# Time Series Analysis-based Prediction of Dengue Spread using Climate Data

**P. Devendar Babu[1], T. Samyuktha[2], U. Krishnaveni[2], T. Kavya[2], Madhu Priya[2]**

[1,2]Department of Information Technology

[1,2]Malla Reddy Engineering College for Women (A), Maisammaguda, Medchal, Telangana.

## Abstract

Dengue is a human arbovirus disease transmitted by the female mosquito of the genus Aedes, mainly Aedes aegypti and Ae. albopictus. Dengue, the most frequent arthropod-borne viral disease, is prevalent in tropical and subtropical regions. Two major clinical forms of dengue illness involve the mild form of dengue fever and severe form mostly characterized by plasma leakage with or without haemorrhage. Two-fifths of the world population (about 2.5 billion people) is at risk of dengue infection. The prevalence of this disease has grown dramatically in the recent decades. Between 50 and 100 million people are infected each year worldwide and more than 500,000 are hospitalized. The average annual incidence was multiplied by thirty in the last fifty years. Incidence of dengue haemorrhagic fever (DHF) is increasing in many tropical regions inducing 20,000 deaths per year, mostly among children under 15 years. Dengue is endemic in all surrounding countries with the four serotypes circulating in the region within a period of ten years. Countries or territories with the highest number of reported dengue cases were Puerto Rico, the Dominican Republic, Martinique, Trinidad and Tobago and French Guiana. Population movement is an important factor in the virus dissemination. It contributes to carry new virus strains, but it also participates to introduce nonimmune subjects in an endemic area.

This proposed system is built to predict the spread of dengue fever with climate data using the concept of time series analysis. In addition, this project also performs the exploratory data analytics on the dengue dataset over a period of time. Finally, prediction analysis also performed with the usage of advancement rendered by machine learning algorithms.

**Keywords:** Time series analysis, Dengue, Machine learning.

## 1. INTRODUCTION

Dengue is a potentially life-threatening arboviral disease transmitted by female Aedes mosquitoes, especially A. aegypti, A. albopictus, and A. vitattus. These vectors are common tropical hematophagous ectoparasites. This zoonotic disease spread from African or Asian non-human primates 500 to 1000 years ago, but within the last 60 years it has spread from just 9 countries experiencing severe epidemics to become endemic in over 100 countries worldwide, even affecting non-tropical or subtropical areas. Moreover, approximately one hundred million people yearly suffer from the symptomatic disease caused by its four serotypes. Given the significant impact of environmental changes on disease transmission, the One Health approach is urgently needed to implement the integration between human, animal, and ecological health.

The objective of this paper is to provide an insight into techniques that can be used for future predictive models based on the One Health perspective, particularly in respect to Latin America but also elsewhere.

One Health is a multidisciplinary approach that acknowledges the synergy between human and animal health and their shared environment. This idea is not new; the noted nineteenth-century pathologist (and originator of the term zoonosis) Rudolph Virchow famously asserted in 1858 that "between animal and human medicine, there are no dividing lines—nor should there be".

This approach has become increasingly important in the 21st Century with the convergence of the pressures of changing climate, migration of human and animal populations, and the growing human population that increases the proximity between wildlife and humans. Indeed, the term One Health was only coined in the early 2000s with the appearance of the zoonotic SARS and H5N1 influenza diseases.

Whilst the One Health perspective is widely seen as necessary and increasingly used for better disease control, epidemiological approaches have not kept up with this change. Conventional epidemiological perspectives tend to view disease broadly from a human-only perspective, focusing on human demographic conditions with often only climatic/environmental factors accommodating the disease vector health. In contrast, One Health requires the health and lifecycle of the zoonotic disease vectors to be explicitly considered alongside the human environment, demographics, and interaction with the zoonotic host vectors.

For example, whilst environmental and sociological considerations often take a back seat in One Health, they frequently occupy the centre stage in epidemiology. Factors such as mean temperatures and rainfall used in predicting dengue, with a very vague consideration of how they affect the mosquito vectors, are an emergent challenge to be considered. High rainfall, for instance, is beneficial to mosquitoes because it provides water-filled locations for eggs and larvae, whilst the mosquitoes are primarily impervious to strikes by raindrops that might otherwise kill them. In addition, temperature and rain generally affect many other infectious and tropical diseases.

