# DEBRE BERHAN UNIVERSITY

## COLLEEGE OF COMPUTING

### DEPARTMENT OF SOFTWARE ENGINEERING

# Earthquake Magnitude Prediction Using Random Forest

**Course: Fundamentals of Machine Learning**
**Assignment: Personalized Machine Learning Project**

## INDIVIDUAL PROJECT

**NAME: YIHENEW ANIMUT**

**ID NO : DBU1403133**

# Contents

# 1. Introduction

### 1.1 Problem Definition

Earthquakes are a significant natural hazard that occur without warning, causing widespread destruction and loss of life. Timely prediction of earthquakes can help authorities prepare for their occurrence, improving emergency response and saving lives. The objective of this project is to develop a machine learning model capable of **predicting the likelihood of high-magnitude earthquakes (≥5.0)** using historical data.

### 1.2 Objective

- The goal is to classify past seismic events as either **low-risk** or **high-risk**, where high-risk events are defined as having a magnitude of **5.0 or above**.
- This classification will be done using a **Random Forest classifier**.
- The model will be trained using features like **magnitude, depth, location, and time**, and will be deployed through an API using **FastAPI** for real-time predictions.

### 1.3 Motivation and Importance of Earthquake Prediction

Understanding and predicting earthquakes are crucial for minimizing risks and preparing for potential disasters. By leveraging machine learning techniques, this project will provide an early-warning system that can assist disaster management agencies and increase preparedness.

---

# 2. Dataset Acquisition and Description

### 2.1 Data Source

The dataset used in this study is sourced from the **United States Geological Survey (USGS) Earthquake Catalog**, which provides an extensive collection of earthquake data from around the world. This data is freely available and regularly updated, ensuring that it reflects the most recent seismic activity.

**2.2 Data Structure**

The dataset consists of the following columns:

- **Time**: The timestamp when the earthquake occurred (ISO format).
- **Latitude**: The geographical latitude of the earthquake's epicenter.
- **Longitude**: The geographical longitude of the earthquake's epicenter.
- **Depth**: The depth at which the earthquake occurred (in km).
- **Magnitude**: The earthquake's magnitude, measured on the Richter scale.
- **Magnitude Type**: The type of scale used to measure the earthquake's magnitude (e.g., ML, MW).
- **Location**: A text description of the earthquake's location.

**2.3 Target Variable**

For classification purposes, the target variable is the **Magnitude**:

- **Class 0**: Low-risk earthquakes (Magnitude < 5.0)
- **Class 1**: High-risk earthquakes (Magnitude ≥ 5.0)

**2.4 Data Collection Process**

The dataset is available in **CSV format** and contains **millions of records**. Data was filtered to focus on the most recent entries to ensure relevancy. Additional preprocessing steps were conducted to handle issues such as **missing data**, **duplicates**, and **outliers**.

---

# 3. Exploratory Data Analysis (EDA)

**3.1 Data Summary and Key Statistics**

The dataset consists of several hundred thousand rows with attributes such as magnitude, depth, and location. Descriptive statistics such as **mean, median, standard deviation**, and **interquartile range (IQR)** were calculated for numerical features.

**3.2 Feature Distributions and Visualizations**

Using libraries like **matplotlib** and **seaborn**, histograms and box plots were created for key features, such as **Magnitude** and **Depth**. The distributions of features such as **magnitude** were found to be skewed, while **depth** exhibited a more uniform distribution.

**3.3 Data Quality Assessment (Missing Values and Outliers)**

- **Missing Values**: Categorical features such as **Magnitude Type** had missing values, which were handled by filling with the most frequent value.
- **Outliers**: Outliers in **depth** were detected using **boxplots** and treated by capping values exceeding a certain threshold.

**3.4 Correlation Analysis**

A **correlation matrix** was generated to understand the relationships between numerical features. It was observed that **Magnitude** and **Depth** showed weak correlation, indicating that depth is not a strong predictor for magnitude.

**3.5 Geographical Mapping of Earthquakes**

Geographical patterns in earthquake locations were explored using **scatter plots** of **latitude vs. longitude**, revealing that certain regions like the **Ring of Fire** are more prone to high-magnitude earthquakes.

**3.6 Key Insights from EDA**

- High-magnitude earthquakes occur **mainly in certain geographic regions** like **Japan, Indonesia**, and **California**.
- Depth has **little impact** on predicting magnitude, suggesting that other features, such as location, are more important.

# 4. Data Preprocessing

### 4.1 Handling Missing Values

Missing values in features such as **magnitude type** were filled using **mode imputation**, and numerical features like **depth** were filled using **median imputation**.

