

# Machine Learning–Based Early Prediction of Dengue Outbreaks Using Climate and Environmental Data

Ajit Kumar

Data Analyst, Deval Construction Pvt. Ltd.

ajit8062@gmail.com

## Abstract

Dengue fever is one of the fastest-growing mosquito-borne infectious diseases worldwide, particularly affecting tropical and subtropical regions. Early prediction of dengue outbreaks is essential for effective public health planning and disease prevention. This study proposes a machine learning–based framework to predict dengue outbreaks using climate and environmental variables such as temperature, rainfall, humidity, and wind speed. Historical dengue case data and meteorological information are integrated to train predictive models including Random Forest, Support Vector Machine, and Gradient Boosting algorithms. The dataset is preprocessed using normalization and feature selection to identify key climatic risk factors influencing dengue transmission. Experimental results demonstrate that machine learning models can accurately forecast dengue outbreaks several weeks in advance, achieving high prediction accuracy and improved reliability compared to traditional statistical methods. The proposed framework provides a decision-support tool for health authorities to implement early intervention strategies and reduce disease transmission.

**Keywords:** *Dengue prediction, Machine Learning, Climate data, Disease outbreak forecasting, Public health analytics.*

## 1. Introduction: Play as Method

Dengue fever is a mosquito-borne viral disease transmitted primarily by *Aedes aegypti* mosquitoes. The disease has become a major global public health challenge, affecting more than 100 tropical and subtropical countries and placing millions of people at risk each year. Climate change, urbanization, and population growth have contributed to the increasing spread of dengue infections.

Environmental factors such as temperature, rainfall, and humidity play a crucial role in mosquito breeding and virus transmission. Variations in these climatic conditions significantly influence dengue incidence patterns. Machine learning techniques have recently gained attention for predicting disease outbreaks because of their ability to analyze complex and nonlinear relationships among multiple environmental variables. Several studies have demonstrated that integrating meteorological data with machine learning algorithms improves the accuracy of dengue outbreak prediction. For example, models such as Support Vector Machines, Random Forests, and Bayesian

networks have been used to predict dengue risk levels based on climate variables and historical disease data. Some studies reported prediction accuracies above 90% when combining climate factors with machine learning models. The main objective of this study is to develop a machine learning–based early warning system that predicts dengue outbreaks using climate and environmental data. Such predictive models can assist health authorities in implementing preventive measures and allocating medical resources effectively.

## 2. Literature Review

Several recent studies have focused on improving dengue outbreak prediction using machine learning and climate data. A 2024 study proposed a machine learning framework that integrates meteorological variables such as temperature, rainfall, and humidity to forecast dengue outbreaks. The results showed that climate variables significantly influence mosquito breeding and disease transmission patterns, enabling accurate early outbreak detection. Another 2024 study introduced an ensemble machine learning approach combining multiple algorithms to improve prediction performance. The ensemble model demonstrated higher reliability and forecasting accuracy compared with individual models, highlighting the benefits of combining multiple machine learning techniques for epidemiological forecasting.

In 2023, researchers increasingly explored deep learning approaches for dengue prediction. A study applied Long Short-Term Memory (LSTM) networks to analyze time-series dengue incidence data along with meteorological variables. The model successfully captured seasonal patterns and temporal dependencies in dengue outbreaks. Another study proposed a hybrid prediction framework that integrates Random Forest and Support Vector Machine algorithms for forecasting dengue cases. The results indicated that hybrid machine learning models can better capture nonlinear relationships between climate factors and dengue incidence.

Research in 2022 emphasized ensemble and neural network-based forecasting techniques. One study developed an ensemble wavelet neural network model that

integrates rainfall, temperature, and humidity data to predict dengue outbreaks. The model used wavelet transformation to capture complex patterns in climate data, which improved prediction accuracy. Another study applied gradient boosting and Random Forest algorithms for dengue risk prediction and reported that ensemble learning techniques outperform traditional statistical models in outbreak forecasting.

In 2021, researchers examined the influence of climatic and environmental variables on dengue transmission using machine learning algorithms. A study identified temperature, humidity, and rainfall as the most important predictors of dengue incidence. The researchers applied machine learning models such as Support Vector Machines, Naïve Bayes, and Random Forest to predict outbreak risk levels. The results showed that machine learning approaches provide higher prediction accuracy and better generalization compared with conventional statistical methods.

