# Flood Probability Prediction Analysis

## 1. Introduction

This project aims to predict flood probabilities using machine learning models. The analysis involves data preprocessing, feature engineering, model training, and evaluation. The datasets used include training and testing data with various features related to flood probability. The goal is to build robust models that can accurately predict flood probabilities, which can be crucial for disaster management and preparedness.

## 2. Data Loading and Initial Inspection

### 2.1. Import Libraries

I start by importing the necessary libraries for data manipulation, visualization, and machine learning. Libraries such as pandas, numpy, seaborn, matplotlib, and machine learning libraries like xgboost, catboost, and lightgbm are essential for this analysis.

### 2.2. Load Data

The training and testing datasets are loaded from CSV files. These datasets contain various features that are potential indicators of flood probability.

### 2.3. Inspect Data

I inspect the shape of the datasets to understand the number of rows and columns. The training dataset has 1,117,957 rows and 22 columns, while the testing dataset has 745,305 rows and 21 columns. I also print the first few rows of each dataset to get a sense of the data structure and content.

### 2.4. Check for Null Values

I check for null values in both datasets to ensure data integrity. Both datasets have no null values, which is a good sign for the quality of the data.

## 3. Data Preprocessing

### 3.1. Remove Outliers

Outliers can significantly affect the performance of machine learning models. I define a function to remove outliers using the Interquartile Range (IQR) method. This method identifies and removes data points that fall below the lower bound or above the upper bound, which are calculated based on the first and third quartiles. After removing outliers, the shape of the cleaned training dataset is reduced to 849,100 rows and 20 columns.

### 3.2. Feature Engineering

Feature engineering involves creating new features from the existing ones to improve the performance of machine learning models. I define a function to calculate statistical measures such

as mean, standard deviation, maximum, minimum, range, variance, skewness, kurtosis, and median absolute deviation for each row. These new features can capture additional patterns and relationships in the data.

### 3.3. Data Scaling

Data scaling is essential for machine learning models that are sensitive to the scale of the input features. I use the StandardScaler from sklearn to standardize the features by removing the mean and scaling to unit variance. This ensures that all features contribute equally to the distance calculations in the models.

### 4. Model Training and Evaluation

### 4.1. CatBoost Regressor

CatBoost is a gradient boosting library that is efficient and provides state-of-the-art results. I split the data into training and testing sets, train the CatBoost model, and evaluate its performance using the R-squared metric. The R-squared value for the CatBoost model is approximately 0.9267, indicating a high level of accuracy in predicting flood probabilities.

### 4.2. LightGBM Regressor

LightGBM is another gradient boosting framework that is designed to be distributed and efficient. I train the LightGBM model and evaluate its performance using the R-squared metric. The R-squared value for the LightGBM model is also high, indicating its effectiveness in predicting flood probabilities.

### 4.3. XGBoost Regressor

XGBoost is a popular gradient boosting library that is known for its performance and speed. I train the XGBoost model and evaluate its performance using the R-squared metric. The R-squared value for the XGBoost model is approximately 0.9266, which is comparable to the other models.

### 5. Visualization

### 5.1. Residual Plots

Residual plots are used to assess the performance of the models. I create residual plots for the CatBoost model to visualize the differences between the predicted and actual values. The residual plot for CatBoost shows that the residuals are randomly distributed around zero, indicating a good fit. I also create an actual vs. predicted plot to visualize the relationship between the actual and predicted values.

### 5.2. Feature Importance

Feature importance plots help identify the most influential features in the model. I create a feature importance plot for the CatBoost model to visualize the importance of each feature. This plot helps in understanding which features contribute the most to the prediction of flood probabilities.

### 5.3. Q-Q Plots

Q-Q plots are used to assess the normality of the residuals. I create Q-Q plots for the LightGBM and CatBoost models to visualize the distribution of the residuals. The Q-Q plots help in understanding whether the residuals follow a normal distribution, which is an important assumption for many statistical tests.

**6. Predicted Values File**

I create a DataFrame containing the predicted flood probabilities for the testing dataset. The predictions are a lighted average of the predictions from the LightGBM, CatBoost, and XGBoost models. I save the predicted values to a CSV file for further analysis. I also create a histogram to visualize the distribution of the predicted flood probabilities.

