hyperparameter adjustment. The method in question is widely

favoured due to its inherent simplicity. In the context of a classification problem, a common approach involves the generation of many decision trees using diverse samples. Subsequently, a voting mechanism is employed to determine the final classification based on the most prevalent outcome among the decision trees.

5.4.1 Ensemble Techniques

Ensemble techniques primarily involve the amalgamation of many models. Consequently, instead of relying on a singular model for prediction generation, a compilation of models is utilised, augmenting the total performance. There are two main ensemble techniques commonly employed in the field of Machine Learning:

1. **Bagging** – Bagging combats the limitations of relying on a single model by creating diverse subsets of the original dataset through bootstrapping. Bootstrapping involves repeatedly sampling the dataset with replacement, resulting in different data points for each model. In the Random Forest Classifier context, this translates to constructing multiple decision trees, each trained on a distinct subset. The final prediction is then made by aggregating the individual predictions through a majority vote in classification tasks or averaging in regression tasks
2. **Boosting** – Boosting is a sequential ensemble technique that focuses on correcting the errors made by preceding models. Unlike Bagging, where models are independent, Boosting creates a series of models where each model attempts to improve the shortcomings of its predecessors. The final model's prediction is generated through a weighted combination of these sequential models.

5.4.2 Decision Trees: Random Forest's Foundation

To understand the functioning of a random forest algorithm, it is imperative first to grasp the concept of Decision Trees. Decision Trees are an alternative method employed in Supervised Machine Learning for classification and regression problems.

Decision trees are graphical representations that resemble a tree structure and illustrate predictions made by a sequence of divisions depending on features. The process commences with a primary node and concludes with a terminal selection.

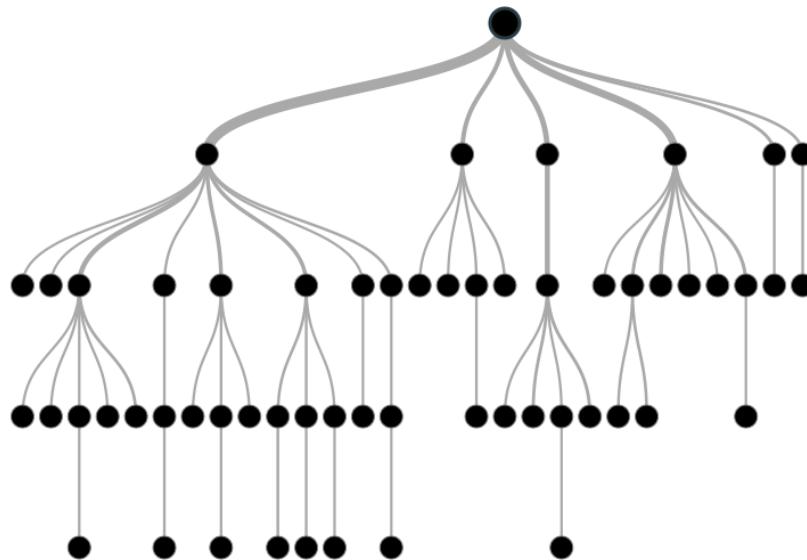


Figure 5.2: <https://wiki.pathmind.com/decision-tree>

The root node, the decision node, and the leaf node are the three separate parts of the structure. The root node acts as the starting point for population subdivision. The nodes that result from splitting a root node are known as decision nodes, whereas the node that cannot be split further is known as a leaf node. Determining which property should be the root node is consequently a difficult task. It is interesting to choose the root node from a dataset with several properties. We must comprehend the "Gini Index" in order to respond to this query.

5.4.3 Gini Index

The ability to discern whether a feature's division will be contaminated or uncontaminated can aid in deciding which parts should be further subdivided. A pure sub-split is characterized by producing binary outcomes expressly limited to "yes" or "no" responses. Suppose this is our dataset.

Feature 1	Feature 2	Feature 3
2	7	Yes
3	3	Yes
6	5	No
4	1	Yes

Table 5.1: A 3x5 Table

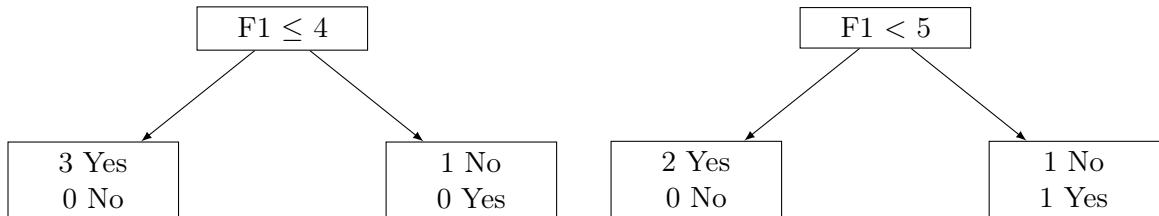


Figure 5.3: Tree Diagrams

When feature 1 is selected as the root node, a pure split is obtained, while selecting feature 2 results in an asymmetrical division. Therefore, it is necessary to ascertain the impurity level of the particular node. Utilizing the "Gini Index" will facilitate your understanding of this concept.

It is imperative to ascertain the impurity of our dataset, which will serve as the basis for identifying the root node with the lowest impurity, specifically through the utilization of the Gini index. Mathematically, the Gini index can be written as:

$$\begin{aligned} \text{Gini Index} &= 1 - \sum_{i=1}^n (P_i)^2 \\ &= 1 - [(P_+)^2 + (P_-)^2] \end{aligned} \tag{5.1}$$

P_+ is the probability of a positive class, and P_- is the probability of a harmful category.

ID	Loan Amount	Loan Status
1	100	Bad
2	200	Good
3	250	Bad
4	400	Good
5	300	Bad

Table 5.2: A 3x6 Table

Using Loan Amount as the root node:

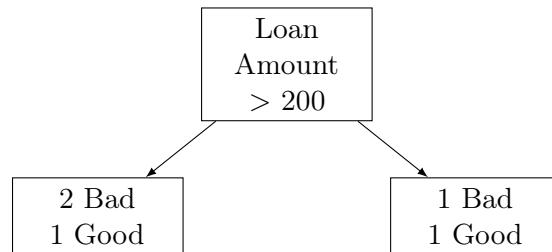


Figure 5.4: Tree Diagrams

using a left split's values in the formula

$$\begin{aligned}
 \text{Gini Index} &= 1 - \sum_{i=1}^n (P_i)^2 \\
 &= 1 - [(P_+)^2 + (P_-)^2] \\
 &= 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] \\
 &= 1 - [0.1089 + 0.4356] \\
 &= 1 - 0.5445 \\
 &= 0.4555
 \end{aligned}$$

Now using a Right split's values in the formula

$$\begin{aligned}
 \text{Gini Index} &= 1 - \sum_{i=1}^n (P_i)^2 \\
 &= 1 - [(P_+)^2 + (P_-)^2]
 \end{aligned}$$

$$\begin{aligned}
&= 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right] \\
&= 1 - [0.25 + 0.25] \\
&= 1 - 0.5 \\
&= 0.5
\end{aligned}$$

To compute the weighted Gini index, which is the total Gini index of the split:

$$\begin{aligned}
\text{Weighted Gini Index} &= \frac{3}{5} \cdot 0.4555 + \frac{2}{5} \cdot 0.5 \\
&= 0.6 \cdot 0.4555 + 0.4 \cdot 0.5 \\
&= 0.2733 + 0.2 \\
&= 0.4733
\end{aligned}$$

Similarly, this approach will attempt to determine the Gini index for all conceivable divisions and select the root node feature with the lowest Gini index. The Gini index with the lowest value indicates a low impurity level.