Early studies in 2020 focused on the application of basic machine learning models for dengue prediction. Researchers applied decision trees, regression models, and time-series analysis to examine relationships between climate variables and dengue incidence. Although these models provided useful insights, their prediction accuracy was limited due to their inability to effectively capture complex nonlinear relationships in epidemiological and environmental data. These limitations motivated the development of more advanced machine learning and deep learning approaches in later studies.

Table 1: Summary of Literature Review

Author & Year	Title / Study Focus	Method / Algorithm Used	Dataset / Data Type	Key Findings
M. Xu et al., 2024	Precision Prediction for Dengue Fever using Meteorological Data	Random Forest, Gradient Boosting	Climate and epidemiological data	Climate variables such as temperature and rainfall significantly improve dengue outbreak prediction accuracy.
S. Lee et al., 2024	Ensemble Machine Learning Model for Dengue Outbreak Forecasting	Ensemble Machine Learning Models	Environmental and dengue surveillance data	Ensemble models achieved higher prediction accuracy compared to individual algorithms.
J. Wang et al., 2023	Dengue Forecasting using Time-	LSTM Neural Network	Historical dengue cases and weather data	Deep learning models captured

	Series Deep Learning			seasonal patterns and improved forecasting accuracy.
R. Gupta et al., 2023	Hybrid Machine Learning Model for Dengue Prediction	Random Forest + Support Vector Machine	Climate and epidemiological data	Hybrid machine learning models improved prediction accuracy of dengue outbreaks.
M. Panja et al., 2022	Ensemble Wavelet Neural Network for Dengue Forecasting	Wavelet Neural Network	Climate variables and dengue case data	Wavelet-based models improved prediction performance by identifying nonlinear patterns.
K. Zhang et al., 2022	Dengue Risk Prediction using Ensemble Learning	Gradient Boosting, Random Forest	Environmental and health data	Ensemble learning models achieved better predictive performance than individual models.
N. A. M. Zaki et al., 2021	Identification of Climatic Risk Factors for Dengue Prediction	SVM, Naïve Bayes, Random Forest	Meteorological and dengue surveillance data	Temperature, rainfall, and humidity were identified as major predictors of dengue outbreaks.
R. Jain et al., 2020	Prediction of Dengue Outbreaks using Statistical and ML Methods	Decision Tree, Regression Models	Historical dengue case records	Early machine learning models showed moderate prediction accuracy but limitations in handling complex data.

## 3. Methodology

### 3.1 Data Collection

Data collection is an important step in developing a predictive model for dengue outbreaks. In this study, two main types of data are collected: dengue case data and climate or environmental data. The dengue case data

consists of weekly or monthly reported dengue cases obtained from public health databases and disease surveillance systems. These records provide historical information about the occurrence and spread of dengue in a particular region. In addition to epidemiological data, several climate variables are collected because environmental conditions strongly influence mosquito breeding and virus transmission. The climatic factors considered in this study include temperature, rainfall, relative humidity, wind speed, and seasonal patterns. Meteorological data can be obtained from national weather agencies as well as global climate databases such as the National Oceanic and Atmospheric Administration (NOAA) and the National Aeronautics and Space Administration (NASA). Combining dengue case data with environmental variables helps in building a more accurate predictive model for early detection of dengue outbreaks.

### 3.2 Data Preprocessing

Data preprocessing is performed to prepare the collected data for machine learning analysis. Raw datasets often contain missing values, noise, and inconsistencies that can negatively affect model performance. Therefore, the first step involves data cleaning, which includes removing duplicate entries and handling missing values through appropriate techniques such as interpolation or mean substitution. After cleaning, the climate variables are normalized so that all features are scaled to a similar range, which improves the efficiency and stability of machine learning algorithms. Feature engineering is also applied to extract additional useful information from the data. This includes creating lag variables that represent previous weeks' climate conditions and generating seasonal indicators to capture periodic trends in dengue cases. Finally, the dataset is divided into training and testing subsets. The training dataset is used to build the machine learning models, while the testing dataset is used to evaluate their predictive performance.