This article focusses on the assessment of the risk factors for dengue, with particular emphasis on South America, in an attempt to start to put the broader environmental considerations into a more detailed understanding and examination of the small-scale processes as they affect disease incidence.

In addition, the paper provides insight into techniques that can be used for future predictive models, particularly in Latin America and elsewhere. Techniques such as Machine Learning (ML) have experienced explosive growth that promises to revolutionize epidemiology and public health and offer new understandings of dengue and other infectious diseases.

Dengue outbreaks may occur when a high proportion of naïve subjects are concentrated in the same area. As the social and economic impacts are worsening and outbreaks are increasing, it becomes urgent to reinforce an integrated management for the surveillance, control, and prevention of dengue. One key aspect of this strategy is the ability to predict the occurrence of dengue outbreaks. An early warning of dengue outbreaks could improve the efficiency of vector control campaigns and help to target prevention actions. Such early interventions could delay or spread out the epidemic, thus reducing its impact on health system. Health facilities could adapt their response in terms of availability of beds and mobilization of human and material resources. Dengue morbidity and mortality would be minimized through earlier and more appropriate public health response.

## 2. LITERATURE SURVEY

Majeed, M.A.; Shafri, H.Z.M.; [1] dengue fever cases in Malaysia using machine learning techniques. A dataset consisting of weekly dengue cases at the state level in Malaysia from 2010 to 2016 was obtained from the Malaysia Open Data website and includes variables such as climate, geography, and demographics. Six different long short-term memory (LSTM) models were developed and compared for dengue prediction in Malaysia: LSTM, stacked LSTM (S-LSTM), LSTM with temporal attention (TA-LSTM), S-LSTM with temporal attention (STA-LSTM), LSTM with spatial attention (SA-LSTM), and S-LSTM with spatial attention (SSA-LSTM).

Cabrera, M.; Leake, J.; [2] epidemiological prediction of dengue fever using the One Health perspective, including an analysis of how Machine Learning techniques have been applied to it and focuses on the risk factors for dengue in Latin America to put the broader environmental considerations into a detailed understanding of the small-scale processes as they affect disease incidence. Determining that many factors can act as predictors for dengue outbreaks, a large-scale comparison of different predictors over larger geographic areas than those currently studied is lacking to determine which predictors are the most effective.

Dey, Samrat Kumar, et al. [3] develop a machine learning model that can use relevant information about the factors that cause Dengue outbreaks within a geographic region. To predict dengue cases in 11 different districts of Bangladesh, we created a DengueBD dataset and employed two machine learning algorithms, Multiple Linear Regression (MLR) and Support Vector Regression (SVR). This research also explores the correlation among environmental factors like temperature, rainfall, and humidity with the rise and decline trend of Dengue cases in different cities of Bangladesh. The entire dataset was divided into an 80:20 ratio, with 80 percent used for training and 20% used for testing. The research findings imply that, for both the MLR with 67% accuracy along with Mean Absolute Error (MAE) of 4.57 and SVR models with 75% accuracy along with Mean Absolute Error (MAE) of 4.95, the number of dengue cases reduces throughout the winter season in the country and increases mainly during the rainy season in the next ten months, from August 2021 to May 2022.

Kakarla, S.G., Kondeti, P.K., et al. [4] applied vector auto regression, generalized boosted models, support vector regression, and long short-term memory (LSTM) to predict the dengue prevalence in Kerala state of the Indian subcontinent. Consider the number of dengue cases as the target variable and weather variables viz., relative humidity, soil moisture, mean temperature, precipitation, and NINO3.4 as independent variables. Various analytical models have been applied on both datasets and predicted the dengue cases. Among all the models, the LSTM model was outperformed with superior prediction capability (RMSE: 0.345 and R2:0.86) than the other models.

Roster, Kirstin, et al. [5] developed a model for predicting monthly dengue cases in Brazilian cities 1 month ahead, using data from 2007–2019. We compared different machine learning algorithms and feature selection methods using epidemiologic and meteorological variables. They found that different models worked best in different cities, and a random forests model trained on monthly dengue cases performed best overall. It produced lower errors than a seasonal naive baseline model, gradient boosting regression, a feed-forward neural network, or support vector regression.