**7. Conclusion**

The analysis involved data preprocessing, feature engineering, model training, and evaluation. The CatBoost, LightGBM, and XGBoost models Ire trained and evaluated, with all models achieving high R-squared values. The residual plots, feature importance plots, and Q-Q plots provided insights into the model performance and the distribution of the predictions. The predicted flood probabilities Ire saved to a CSV file for further analysis. This project demonstrates the effectiveness of machine learning models in predicting flood probabilities, which can be crucial for disaster management and preparedness.

Correlation Matrix of Train Data

Boxplot after data cleaned


Distribution of Numeric Columns

| id | MonsoonIntensity | TopographyDrainage | RiverManagement | Deforestation | Urbanization |
|---|---|---|---|---|---|
| Train skewness:0.00 | Train skewness:0.44 | Train skewness:0.46 | Train skewness:0.43 | Train skewness:0.43 | Train skewness:0.44 |

| ClimateChange | DamsQuality | Siltation | AgriculturalPractices | Encroachments | IneffectiveDisasterPreparedness |
|---|---|---|---|---|---|
| Train skewness:0.43 | Train skewness:0.44 | Train skewness:0.45 | Train skewness:0.42 | Train skewness:0.46 | Train skewness:0.44 |

| DrainageSystems | CoastalVulnerability | Landslides | Watersheds | DeterioratingInfrastructure | PopulationScore |
|---|---|---|---|---|---|
| Train skewness:0.44 | Train skewness:0.44 | Train skewness:0.43 | Train skewness:0.45 | Train skewness:0.44 | Train skewness:0.45 |

| WetlandLoss | InadequatePlanning | PoliticalFactors | FloodProbability |
|---|---|---|---|
| Train skewness:0.44 | Train skewness:0.46 | Train skewness:0.44 | Train skewness:0.05 |

id MonsoonIntensity TopographyDrainage RiverManagement Deforestation

Urbanization ClimateChange DamsQuality Siltation AgriculturalPractices

Encroachments IneffectiveDisasterPreparedness DrainageSystems CoastalVulnerability Landslides

Watersheds DeterioratingInfrastructure PopulationScore WetlandLoss InadequatePlanning

PoliticalFactors FloodProbability

## Residual Plot (XGBoost)
Predicted values
Residuals

## Actual vs Predicted Plot(XGBoost)
Predicted Values
Actual Values

## Feature Importance(XGBoost)
harmonic_mean 23312.0
geometric_mean 22110.0
coeff_variation 19765.0
id 17447.0
MonsoonIntensity 10384.0
TopographyDrainage 9310.0
RiverManagement 8676.0
Deforestation 8221.0
Urbanization 8224.0
DamsQuality 7763.0
Siltation 7760.0
AgriculturalPractices 7190.0
IneffectiveDisasterPreparedness 7621.0
Encroachments 7253.0
PoliticalFactors 7433.0
CoastalVulnerability 4321.0
ClimateChange 4324.0
DrainageSystems 1342.0
Landslides 1392.0
InadequatePlanning 1215.0
WetlandLoss 1213.0
DeterioratingInfrastructure 7027.0
PopulationScore 6940.0
Watersheds 6894.0
min_features 4693.0
skewness_features 1628.0
kurtosis_features 1036.0
mean_features 388.0
std_features 1171.0
range_features 411.0
median_absolute_deviation 183.0
variance_features 182.0
max_features 176.0
mean_absolute_deviation 157.0
range_abs_diff
F score

## Residual Distribution (XGBoost)
Frequency
Residuals

Residual Plot (CatBoost)

Actual vs. Predicted Plot (CatBoost)

Feature Importance (CatBoost)

Residual Distribution (CatBoost)

Residual Plot (LGBMRegressor)

Actual vs. Predicted Plot (LGBMRegressor)

Feature Importance (LGBMRegressor)

Residual Distribution (LGBMRegressor)

Distribution of Flood Probability Predictions: LightGBM vs. XGBoost

Q-Q Plot for LightGBM Predictions

R-squared: 0.9268

Q-Q Plot for CatBoost Predictions

R-squared: 0.9267

Flood Probability Distribution