The "Entropy" idea refers to a statistical metric used to assess the impurity level in a partition. Entropy may be calculated mathematically as follows:

$$E(S) = -p_{(+)} \log p_{(+)} - p_{(-)} \log p_{(-)}$$

The Gini index is often used since it does not include a logarithmic term like entropy, which results in speedier calculations. Logarithmic calculations are distinguished by their time-consuming nature. The Gini index is a frequent parameter in several boosting methods. The Gini Index will attain its highest possible value when a node contains an equal number of instances from both classes, suggesting a significant impurity level inside the node.

5.4.4 The Implementation of Decision Trees in the Random Forest Algorithm

The primary distinction between these two is that Random Forest is a bagging technique that makes predictions using a portion of the original dataset; this feature of Random Forest helps to avoid Overfitting. Random forest creates many decision trees (DTs) using various sets of observations rather than just one. This approach has the significant benefit of being applicable to both classification and regression issues.

Random Forest Algorithm Procedures

Step 1 – First, we create subgroups of our initial data. Row sampling and feature sampling will be used to build subsets of the training dataset by choosing rows and columns with replacements.

Step 2 – We create a separate decision tree for every single subset we take

Step 3 – An output will be provided by each decision tree.

Step 4 – The final result is evaluated using either the average or majority voting, depending on whether it is a classification problem or a regression problem.

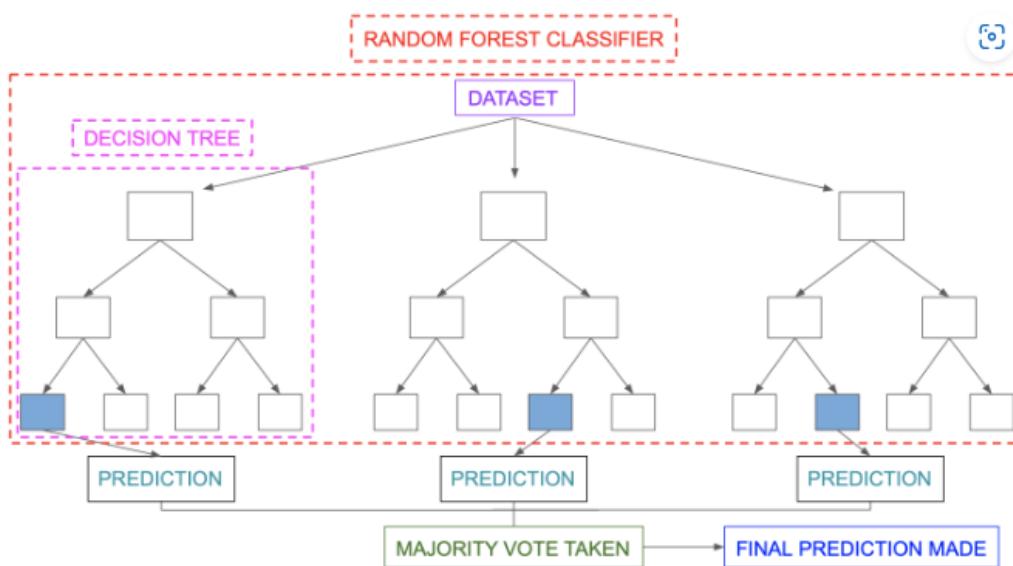


Figure 5.5: Image Source: <https://www.section.io/engineering-education/introduction-to-random-forest-in-machine-learning/>

When we extend our decision tree to full depth, we obtain Low Bias and High Variance. This means that while our model will excel on our training dataset, it will fail when our new data point enters the scene. In order to deal with this high variance condition, we employ random forest, merging many DTs rather than relying just on one. This will help us to reduce our variance and solve our overfitting issue.

Chapter 6

Result

In the spotlight of Chapter 6, we illuminate the outcomes of our research journey, delving into the practical application of our predictive model. Unveiling the power of disaster-type ranking and the model's real-world efficacy, this chapter exposes the intricacies of our predictions. Through interactive visual analyses and validation against ground realities, we bridge the gap between theory and application, showcasing the transformative potential of our approach to disaster risk assessment and management.

6.1 Predictive and Visual Analysis:

`predict_and_visualize_disaster()`

This section embarks on an exploration of the `predict_and_visualize_disaster()` function—a cornerstone of our research that converges predictive insights and dynamic visualizations. By employing geographical coordinates as inputs, this function unveils a two-fold revelation: disaster-type prediction and its intensity visualization through an interactive map overlay.

Applying the `predict_and_visualize_disaster()` Function

The function's mechanism is initiated with latitude and longitude inputs, representing geographic coordinates. In this process, our model springs into action, invoking predictive algorithms to discern probable disaster types. The resulting predictions are

further amplified through ranking, revealing a hierarchy of potential disasters based on calculated probabilities.