Sarder, Faysal, et al. [6] predict the accuracy of dengue outbreak from climate data. A dengue dataset, containing information of climate variables, dengue cases during 2019 to 2021 from Meteorology Department and Directorate General of Health Services (DGHS), Bangladesh. We split the whole dataset into 70:30 ratios were 70% considered as training and 30% for testing purposes. Such, prediction of accuracy we apply various supervised machine learning (ML) algorithms like Support Vector Machine (SVM), Decision Tree (DT), Logistic Regression (LR), Naïve Bayes (NB), AdaBoostClassifier (AdaBoost), XGBRegressor, GradientBoostingClassifier and Random Forest (RF). Finally, from these algorithms, SVM provide the highest accuracy of 96.73%.

Ochida, N., Mangeas, M., et al. [7] proposed statistical estimation of the effective reproduction number (Rt) based on case counts to create a categorical target variable: epidemic week/non-epidemic week. A machine learning classifier has been trained using relevant climate indicators in order to estimate the probability for a week to be epidemic under current climate data and this probability was then estimated under climate change scenarios.

Anuranjan, M. B., et al. [8] considered three different modelling techniques: interpolation, gradient boosting regression and random forest regression. Parameters were tuned and adjusted for optimal performance. Results are based on prediction accuracy and mean absolute error (MAE). The performance was analysed, and the result points out that the gradient boosting regression performs significantly better than the other models and is therefore considered to be a better approach. Future results can be improved by obtaining large amounts of meaningful data and implementing better models associated with time series predicting.

Gupta, G.; Khan, S.; et al. [9] developed dengue predictive models, data from microarrays and RNA-Seq have been used significantly. Bayesian inferences and support vector machine algorithms are two examples of statistical methods that can mine opinions and analyze sentiment from text. In general, these methods are not very strong semantically, and they only work effectively when the text passage inputs are at the level of the page or the paragraph; they are poor miners of sentiment at the level of the sentence or the phrase.

Rocha, F.P., Giesbrecht, M., et al. [10] models were trained with data from the municipality of São Luís do Maranhão, state of Maranhão, Brazil. The majority of related works analyze states, countries, or continental levels, with greater availability of data. To apply the approach to such a small region, some oversampling techniques were used. The number of cases per neighborhood from 2014 to and 2020 and climatic, territorial, and environmental data was used as input variables to estimate the probability of dengue occurrence in the municipality. Due to the unbalanced database, we used the SMOTE, ADASYN, and DBSMOTE oversampling techniques. The DBSMOTE-trained Random Forest classifier achieved the best results with a 75.1% AUC, 75.43% sensitivity and a 60.53% specificity.

[11] proposed a method that combines these three factors with data of Taiwanese dengue fever and uses the secondary area divided by the population as the granularity. Random Forest (RF) and XGBoost (XGB) are used for prediction model of weekly dengue fever infection area. Experimental results showed that the Receiver Operator Characteristic (ROC)/Area Under the Curve (AUC) of RF and XGB are both higher than 93%, and the Recall rate is higher than 80%. With the result, government can determine which area should do disinfection process to reduce the infection rate of dengue infection. Because of accurate prediction and disinfection process, the personnel cost can be reduced and it can prevent waste of medical recourse.

Pacheco, Paolo Ramon DC et al. [12] developed from machine learning algorithms were often used to provide accurate predictions as it can analyze trends from historical dengue data. Currently, the basis for the predictions of machine learning algorithms is unknown, which is why it is termed a "black-box". Climate-based random forest models were created through two implementations: randomForest package and ranger package. The ranger package slightly outperformed the conventional randomForest package, making it a better alternative in implementing random forest. Furthermore, models yielded more accurate predictions of dengue incidences with a delayed effect on the datasets. For local and global interpretations, most of the best models had relative humidity as the most influential to dengue incidence in Metropolitan Manila at all spatial scales.

## 3. PROPOSED SYSTEM

### 3.1 Pre-processing

**Data Preprocessing in Machine learning**

Data pre-processing is a process of preparing the raw data and making it suitable for a machine learning model. It is the first and crucial step while creating a machine learning model.

When creating a machine learning project, it is not always a case that we come across the clean and formatted data. And while doing any operation with data, it is mandatory to clean it and put in a formatted way. So, for this, we use data pre-processing task.
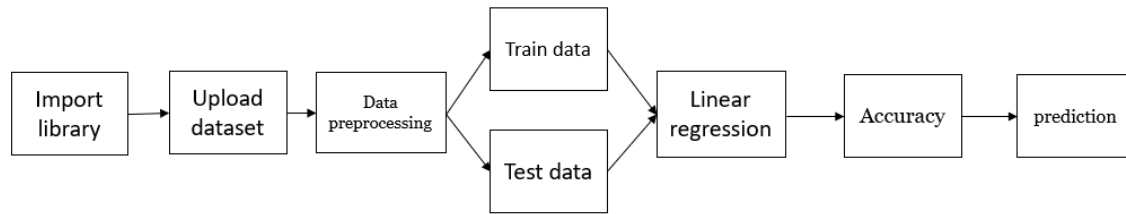


Fig. 1: Block diagram of proposed system.