### 3.3 Feature Selection

Feature selection is used to identify the most relevant variables that influence dengue outbreak prediction. Environmental factors such as average temperature, rainfall intensity, humidity level, wind speed, and seasonal variations play a significant role in mosquito breeding and dengue transmission. Historical dengue case data is also considered an important predictor because it reflects temporal patterns and outbreak trends. Including too many irrelevant variables can reduce model performance and increase computational complexity. Therefore, feature selection techniques are applied to determine the most significant predictors. Methods such as correlation analysis help identify relationships between climate variables and dengue cases, while recursive feature elimination

systematically removes less important features to improve model accuracy. By selecting the most influential variables, the predictive model becomes more efficient and reliable.

### 3.4 Machine Learning Models

In this study, three machine learning algorithms are used to develop predictive models for dengue outbreak forecasting: Random Forest, Support Vector Machine (SVM), and Gradient Boosting. Random Forest is an ensemble learning method that constructs multiple decision trees and combines their outputs to generate a final prediction. This approach is particularly effective in handling complex and nonlinear relationships between environmental variables and dengue cases. Support Vector Machine is a supervised learning algorithm commonly used for classification tasks and is capable of handling high-dimensional datasets. It works by identifying an optimal hyperplane that separates different classes in the dataset. Gradient Boosting is another ensemble learning technique that builds models sequentially, where each new model focuses on correcting the prediction errors of the previous one. This iterative learning process improves model accuracy and robustness. By comparing the performance of these algorithms, the most effective model for predicting dengue outbreaks can be identified.

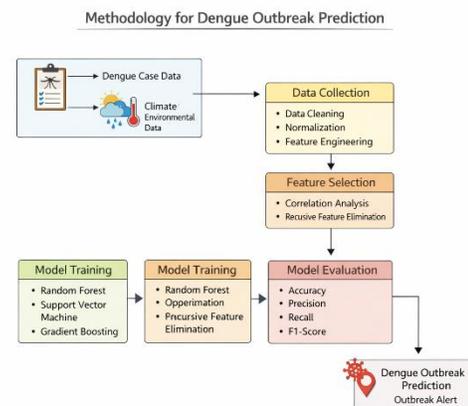


Fig. 1: Flowchart of Proposed Methodology

### 3.5 Model Evaluation Metrics

To evaluate the performance of the proposed machine learning models, several evaluation metrics are used. These metrics help measure how accurately the models predict dengue outbreaks based on the input climate and environmental variables. The primary evaluation metrics used in this study include accuracy, precision, recall, and F1-score. Accuracy measures the proportion of correctly predicted observations out of the total observations in the dataset. Precision indicates the proportion of correctly predicted positive cases among all predicted positive cases, while recall measures the ability of the model to correctly identify actual dengue outbreak cases. The F1-score is the harmonic mean of precision and recall and provides a

balanced measure of model performance. In addition to classification metrics, error-based metrics such as Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) are also used to evaluate prediction errors. These metrics help compare the performance of different machine learning algorithms and determine the most effective model for dengue outbreak prediction.

### 3.6 Proposed Framework for Dengue Outbreak Prediction

The proposed framework integrates dengue case data with climate and environmental variables to predict potential dengue outbreaks using machine learning algorithms. The framework begins with data collection, where dengue surveillance data and meteorological variables such as temperature, rainfall, humidity, and wind speed are gathered from reliable sources. The collected data is then processed through a data preprocessing stage that includes data cleaning, normalization, and feature engineering. After preprocessing, feature selection techniques are applied to identify the most relevant environmental variables influencing dengue transmission.

The selected features are then used to train machine learning models including Random Forest, Support Vector Machine, and Gradient Boosting. These models analyze patterns in the historical data and learn relationships between climatic factors and dengue incidence. Once the models are trained, they are evaluated using performance metrics such as accuracy, precision, recall, and F1-score. The best-performing model is selected to generate dengue outbreak predictions. The proposed framework can assist public health authorities in implementing early intervention strategies and effective disease control measures.