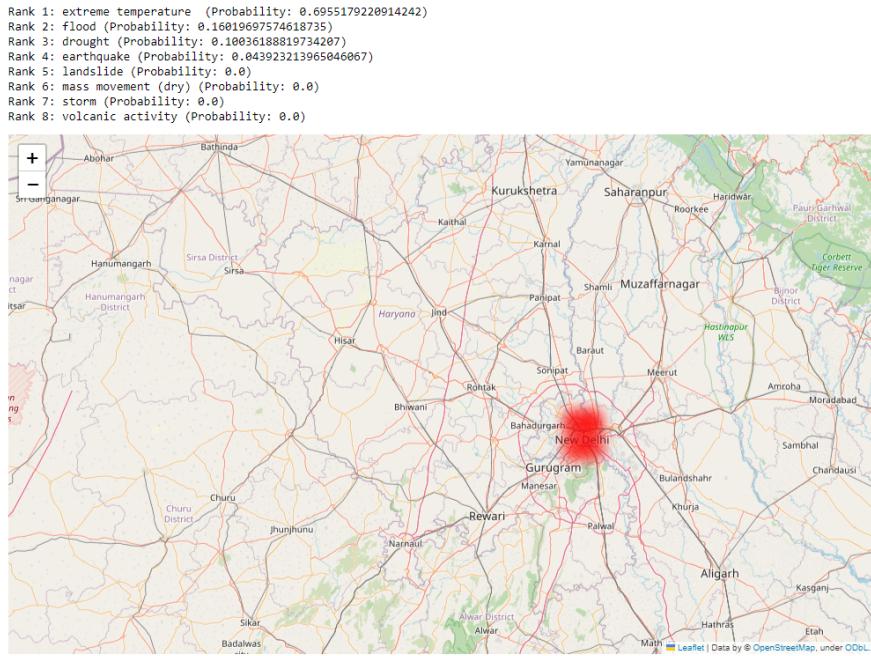


Figure 6.1: Prediction and Visualization

Visualizing Disaster Intensity Through HeatMap Overlay

The predictive prowess of the `predict_and_visualize_disaster()` function materializes in an interactive map visualization. By superimposing a HeatMap overlay, this approach visually encapsulates the intensity and spatial distribution of potential disasters. The map becomes an intuitive insight portal, engaging stakeholders with a tangible perspective on disaster vulnerability and enabling informed decision-making.

In this intricate interplay of prediction and visualization, the `predict_and_visualize_disaster()` function bridges theoretical algorithms and practical insights, embodying the potential to reshape disaster risk assessment dynamically and engagingly.

6.2 Disaster Type Ranking: A Pragmatic Approach

In this section, we delve into the significance of disaster-type ranking—a pragmatic approach that augments the depth and accuracy of our predictive model. By exploring the rationale behind this ranking methodology, we uncover its relevance in capturing the nuanced dynamics of geographical vulnerability and disaster frequency.

Why Disaster Type Ranking Matters

The choice to embrace disaster-type ranking goes beyond predictive accuracy—it aligns with the diverse geospatial dynamics of disaster-prone regions. From the Himalayan range to urban centers, the distribution of potential disasters varies extensively. This approach acknowledges this variability, catering to distinct vulnerabilities by offering a ranked hierarchy of disaster types. Whether it's a monsoon-induced flood or a seismic tremor, our methodology adapts to regional realities, enhancing preparedness and response strategies.

Disaster Ranking as an Amplifier of Accuracy

Within the realm of machine learning, a single accuracy score can belie the intricacies of real-world predictions. Disaster-type ranking leverages probability-based insights to provide a more comprehensive evaluation. By acknowledging that not all locations are equally vulnerable to each disaster, we gain a refined perspective on accuracy. This adaptive ranking ensures that predictions resonate with the local context, fostering a more precise understanding of potential threats and response strategies.

In adopting a pragmatic disaster-type ranking approach, we navigate beyond uniformity and embrace the diverse tapestry of disasters. This methodology not only bolsters prediction accuracy but also resonates with the real-world intricacies of geographical vulnerability and disaster prevalence.

6.3 Validation in the Real World: A Case for Model Efficacy

This section immerses us in the practical realm, where our predictive model faces the crucible of real-world observations. By aligning our predictions with ground-level realities, we affirm the efficacy of our approach and its potential to reshape disaster risk assessment and management strategies.

Model Validation Beyond Data: Real-world Observations

Transcending the confines of data, we embark on a validation journey that juxtaposes our model predictions with tangible ground-level scenarios. By selecting a range of cities, including well-known urban centers and disaster-prone regions, we assess how our model fares against actual occurrences. This exercise unveils the true potential of our predictive framework beyond the realm of numbers, grounding it in the empirical reality of disaster occurrences.

Mapping Model Predictions to Ground Reality

Our endeavor to validate the model leads us to a captivating discovery—our predictions mirror the real world. Cities experiencing recurrent extreme temperature fluctuations align with our model's insights. Urban centers consistently grappling with drought find resonance in our predictions. Himalayan regions, where flooding and seismic activity prevail, converge seamlessly with our model's priority ranking. This alignment underscores the remarkable congruence between predictive algorithms and the ground truth.

Insights and Findings in Alignment with Expectations

As we traverse this validation journey, we unearth a narrative of validation, where machine learning insights and human experience harmonize. Our model encapsulates not just the art of prediction, but also the wisdom of observation. It bridges the gap between scientific algorithms and tangible realities, embodying the potential to elevate disaster risk assessment to unprecedented heights of precision and preparedness.

In this chapter's final stride, we bear witness to the triumphant validation of our

model, a validation that extends beyond data points and converges with the lived experience of disasters. Through this harmonious convergence, our approach materializes as a beacon of efficacy—a transformative tool poised to navigate the ever-evolving landscape of disaster risk management.

Chapter 7

Visualization of Geo-Spatial Data

In this chapter, we employ sophisticated data visualization techniques to derive meaningful insights from extensive geospatial datasets. The armament at our disposal comprises the Plotly, Pydeck, and Folium libraries, which have been carefully chosen for their versatility and interactive capabilities. By incorporating 3D globe representations, Plotly enhances the realm of interactive visualizations, enriching our geographical studies with a sense of depth. The Pydeck framework is based on the Deck.gl foundation facilitates the development of captivating maps and visualizations by utilizing high-level application programming interfaces (APIs). Folium is widely regarded as the premier tool for generating interactive leaflet maps that incorporate annotations, pop-ups, and dynamic components. Utilizing this collaborative endeavor, we integrate geospatial data into compelling and enlightening formats, augmenting comprehension of disaster incidents and patterns across diverse geographical regions.

I suggest referring to my Python file "Data_Analytics" to thoroughly demonstrate the Geo Maps and Graphs discussed below to enhance the visual representation.

7.1 Global Disasters Overview

The visual representation offers a compelling display of global aggregate figures for various disasters across multiple nations. Upon executing the given code, a collection

of choropleth maps was generated, depicting the prevalence of certain natural disasters across various countries. The color spectrum employed in this study is based on the logarithm (\log_{10}) of the count, where deeper shades correspond to higher numbers. Hover text is utilized to provide more contextual information, such as the country's name and its ranking concerning the frequency of the specified disaster category.



Figure 7.1: Global Disaster Count by Country

7.2 Temporal Patterns

Through the visualisation analysis process, many insights were acquired about the temporal distribution of different categories of disasters. As illustrated in the picture provided, there have been fluctuations in the intensity of natural disasters throughout

history. Identifying these observable inclinations or regularities contributes to a more comprehensive comprehension of changing dynamics and variations of disaster events. The capacity to analyse and comprehend diverse trends and patterns across time is a handy tool in planning, decision-making, and emergency management. Using animated lines as a visual representation of data offers a valuable approach for efficiently analysing and interpreting differences across different classifications of natural catastrophes. This technique effectively aids in the identification and comprehension of distinctions, hence enhancing the process of searching for and analysing relevant information.