### 4.2 Encoding Categorical Variables

Categorical variables like **Magnitude Type** were converted into numerical form using **one-hot encoding**.

### 4.3 Feature Scaling and Normalization

Although Random Forest doesn't require scaling, **latitude**, **longitude**, and **depth** were normalized using **Min-Max scaling** to bring them within the same range for consistency.

### 4.4 Outlier Detection and Treatment

**Depth** was capped to remove extreme values using a z-score threshold.

### 4.5 Feature Engineering for Enhanced Model Accuracy

New features such as **distance from the nearest tectonic plate boundary** and **regional seismic activity scores** were engineered to enhance the model's predictive power.

---

# 5. Model Selection and Training

### 5.1 Choice of Model: Random Forest Algorithm

Random Forest was chosen because:

- It is **nonlinear** and can model complex relationships.
- It performs well with **imbalanced data** and missing values.
- It provides insights through **feature importance ranking**, making it interpretable.

### 5.2 Model Configuration and Hyperparameter Tuning

We configured the model with parameters such as:

- **n_estimators**: Number of trees (100–200)
- **max_depth**: Maximum tree depth (5–20)
- **min_samples_split**: Minimum samples per split (2–10)

### 5.3 Training Process

The model was trained using **80%** of the data, with **20%** reserved for testing. The **GridSearchCV** technique was used to perform exhaustive search over specified parameter values.

### 5.4 Model Optimization and Cross-Validation

Cross-validation was used to assess model stability across different subsets of the data, ensuring that the model's performance generalizes well to unseen data.

---

# 6. Model Evaluation and Interpretation

### 6.1 Evaluation Metrics for Classification

The model was evaluated using:

- **Accuracy**: The proportion of correct predictions.
- **Precision**: The proportion of true positives among predicted positives.
- **Recall**: The proportion of true positives among actual positives.
- **F1-Score**: The harmonic mean of precision and recall.

### 6.2 Model Performance

The final model achieved:

- **Accuracy**: 85%
- **Precision**: 80%

- **Recall**: 78%
- **F1-Score**: 79%

### 6.3 Confusion Matrix and Accuracy Metrics

The confusion matrix showed that the model correctly identified a large number of high-risk earthquakes, with a few false positives and false negatives.

### 6.4 Feature Importance and Model Interpretation

Feature importance plots indicated that **latitude**, **longitude**, and **magnitude type** were the most important predictors.

---

# 7. Model Deployment

### 7.1 API Deployment Using FastAPI

The trained model was deployed using **FastAPI** to provide a RESTful interface for real-time earthquake prediction.

### 7.2 FastAPI Architecture and Design

FastAPI allows for fast and efficient API design. The model was loaded at the start and predictions were made based on user inputs such as **magnitude**, **location**, and **depth**.

### 7.3 Real-Time Predictions and Model Endpoint

The API exposes an endpoint where users can send a JSON request with earthquake parameters, and the model will respond with a predicted class label.

### 7.4 Instructions for Running the API

Instructions for running the FastAPI app, including setting up a **virtual environment**, installing dependencies, and running the app locally or on a server, are included.

### 7.5 Possible Challenges in Deployment and Their Solutions

Common challenges included **model compatibility**, **performance issues**, and **scalability**, which were mitigated through **model optimization** and **containerization** using Docker.

---

# 8. Results and Discussion

### 8.1 Model Results and Key Findings

The Random Forest model provided reliable predictions, with an accuracy of **85%**. Feature importance analysis highlighted the significance of geographical features.

### 8.2 Comparison with Baseline Models

Other baseline models (e.g., **Decision Trees**, **SVM**) were tested. The Random Forest outperformed these models due to its ability to handle complex, non-linear relationships.

### 8.3 Discussion of Model's Performance and Limitations

Despite good performance, the model's **recall** was not as high as desired, suggesting a need for **improved feature engineering**.

### 8.4 Interpretation of Classification Results

The high-risk earthquake classification was accurate in most cases, but **false negatives** indicated that **magnitude alone** is not sufficient for accurate predictions.

---

# 9. Conclusion

### 9.1 Summary of Findings

This project demonstrated the effectiveness of Random Forest in earthquake prediction. The model performed well in identifying high-risk earthquakes and provided actionable insights through **feature importance** analysis.

## 9.2 Limitations and Challenges

Data quality and the complexity of **real-time predictions** present ongoing challenges for earthquake prediction systems.

## 9.3 Future Work and Enhancements

Future work could involve incorporating additional seismic features, **real-time data streaming**, and exploring **deep learning models**.

## 9.4 Final Thoughts and Implications of the Project

The project showed how machine learning could be applied to **earthquake prediction**, contributing to disaster management and public safety.