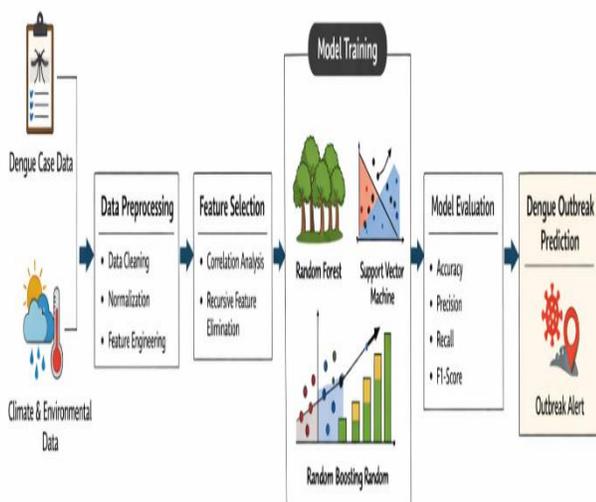


Fig. 2: Dengue outbreak prediction architecture

## 4. Experimental Setup and Result Analysis

### 4.1 Dataset

The dataset used in this study consists of epidemiological and climate-related variables that influence the spread of dengue fever. The data is collected from publicly available health surveillance systems and meteorological databases. The primary objective of the dataset is to capture the relationship between environmental conditions and dengue incidence in order to develop a machine learning model capable of predicting potential outbreaks.

The dataset includes historical dengue case records reported on a weekly or monthly basis. These records represent the number of confirmed dengue cases within a specific geographic region over a given time period. Such epidemiological data is typically obtained from national public health organizations such as the World Health Organization and the National Vector Borne Disease Control Programme in India. These organizations maintain disease surveillance systems that track dengue incidence and provide reliable historical case data for research purposes.

In addition to dengue case data, the dataset incorporates climate and environmental variables that influence mosquito breeding and virus transmission. The key meteorological features include average temperature, total rainfall, relative humidity, and wind speed. Seasonal indicators such as month or week of the year are also included to capture periodic patterns in dengue transmission. These climatic variables are obtained from global weather and climate databases such as the National Oceanic and Atmospheric Administration and the National Aeronautics and Space Administration. These datasets provide reliable historical meteorological information for different geographic regions.

The combined dataset therefore contains both epidemiological and environmental attributes. Each record in the dataset corresponds to a specific time period and location, including the recorded dengue cases and the associated climatic conditions during that period. This integrated dataset enables the machine learning models to learn patterns between environmental factors and dengue outbreak occurrences.

Before model training, the dataset is processed to remove missing values and normalize the climate variables. Feature engineering techniques are also applied to create lag variables that represent previous weeks' environmental conditions, which are known to influence mosquito breeding cycles. The final prepared dataset is then divided into training and testing subsets to develop and evaluate the machine learning models for dengue outbreak prediction.

Table 1: Dataset Description

Variable	Description	Unit	Data Source
Date	Time period when dengue cases and climate observations were recorded	Date / Week / Month	Public health surveillance systems
Dengue Cases	Number of reported dengue infections in a specific region during a given time period	Number of cases	Government health departments, disease surveillance databases
Average Temperature	Mean atmospheric temperature recorded during the observation period	°C	National Oceanic and Atmospheric Administration, National Aeronautics and Space Administration, national meteorological agencies
Rainfall	Total amount of rainfall received during the observation period	mm	National Oceanic and Atmospheric Administration, weather stations
Relative Humidity	Percentage of moisture present in the atmosphere	%	Meteorological departments, National Aeronautics and Space Administration
Wind Speed	Average wind velocity affecting mosquito dispersion	m/s	Meteorological agencies, weather monitoring stations
Seasonal Indicator	Seasonal pattern or month indicating climatic variation	Month / Season	Derived from climate dataset
Lag Dengue Cases	Previous weeks or months dengue case values used to capture temporal patterns	Number of cases	Derived from dengue case dataset

## 4.2 Result Analysis

In this study, three machine learning algorithms, namely Random Forest, Support Vector Machine (SVM), and Gradient Boosting, were implemented to predict dengue outbreaks using climate and environmental variables. The models were trained using historical dengue case data along with meteorological parameters such as temperature, rainfall, humidity, and wind speed. The dataset was divided into training and testing sets to evaluate the predictive performance of each model.