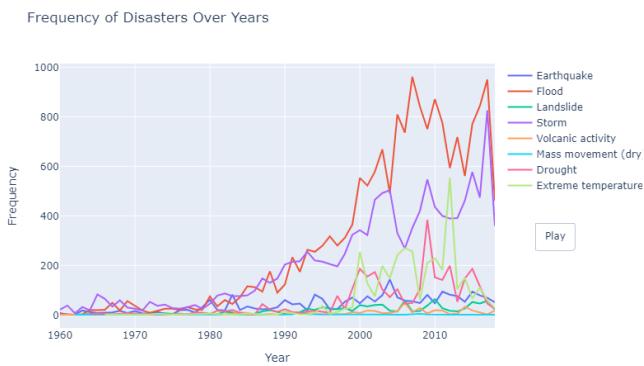


Figure 7.2: Chart of Disasters over Time

7.3 Disaster Type Analysis

This bar graph visually represents the distribution of disasters by category. The code processes a data file containing catastrophe-related information categorized within a "type of catastrophe" column, followed by a list of events corresponding to each type of catastrophe.

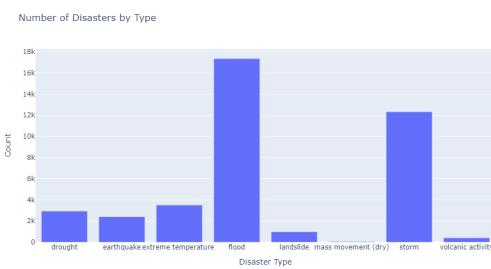


Figure 7.3: Bar Chart of Disasters by Type

7.4 Affected Population by Country

In the context of visualizing disaster impact, the choropleth map takes center stage, accentuating the gravity of the situation through a logarithmically scaled representation of the affected population. This visual narrative is woven through a meticulous fusion of data and code, where a "Total Affected" column finds its place within the DataFrame df. The result is a dynamic portrayal that amplifies the comprehension of disaster impact in a nuanced and informative manner.

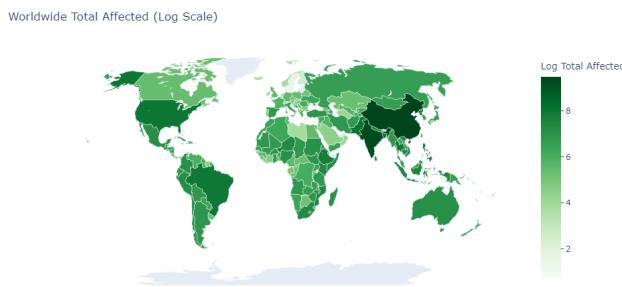


Figure 7.4: Visualisation of Total Affected Population by Country

7.5 Impact and Fatality Analysis

A choropleth map portraying the global distribution of fatalities caused by disasters emerges as the focal point. This function's input entails a DataFrame named "df," encapsulating vital "Total Deaths" column data alongside comprehensive disaster and associated country information.

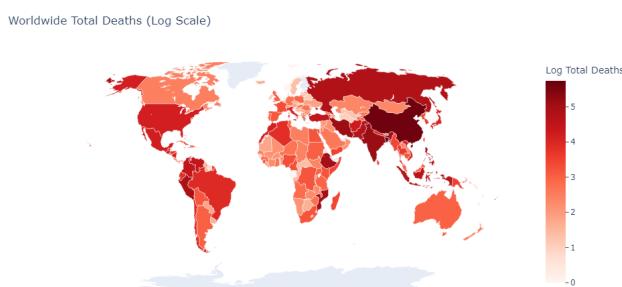


Figure 7.5: Choropleth Map of Total Global Deaths

7.6 Geographic Distribution

This Folium portrayal breathes life into a map, vividly illustrating the global expanse of various disasters using their latitude and longitude coordinates. On this canvas, each disaster is artfully represented by a CircleMarker, a symbol of its geographical presence and underlying importance.

```
color_palette = {'earthquake': 'red', 'storm': 'blue', 'flood': 'green', 'landslide': 'orange',  
'volcanic activity': 'purple', 'drought': 'yellow', 'mass movement (dry)': 'brown', 'extreme  
temperature': 'pink'}
```

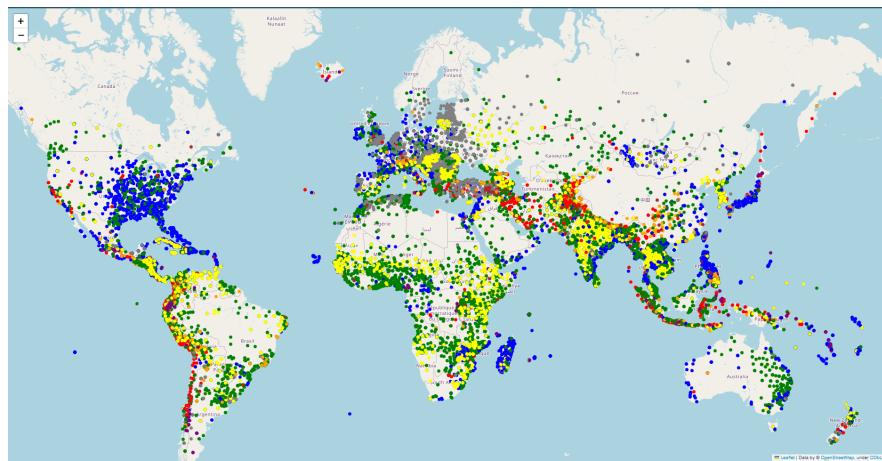


Figure 7.6: Geographic Distribution

7.7 Spatial Concentration

This 3D representation, dynamically visualized through the 'generate_disaster_plot()' function, illuminates the landscape with pinpointed disaster locations. Tailored to each catastrophe, the function integrates disaster names, filtering the DataFrame 'df_gdis' to isolate relevant data. With immersive control, users can rotate this 3D globe, zooming in and out, unraveling intricate disaster distributions across the world's canvas.