## Why do we need Data Pre-processing?

A real-world data generally contains noises, missing values, and maybe in an unusable format which cannot be directly used for machine learning models. Data pre-processing is required tasks for cleaning the data and making it suitable for a machine learning model which also increases the accuracy and efficiency of a machine learning model.

- Getting the dataset
- Importing libraries
- Importing datasets
- Finding Missing Data
- Encoding Categorical Data
- Splitting dataset into training and test set
- Feature scaling

## Splitting the Dataset into the Training set and Test set

In machine learning data pre-processing, we divide our dataset into a training set and test set. This is one of the crucial steps of data pre-processing as by doing this, we can enhance the performance of our machine learning model.

Suppose if we have given training to our machine learning model by a dataset and we test it by a completely different dataset. Then, it will create difficulties for our model to understand the correlations between the models.

If we train our model very well and its training accuracy is also very high, but we provide a new dataset to it, then it will decrease the performance. So we always try to make a machine learning model which performs well with the training set and also with the test dataset. Here, we can define these datasets as:

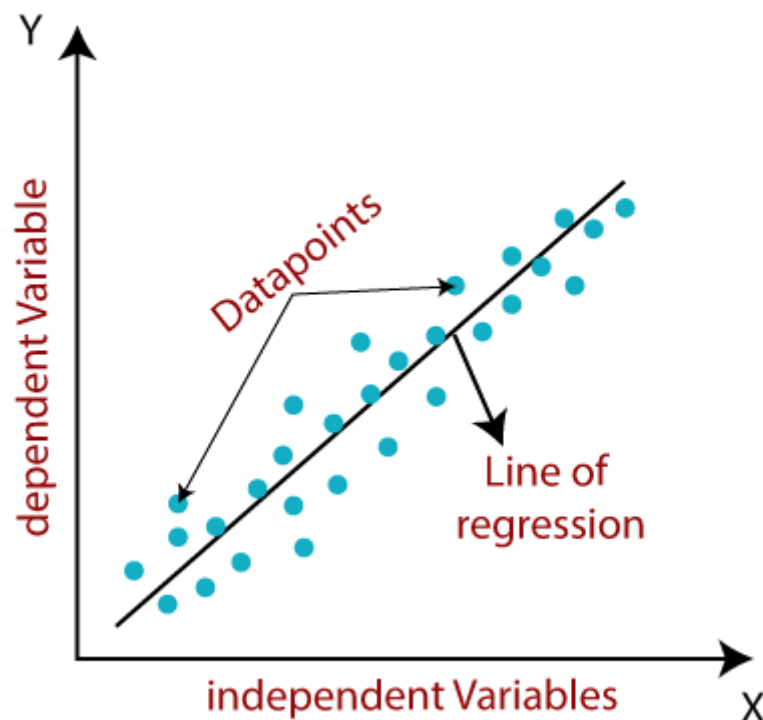**Training Set**: A subset of dataset to train the machine learning model, and we already know the output.

**Test set**: A subset of dataset to test the machine learning model, and by using the test set, model predicts the output.

### 3.2 Linear Regression

Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables such as **sales, salary, age, product price,** etc.

Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (y) variables, hence called as linear regression. Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

The linear regression model provides a sloped straight line representing the relationship between the variables. Consider the below image:



Mathematically, we can represent a linear regression as:

$$y = a_0 + a_1 x + \varepsilon$$

**Here,**

Y= Dependent Variable (Target Variable)
X= Independent Variable (predictor Variable)
a0= intercept of the line (Gives an additional degree of freedom)

a1 = Linear regression coefficient (scale factor to each input value).