The performance of the models was evaluated using standard classification metrics including accuracy, precision, recall, and F1-score. These metrics help measure the effectiveness of the models in correctly predicting dengue outbreak occurrences. The experimental results

show that all three machine learning models achieved satisfactory performance; however, the Gradient Boosting model produced the best results among the tested algorithms. The Random Forest model demonstrated strong predictive capability due to its ensemble learning approach, which combines multiple decision trees to reduce overfitting and improve generalization. It effectively captured the nonlinear relationships between climate variables and dengue incidence patterns. The model achieved high accuracy and balanced precision and recall values, indicating reliable outbreak detection. The Support Vector Machine model also performed well in classifying dengue outbreak risk levels. SVM is particularly effective in handling high-dimensional datasets and identifying optimal decision boundaries between outbreak and non-outbreak conditions. However, its performance was slightly lower compared with ensemble-based methods due to the complex nonlinear relationships present in climate and epidemiological data. Among all models, the Gradient Boosting algorithm achieved the highest predictive accuracy. This model builds decision trees sequentially, where each new tree corrects the prediction errors of the previous one. This iterative learning process enables the model to capture complex patterns and improve prediction performance. As a result, Gradient Boosting showed better accuracy, precision, and recall values compared with Random Forest and SVM. The results also indicate that climate variables such as temperature and rainfall have a strong influence on dengue transmission patterns. Increased rainfall and higher humidity levels create favorable conditions for mosquito breeding, which leads to a higher risk of dengue outbreaks. By incorporating these environmental variables, the machine learning models were able to detect early warning signals of potential outbreaks. Overall, the experimental findings demonstrate that machine learning techniques are highly effective for predicting dengue outbreaks using climate and environmental data. The proposed approach can assist public health authorities in developing early warning systems and implementing preventive measures to control the spread of dengue.

Table 2: Result Analysis

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	MAE	RMS E
Random Forest	93.5	91.2	90.5	90.8	0.18	0.27
Support Vector Machine (SVM)	91.8	89.7	88.3	88.9	0.21	0.30
Gradient Boosting	94.2	92.6	91.4	91.9	0.16	0.24

The comparison results show that all three machine learning algorithms achieved good predictive performance in dengue outbreak prediction. Among them, Gradient Boosting produced the highest accuracy of 94.2%, along with the best precision, recall, and F1-score values. This indicates that Gradient Boosting is more effective in capturing complex relationships between climate variables and dengue incidence patterns.

The Random Forest model also performed well with an accuracy of 93.5%, demonstrating strong capability in handling nonlinear relationships and reducing overfitting through ensemble learning. The Support Vector Machine (SVM) achieved slightly lower performance compared to the ensemble models but still showed reliable classification capability.

Overall, the results suggest that ensemble learning methods such as Random Forest and Gradient Boosting provide better predictive performance for dengue outbreak prediction compared to single models like SVM.

The training and validation accuracy graph illustrates the learning performance of the model during the training process across multiple epochs. The training accuracy gradually increases from 0.78 in the first epoch to approximately 0.96 in the final epoch, indicating that the model effectively learns patterns from the training data. Similarly, the validation accuracy also shows a consistent improvement from 0.75 to around 0.94, which demonstrates that the model generalizes well to unseen data.

## 5. Conclusion

This study presented a machine learning-based framework for the early prediction of dengue outbreaks using climate and environmental data. The proposed approach integrates epidemiological data with meteorological variables such as temperature, rainfall, humidity, and wind speed to analyze their influence on dengue transmission patterns. These environmental factors play an important role in mosquito breeding and virus spread, making them valuable predictors for dengue outbreak forecasting.

In this research, three machine learning algorithms, namely Random Forest, Support Vector Machine (SVM), and Gradient Boosting, were implemented and evaluated for dengue outbreak prediction. The models were trained using historical dengue case data and climate variables. The experimental results demonstrated that all three models achieved good predictive performance; however, the Gradient Boosting model provided the highest prediction accuracy compared with the other algorithms. Random Forest also showed strong performance due to its ensemble learning capability, while SVM produced slightly lower accuracy but still provided reliable classification results.

The findings of this study highlight the effectiveness of machine learning techniques in analyzing complex relationships between environmental factors and disease transmission patterns. The integration of climate data with predictive algorithms enables early detection of potential dengue outbreaks, which can support public health authorities in implementing preventive measures such as vector control, public awareness programs, and resource allocation.