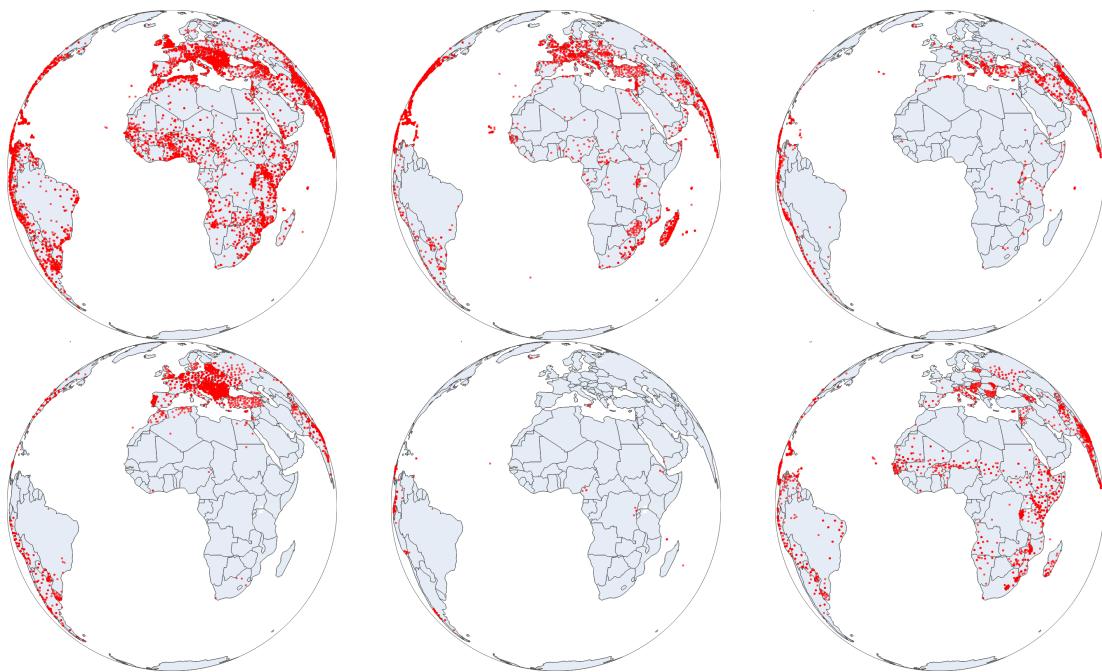


Figure 7.7: 3D Representation of the Geographic Concentration of a Specific Disaster Type

7.8 Earthquake Mapping

Graphically illustrate the specific locations of earthquake occurrences by applying a filter to the dataset (df) to exclusively showcase these instances.

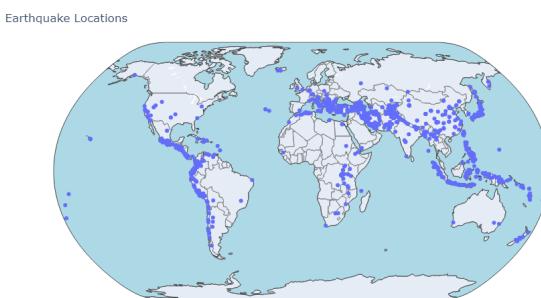


Figure 7.8: Earthquake Mapping

7.9 Temporal Evolution

This Animated scatterplot brings forth a dynamic visualization using Plotly Express, revealing the evolution of specific emergency occurrences across time. Within the

code, the 'plot_disaster_locations' function gracefully accepts DataFrame data and a designated catastrophe type as inputs.

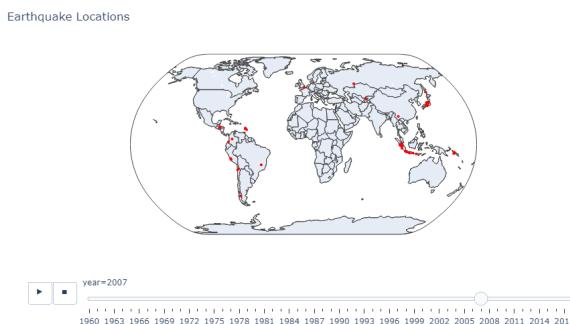


Figure 7.9: Animated Scattergeo Plot of Specific Disasters over Years

This Animated color map vividly portrays the annual occurrence of a specific natural disaster type across countries. Orchestrated by the plot_yearly_disaster function, this visualization springs to life through DataFrame data and catastrophe-type inputs.

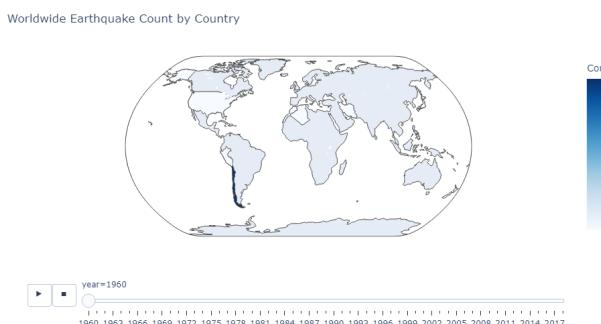


Figure 7.10: Animation of Specific Disasters over Years in Different Countries

7.10 Country-Specific Analysis: India and Indonesia

Based on the findings derived from the Global Disasters Overview, which emphasized the increased occurrence of earthquakes, landslides, and volcanic activity in Indonesia, designating it as a prominent location for such calamities and recognizing India's heightened vulnerability to weather-related disasters such as droughts, floods, and extreme temperatures, defining it as a significant site for such disasters, our decision was made to develop spatial visualizations representing these events.

Possible Reason:

- Indonesia
 1. **Geological Context:** Indonesia's geographical positioning within the Pacific Ring of Fire, known for its high levels of tectonic activity, results in regular seismic events. This region is characterized by the convergence of many tectonic plates increasing seismic and volcanic activity frequency.
 2. **Volcanic Activity:** Indonesia possesses a significant quantity of active volcanoes due to its geographical positioning inside the Pacific Ring of Fire. A considerable number of volcanoes enhances the likelihood of volcanic eruptions.
 3. **Tropical Climate:** Indonesia's susceptibility to landslides is attributed to its tropical environment, characterized by copious precipitation, particularly in regions with rugged topography. The convergence of precipitation, topographical elevation, and porous soil augments the probability of landslides.
- India
 1. **Diversified Geography:** The geographical features of India encompass a diverse range of terrains, including plains, mountains, deserts, and coastal regions, which give rise to a broad spectrum of climates and weather patterns.
 2. **Monsoons:** India experiences several monsoon seasons characterized by heavy rainfall, flash floods, and other extreme weather occurrences. The nation's river systems and basins are paramount in assessing the probability of flooding events.
 3. **Heatwaves and Droughts:** Certain sections of India are subject to intense heatwaves and droughts due to topographical variations, prevailing wind patterns, and arid climates.
 4. **Urbanization and Infrastructure:** Urbanisation coupled with insufficient infrastructure in specific regions of India intensifies the impact of natural calamities, particularly floods in highly populated urban localities.