ε = random error

The values for x and y variables are training datasets for Linear Regression model representation.
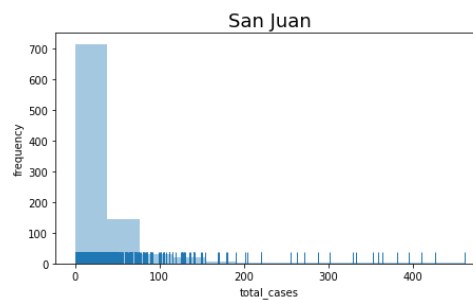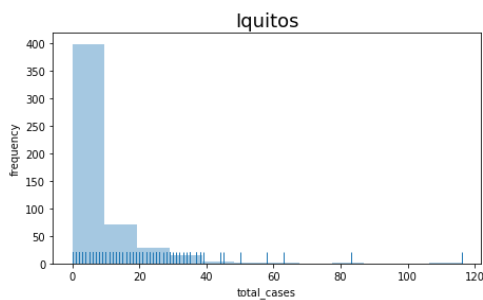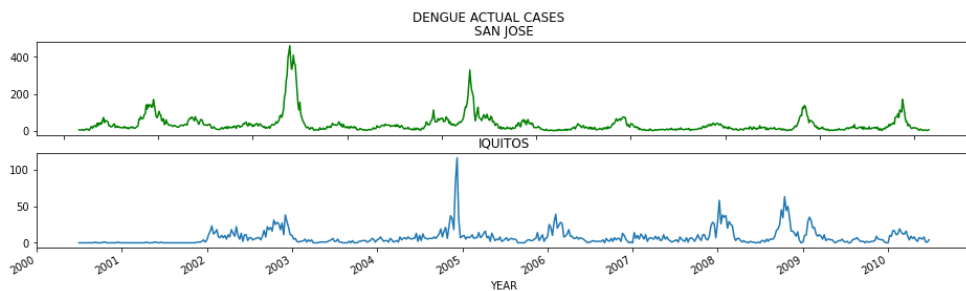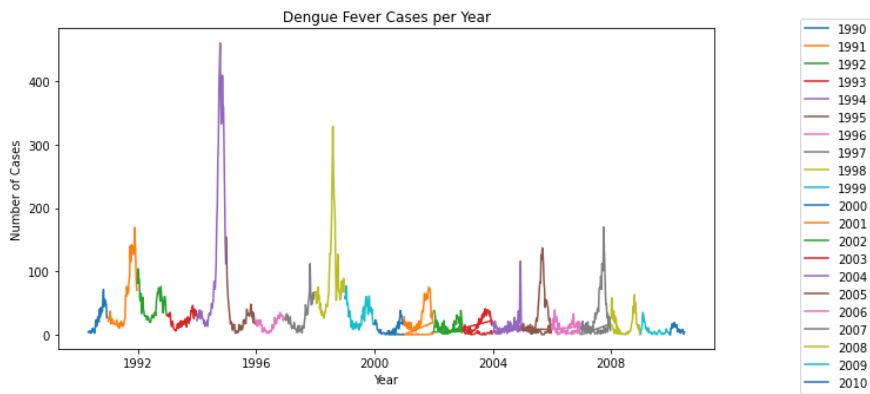
## 4. RESULTS AND DISCUSSION