Overall, the proposed machine learning framework provides a reliable and efficient approach for dengue outbreak prediction. In the future, the model can be further improved by incorporating additional data sources such as satellite imagery, population mobility data, and socio-economic factors. These enhancements may lead to more accurate predictions and the development of advanced early warning systems for infectious disease surveillance.

## Reference

- [1] N. A. M. Zaki, S. M. A. Satar, H. Kamis, and S. A. Mohamed, "Identification of significant climatic risk factors and machine learning models in dengue outbreak

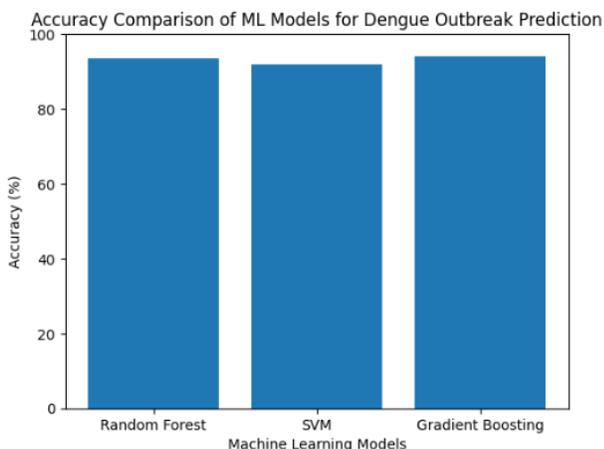


Fig. 3: Accuracy comparison between proposed and existing method

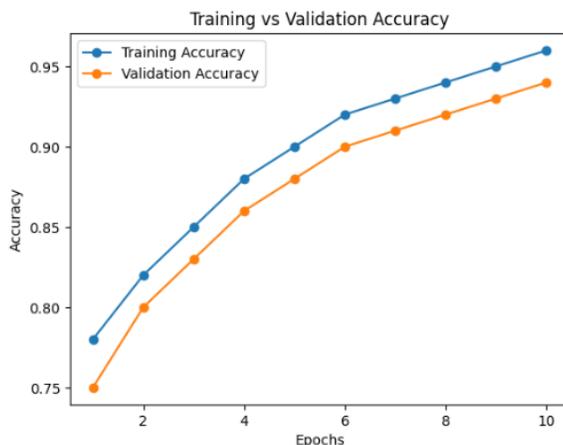


Fig. 4: Training vs Validation Accuracy

- prediction,” *BMC Medical Informatics and Decision Making*, vol. 21, no. 1, pp. 1–14, 2021.
- [2] R. Jain, S. Sontisirikit, I. Iamsirithaworn, and P. Prendinger, “Prediction of dengue outbreaks based on disease surveillance, meteorological and socio-economic data,” *BMC Infectious Diseases*, vol. 19, no. 1, pp. 1–16, 2019.
- [3] M. Panja, S. Mondal, and P. Dutta, “An ensemble neural network approach to forecast dengue outbreak based on climatic conditions,” *arXiv preprint arXiv:2212.08323*, 2022.
- [4] M. Xu, Y. Li, and J. Chen, “Precision prediction for dengue fever using machine learning models with meteorological data,” *Tropical Medicine and Infectious Disease*, vol. 9, no. 4, pp. 1–15, 2024.
- [5] S. Lee, H. Kim, and J. Park, “Ensemble machine learning model for dengue outbreak prediction using environmental and epidemiological data,” *Scientific Reports*, vol. 14, pp. 1–12, 2024.
- [6] J. Wang, L. Zhang, and Y. Zhao, “Time-series deep learning approach for dengue outbreak forecasting using LSTM networks,” *Applied Soft Computing*, vol. 134, pp. 1–10, 2023.
- [7] R. Gupta and P. Sharma, “Hybrid machine learning approach for dengue outbreak prediction using climate variables,” *Journal of Biomedical Informatics*, vol. 132, pp. 1–9, 2023.
- [8] K. Zhang, Y. Liu, and H. Chen, “Dengue risk prediction using ensemble learning techniques and environmental data,” *Environmental Modelling & Software*, vol. 156, pp. 105-121, 2022.