7.10.1 Indonesia

Disaster Occurrences Spatial Visualisation

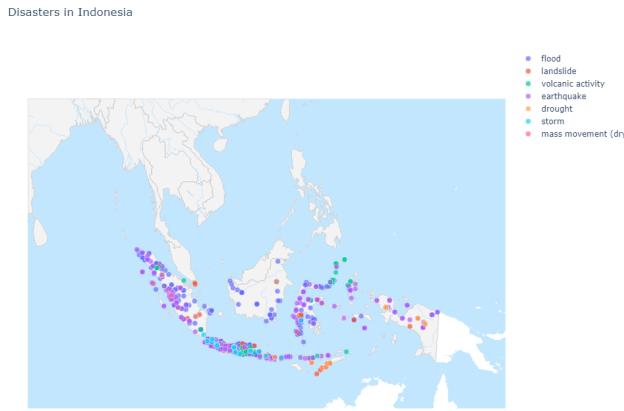


Figure 7.11: Disasters in Indonesia

Individual Types of Disasters

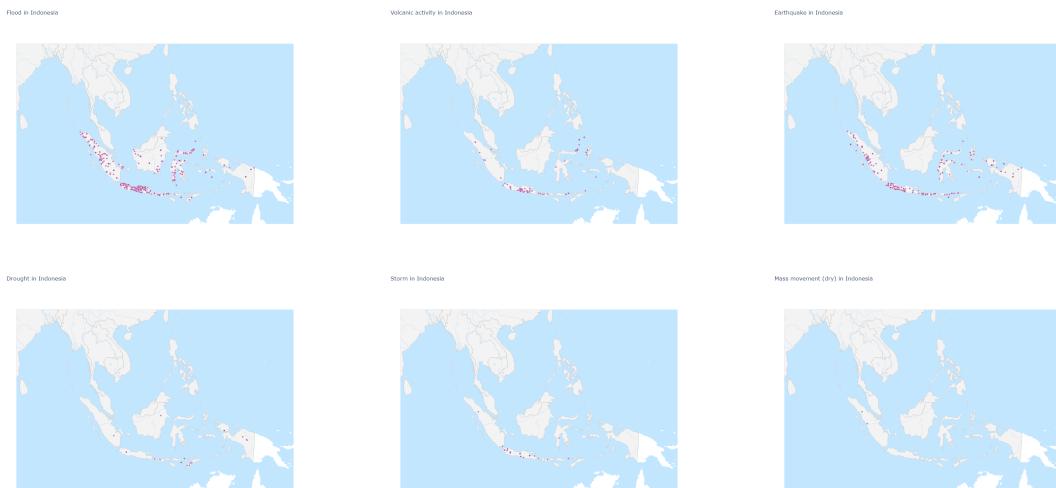


Figure 7.12: Spatial Visualisation of Different Types of Disasters in Indonesia

Temporal Evolution of Indonesia



Figure 7.13: Animated Scattergeo Plot Illustrating Disaster Locations over Time in Indonesia

7.10.2 India

Interactive Scattergeo Plot with Filtering Options for Different Disaster Types

These maps depict the geographical extent of disaster-prone regions within India. The code seamlessly categorizes data across various disaster types, affording users the option to tailor their view for each catastrophe. Through intuitive controls, users can effortlessly adjust the visual display of each disaster type on these dynamic maps.

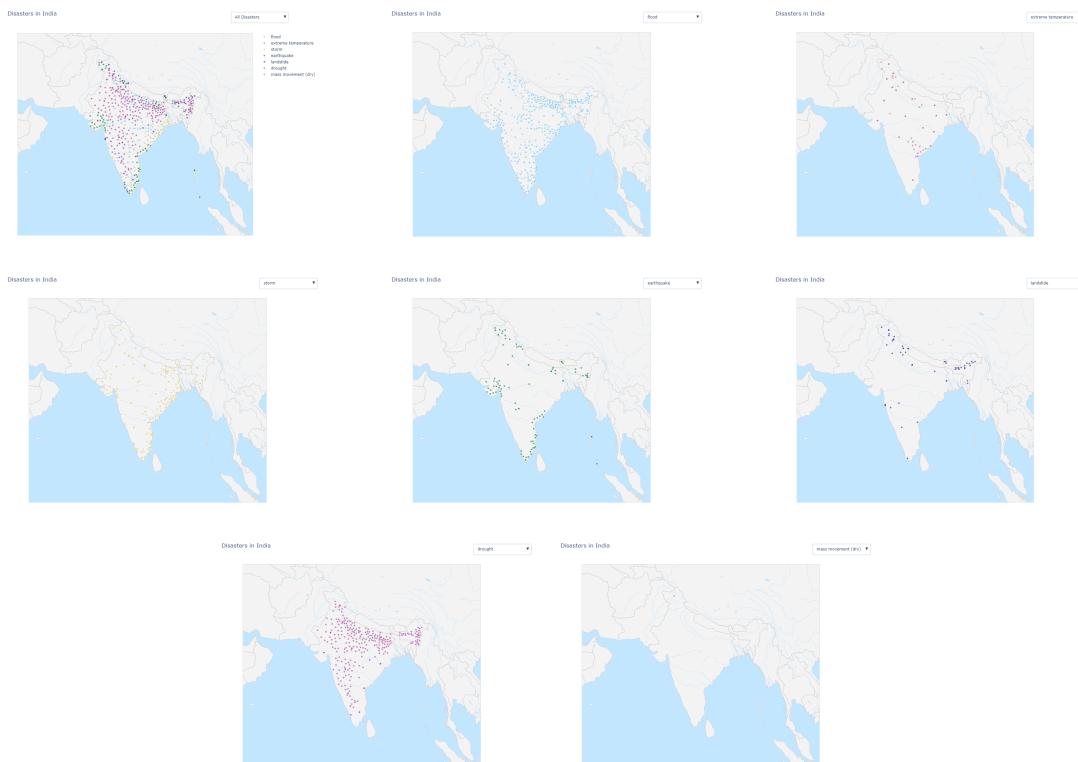


Figure 7.14: Interactive Scattergeo Plot of Spatial Distribution of Different Disasters

Temporal Evolution of India

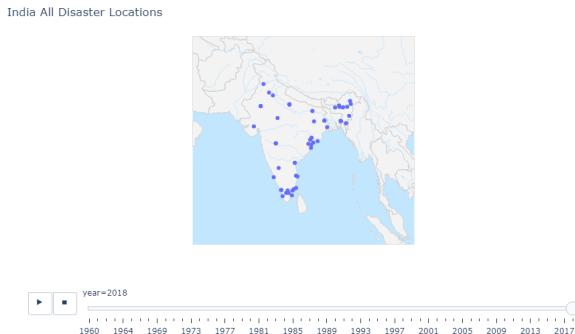


Figure 7.15: Animated Scattergeo Plot Depicting Disaster Locations in India over Time

Total Disaster Comparison between Indonesia and India over Time

These characteristics contribute to the observed gap between Indonesia and India concerning the diversity and frequency of natural disasters, explaining their divergent rankings across different disaster categories.



Figure 7.16: Total Disasters: India Vs Indonesia

7.11 Regional Analysis: Top adm1 (Region/Province) for All Disasters

The bar graph illustrates the top 10 provinces or states with the highest frequency of disasters. The province or state that exhibits the most significant frequency of disaster occurrences is situated at the topmost position of the graph, facilitating direct

comparisons with other provinces or states. This visualization presents disaster data and incorporates valuable insights relevant to catastrophe management, resource allocation, and targeted mitigation methods. This visualization tool enhances the ability to identify places impacted by disasters, hence serving as a significant tool for decision-making.