Sample dataset

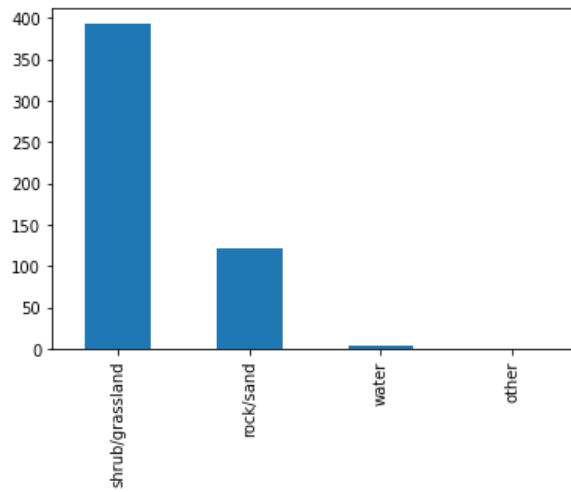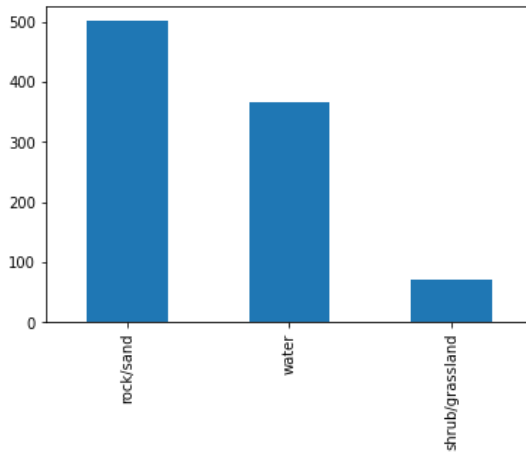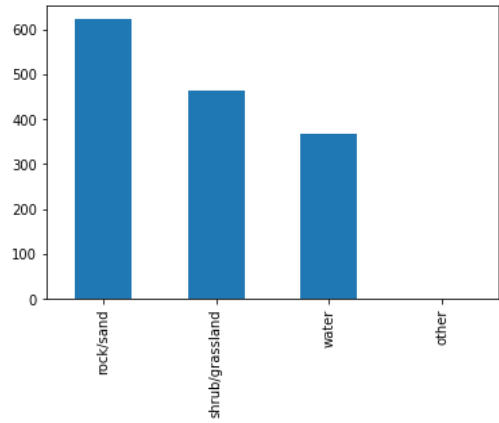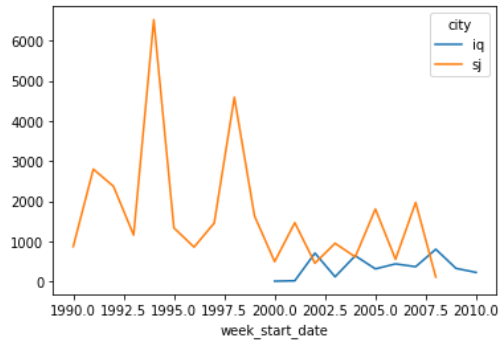| | city | year | weekofyear | week_start_date | ndvi_ne | ndvi_nw | ndvi_se | ndvi_sw | precipitation_amt_mm | reanalysis_air_temp_k | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | sj | 2008 | 18 | 2008-04-29 | -0.0189 | -0.018900 | 0.102729 | 0.091200 | 78.60 | 298.492857 | ... |
| 1 | sj | 2008 | 19 | 2008-05-06 | -0.0180 | -0.012400 | 0.082043 | 0.072314 | 12.56 | 298.475714 | ... |
| 2 | sj | 2008 | 20 | 2008-05-13 | -0.0015 | -0.012400 | 0.151083 | 0.091529 | 3.66 | 299.455714 | ... |
| 3 | sj | 2008 | 21 | 2008-05-20 | -0.0015 | -0.019867 | 0.124329 | 0.125686 | 0.00 | 299.690000 | ... |
| 4 | sj | 2008 | 22 | 2008-05-27 | 0.0568 | 0.039833 | 0.062267 | 0.075914 | 0.76 | 299.780000 | ... |

5 rows × 26 columns

| reanalysis_sat_precip_amt_mm | reanalysis_specific_humidity_g_per_kg | reanalysis_tdtr_k | station_avg_temp_c | station_diur_temp_rng_c | station_max_temp_c |
|---|---|---|---|---|---|
| 78.60 | 15.918571 | 3.128571 | 26.528571 | 7.057143 | 33.3 |
| 12.56 | 15.791429 | 2.571429 | 26.071429 | 5.557143 | 30.0 |
| 3.66 | 16.674286 | 4.428571 | 27.928571 | 7.785714 | 32.8 |
| 0.00 | 15.775714 | 4.342857 | 28.057143 | 6.271429 | 33.3 |
| 0.76 | 16.137143 | 3.542857 | 27.614286 | 7.085714 | 33.3 |

| station_min_temp_c | station_precip_mm | month | total_cases_avg |
|---|---|---|---|
| 21.7 | 75.2 | 4 | 10.722222 |
| 22.2 | 34.3 | 5 | 9.944444 |
| 22.8 | 3.0 | 5 | 11.500000 |
| 24.4 | 0.3 | 5 | 11.166667 |
| 23.3 | 84.1 | 5 | 13.777778 |