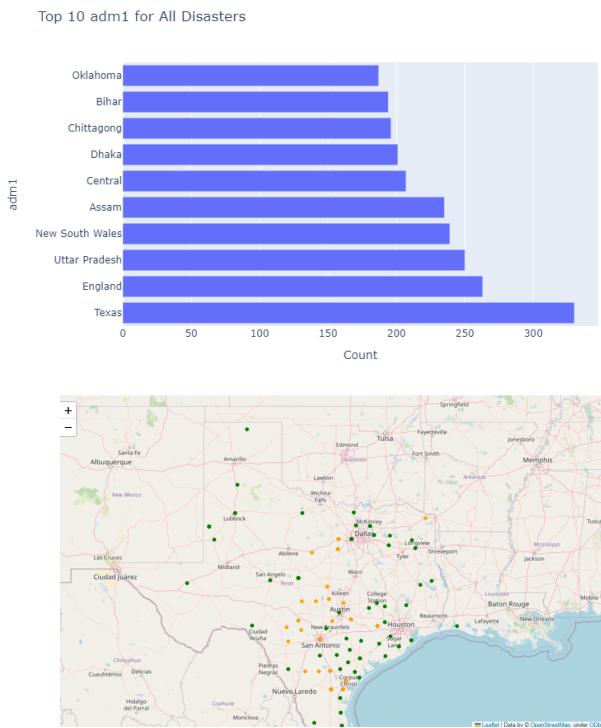


Figure 7.17: Graph of Top 10 Global Provinces and Taxes Scatterplot

7.12 City Analysis: Top 10 Locations for all Disasters

This bar graph portrays ten locations with the highest occurrence of disasters.

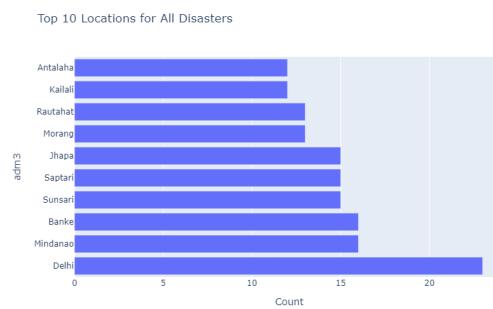


Figure 7.18: Top 10 Locations

Chapter 8

Discussion

Within a comprehensive worldwide investigation utilizing machine learning methodologies, discernible trends emerged about the distribution of disasters across several geographical locations. This undertaking, though, centers its attention on the Asian continent, with a specific concentration on South Asia, East Asia, and Southeast Asia. These regions have been identified as focal points of catastrophic events, underscoring their increased vulnerability to natural disasters. China consistently ranks among the top five countries regarding the number of disaster incidences, suggesting a significant prevalence of diverse natural catastrophes. China's geographical vastness and climatic diversity imply that the occurrences and types of natural disasters are likely to differ across different regions.

Its extensive size and diverse geography significantly influence the country's natural disaster susceptibility. To comprehensively analyze the frequency and persistence of disasters, it is imperative to consider the regional and national attributes, as the results presented. Indonesia has been identified as a region with heightened vulnerability to seismic events, mass movements of land, and volcanic phenomena, primarily originating within its territorial boundaries. The finding mentioned above can be attributed to the geographical positioning of Indonesia within the Pacific Ring of Fire, an area renowned for its heightened seismic and volcanic occurrences.

The frequency of weather disasters, including earthquakes, floods, and abnormally high temperatures, has witnessed a notable escalation in India. Most of these weather-related occurrences have occurred within the nation's territorial boundaries. Considering its expansive coastline, rugged hilly topography, and intricate geographical configuration, India's position might be assessed preliminarily. The unique geographical location of India increases its vulnerability to such catastrophic disasters.

Geographical and regional climate conditions significantly influence the global distribution of natural disasters. The research examined diverse patterns and trends that define the occurrence of several categories of natural disasters over different periods. The observed rise in the frequency of cyclones and floods suggests an elevated susceptibility and probability of these events, which may be linked to factors such as climate change. The potential ramifications of increasing temperatures and drought conditions on ecosystems, agricultural practices, water availability, and human well-being are subjects of significant apprehension.

Comprehending the dynamic patterns and fundamental factors contributing to calamities is paramount in formulating efficacious approaches to disaster management, devising strategies for risk reduction, and establishing comprehensive policy frameworks. This understanding empowers individuals involved in decision-making processes to effectively distribute resources, proactively undertake measures to prevent or reduce the impact of natural catastrophes, and devise flexible strategies to alleviate the adverse consequences of such events. This study employed machine learning techniques to visually illustrate the dispersion of impacted populations and the overall fatality count resulting from several natural calamities. This application is an excellent tool for finding geographical areas with heightened mortality rates.

A choropleth map is crucial in efficiently discerning trends and discrepancies in the spatial distribution of overall mortality. This instrument holds significant importance in the context of future disaster management, preventative analysis, and decision-making

activities. The utilization of the visualization technique effectively serves to accentuate worldwide patterns in deaths associated with natural disasters. By employing data filters, program users can enhance their understanding of the geographical extent of the crisis by highlighting various global areas. Visualization aids in the knowledge of the spatial patterns and the critical regions of global disasters, improving comprehension of their overall dynamics.

Considering Indonesia's high susceptibility to seismic events such as earthquakes, landslides, and volcanic eruptions, it might be regarded as a focal point for natural disasters. In light of this, it is feasible to adapt the current code to produce spatial visualizations for India and Indonesia. This software application allows users to efficiently identify places impacted and various disasters, enhancing their understanding of the current situation. Furthermore, using color-coded markers visually indicates and differentiates the different disaster groups, augmenting the comprehension of distribution patterns.

Spatial visualization is a valuable tool in various domains, such as disaster management, risk assessment, and strategic planning. By employing this approach, stakeholders can effectively identify locations with a high risk of disasters and strategically deploy resources, optimizing efforts in responding to such events.

The project's potential impact includes:

- 1. Disaster Preparedness and Management:** The initiative of Disaster Preparedness and Management holds promise in aiding authorities and communities to bolster their disaster preparedness efforts. It achieves this by facilitating the formulation of more efficient evacuation strategies, optimizing resource allocation, and implementing early warning systems.
- 2. Urban Planning and Infrastructure Development:** The project's significant insights have the potential to impact urban planning decisions and provide guidance for the building of infrastructure that is resilient to disasters.

3. Emergency Response and Resource Allocation: The project's projections have the potential to serve as a valuable tool for emergency responders, enabling them to strategically allocate resources and effectively coordinate relief efforts in areas that have been severely impacted.
4. **Community Awareness and Education:** Through distributing information regarding potential disasters influenced by geographical factors, this project empowers citizens to proactively adopt preventive measures and engage in efforts to strengthen community resilience.
5. **Adaptation to Climate Change:** The ability of the program to spot changes in disaster patterns aligns with the growing recognition of the impacts of climate change. This enables the process of adjusting to changing threats and the formulation of tactics that possess resilience in the face of shifting climatic conditions.
6. **Evidence-Based Policy Formulation:** The study's findings have the potential to assist policymakers in developing evidence-based approaches for disaster management, land use planning, and environmental conservation.
7. **Global Aid and Cooperation:** When a program encompasses numerous nations or areas, it possesses the capacity to facilitate international collaboration to decrease risks associated with disasters. This objective could be achieved by facilitating information and best practice sharing among diverse organizations.
8. **Research and Innovation:** The program serves as a fundamental basis for future research and innovation in catastrophe risk assessment and prediction, therefore playing a pivotal role in advancing the discipline of disaster management.

The program can favorably influence disaster resilience, public safety, and sustainable development. Through integrating geospatial data, machine learning approaches, and advanced visualization techniques, the project has successfully developed a vital tool that empowers communities, governments, and organizations to effectively address the

impacts of natural disasters and foster the development of a more resilient future.

The project's diverse impacts encompass disaster management, urban planning, emergency response, community participation, policy development, and global collaboration in tackling catastrophe-related issues.

Chapter 9

Conclusion

Considerable progress has been achieved in investigating and overseeing catastrophic risk. The present study has made notable progress in disaster prediction and visualization with geospatial attributes by employing machine learning models and advanced data visualization approaches.

The proposed project holds promise for enhancing crisis management at both national and international levels, with a particular focus on countries such as India and Indonesia. Integrating these approaches presents a potentially fruitful method for enhancing disaster response strategies, fostering resilience, and facilitating efficient global catastrophe preparedness and management.

9.1 Summary of Findings

The investigation has unveiled the profound attributes and incidents intrinsic to examining natural calamities. This study's utilization of automated learning models has showcased its efficacy in precisely forecasting natural disasters. End-users can foresee and prepare for various potential risks and weaknesses. Advanced data visualization technologies have facilitated stakeholders in making informed decisions by gaining deeper insights and optimizing resource allocation. Furthermore, the integration of dynamic variables such as swings in temperature, migrations of populations, and developments in

predictive models has enhanced the precision and reliability of disaster forecasts.

9.2 Conclusion Statement

Integrating geographic data analysis, machine learning, and data visualization has led to a significant risk assessment and management shift, enabling a proactive and informed strategy. This program has the potential to act as a catalyst for the transformation of disaster response tactics, leading to the preservation of human lives, the reduction of damages, and the facilitation of community recovery.

9.3 Recommendations

The initiative's potential impact is expected to change various industries significantly. Policymakers and disaster management authorities may efficiently address and reduce disaster risks by harnessing the potential of predictive models and visualization tools. This approach enables them to develop policies and plans grounded in empirical evidence. Urban planners and infrastructure developers will gain crucial guidance in developing resilient cities and essential infrastructure capable of withstanding catastrophic events.

Moreover, this work serves as a fundamental basis for future investigations. Based on this underlying framework, scholars can explore supplementary data sources and broaden the range of predictive models to encompass a more comprehensive array of natural and anthropogenic calamities. Partnerships among academic institutions, government agencies, and non-governmental organizations possess the capacity to cultivate all-encompassing frameworks for managing disaster risks. Consequently, these collaborations facilitate more efficient and coordinated reactions to potential threats and perils from natural calamities. .

9.4 Reflection and Self-Evaluation

As an academic scholar, this dissertation demonstrates my steadfast dedication to exploring advanced techniques in geospatial data analysis and forecasting methods for disaster prediction. My adeptness in utilizing technology to tackle pressing concerns in catastrophe risk assessment and management is demonstrated by the effective amalgamation of machine learning algorithms and data visualization approaches.

During this undertaking, I faced several constraints, including a significant focus on geographical attributes and the necessity for more investigation into regional or cultural variables that impact catastrophe planning and response. The program is noteworthy for its ability to facilitate collaboration across disciplines in geospatial data analytics, machine learning, and data visualization. This integration sets a high benchmark for interdisciplinary cooperation. This initiative offers valuable contributions to the scientific community by employing novel approaches, thereby equipping communities with tools to bolster disaster resilience and safeguard against the increasing risks of natural and anthropogenic disasters.

Bibliography

- [1] Ekeanyanwu, C. V., Obisakin, I., Aduwenye, P., & Dede-Bamfo, N. (2022, January 1). Merging GIS and Machine Learning Techniques: A Paper Review. <https://www.scirp.org/journal/paperinformation.aspx?paperid=119857>
- [2] Emergency Events Database (EM-DAT). Inventorying Hazards & Disasters Worldwide Since 1988. <https://www.emdat.be/>
- [3] Gao, M., Wang, Z., & Yang, H. (2022, July 21). Review of Urban Flood Resilience: Insights from Scientometric and Systematic Analysis. <https://www.mdpi.com/1660-4601/19/14/8837>
- [4] Geocoded Disasters (GDIS) Dataset (1960-2018). Natural Disasters. <https://sedac.ciesin.columbia.edu/data/set/pend-gdis-1960-2018>
- [5] Ghaffarian, S., & Emtehani, S. (2021, April 6). Monitoring Urban Deprived Areas with Remote Sensing and Machine Learning in Case of Disaster Recovery. <https://www.mdpi.com/2225-1154/9/4/58>
- [6] Haddad, R., et al. (2016, January 1). Landslide Cartography at the Region of Nabeul-Hammamet Based on Geographic Information System and Geomatic. <https://www.scirp.org/journal/paperinformation.aspx?paperid=71009>
- [7] Harilal, N., Singh, M., & Bhatia, U. (2021, January 1). Augmented Convolutional LSTMs for Generation of High-Resolution Climate Change Projections. <https://ieeexplore.ieee.org/document/9348885>

- [8] Mohan, C. R., & Wagle, A. A. (2018, October 16). Indonesia and India: Dealing with Disasters Together. <https://www.isas.nus.edu.sg/wp-content/uploads/2018/10/ISAS-Briefs-No.-607-Indonesia-and-India-Dealing-with-disasters-together-1.pdf>
- [9] Pu, R. (2017, January 16). A Special Issue of Geosciences: Mapping and Assessing Natural Disasters Using Geospatial Technologies. <https://www.mdpi.com/2076-3263/7/1/4>
- [10] Rezaei, R., & Ghaffarian, S. (2021, October 19). Monitoring Forest Resilience Dynamics from Very High-Resolution Satellite Images in Case of Multi-Hazard Disaster. <https://www.mdpi.com/2072-4292/13/20/4176>
- [11] Ting, C.-Y., et al. (2021, January 1). Geospatial Analytics for COVID-19 Active Case Detection. <https://www.techscience.com/cmc/v67n1/41163>
- [12] Saini, A. (2021). An Introduction to Random Forest Algorithm for beginners. <https://www.analyticsvidhya.com/blog/2021/10/an-introduction-to-random-forest-algorithm-for-beginners/>