340

San Juan Variable Correlations



Iquitos Variable Correlations



Dengue Fever Cases per Year



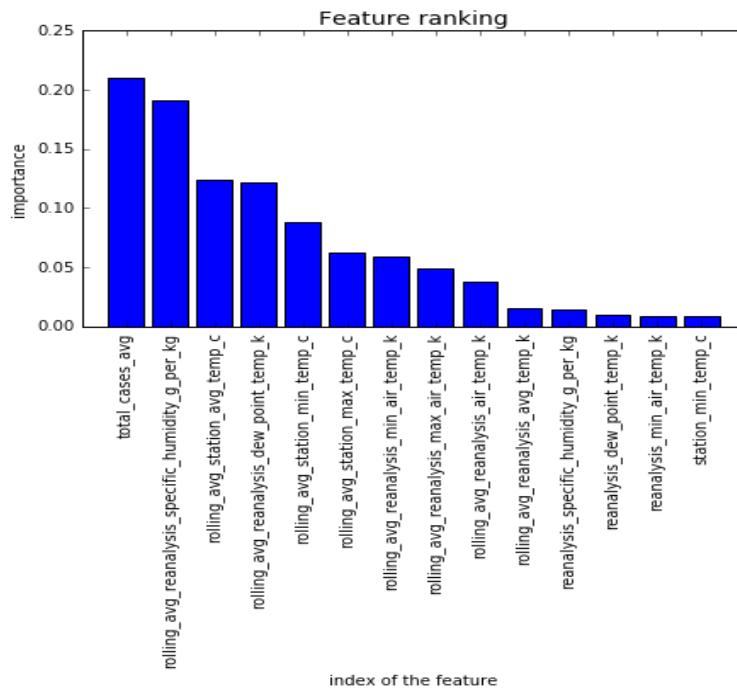DENGUE ACTUAL CASES
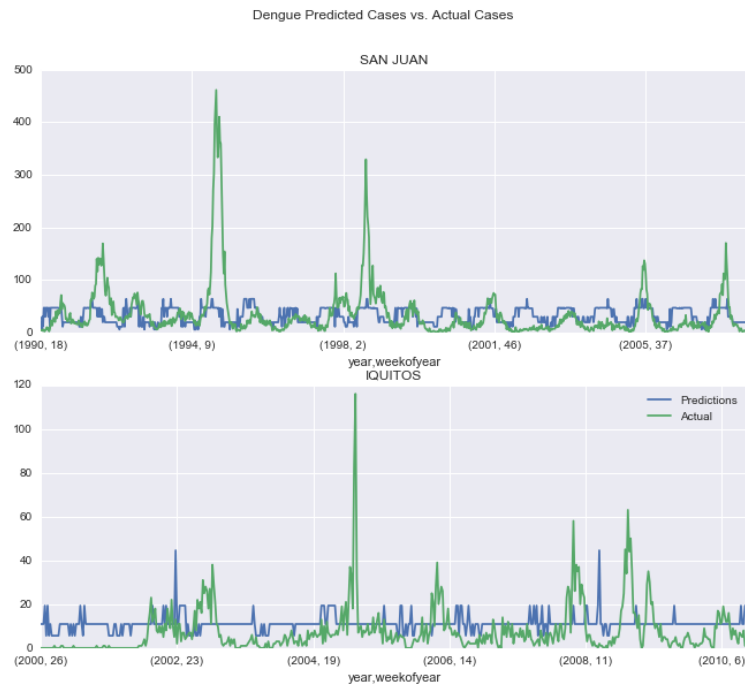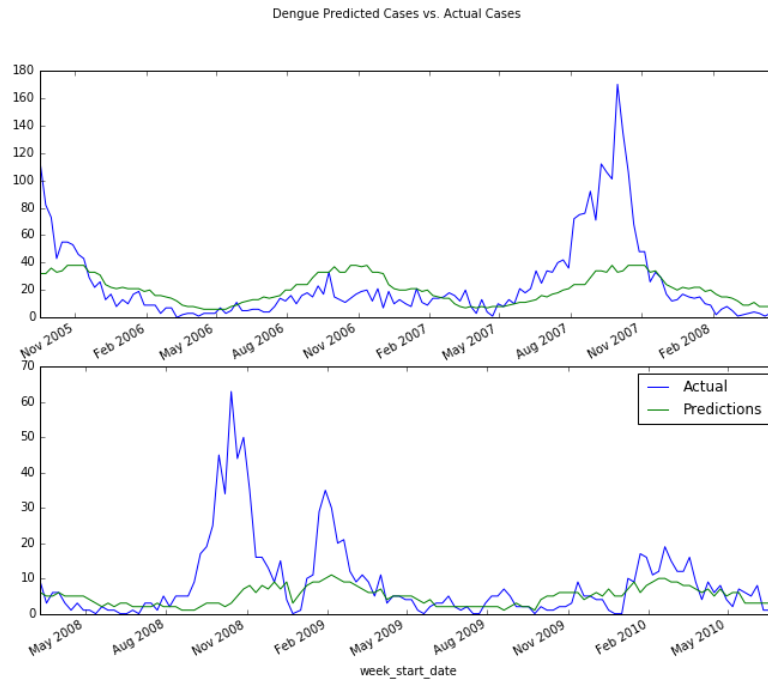SAN JOSE

IQUITOS

Feature ranking on San Juan dataset



Feature ranking on Iquitos dataset



**Prediction using support vector machine:** clearly the model is unable to predict the spikes in disease outbreak.

Dengue Predicted Cases vs. Actual Cases



Dengue Predicted Cases vs. Actual Cases

## 5. CONCLISION

In conclusion, dengue fever is a significant human arbovirus disease transmitted by female mosquitoes of the Aedes genus, primarily Aedes aegypti and Ae. albopictus. It is prevalent in tropical and subtropical regions, posing a considerable risk to approximately 2.5 billion people worldwide. The disease manifests in two major clinical forms: mild dengue fever and severe dengue, characterized by plasma leakage with or without hemorrhage. The prevalence of dengue has dramatically increased in recent decades, with an annual incidence of 50 to 100 million infections and over 500,000 hospitalizations globally.

Dengue haemorrhagic fever (DHF), a severe form of the disease, is on the rise in many tropical regions and is responsible for approximately 20,000 deaths annually, primarily among children under 15 years of age. The disease is endemic in various countries, and population movement plays a crucial role in the dissemination of the virus. It contributes to the introduction of new virus strains and the introduction of nonimmune individuals into endemic areas, further exacerbating the spread of the disease.

To address the challenges associated with the spread of dengue fever, a proposed system utilizes climate data and employs time series analysis to predict its transmission. By leveraging machine learning algorithms and advancements, the system aims to forecast the spread of dengue fever accurately. Additionally, the project incorporates exploratory data analytics to gain insights into the dengue dataset over time.

**Future scope**

The future scope for addressing the spread of dengue fever lies in the integration of advanced technologies such as artificial intelligence, big data analytics, and IoT (Internet of Things). By combining real-time climate data, population movement patterns, and predictive modeling algorithms, it becomes possible to develop sophisticated early warning systems that can accurately forecast dengue outbreaks and enable proactive measures for prevention and control. Furthermore, leveraging mobile applications and wearable devices can facilitate real-time monitoring of mosquito populations, disease surveillance, and targeted interventions, while community engagement and education programs can raise awareness and promote effective preventive measures. This holistic approach, empowered by emerging technologies, has the potential to revolutionize the fight against dengue and significantly reduce the burden of the disease on affected communities.

**REFERENCES**

[1] .      Majeed, M.A.; Shafri, H.Z.M.; Zulkafli, Z.; Wayayok, A. A Deep Learning Approach for Dengue Fever Prediction in Malaysia Using LSTM with Spatial Attention. Int. J. Environ. Res. Public Health 2023, 20, 4130. https://doi.org/10.3390/ijerph20054130

[2] .      Cabrera, M.; Leake, J.; Naranjo-Torres, J.; Valero, N.; Cabrera, J.C.; Rodríguez-Morales, A.J. Dengue Prediction in Latin America Using Machine Learning and the One Health Perspective: A Literature Review. Trop. Med. Infect. Dis. 2022, 7, 322. https://doi.org/10.3390/tropicalmed7100322

[3] .      Dey, Samrat Kumar, et al. "Prediction of dengue incidents using hospitalized patients, metrological and socio-economic data in Bangladesh: A machine learning approach." PLoS One 17.7 (2022): e0270933.

[4] .      Kakarla, S.G., Kondeti, P.K., Vavilala, H.P. et al. Weather integrated multiple machine learning models for prediction of dengue prevalence in India. Int J Biometeorol 67, 285–297 (2023). https://doi.org/10.1007/s00484-022-02405-z

[5] .      Roster, Kirstin, Colm Connaughton, and Francisco A. Rodrigues. "Machine-Learning–Based Forecasting of Dengue Fever in Brazilian Cities Using Epidemiologic and Meteorological Variables." American Journal of Epidemiology 191.10 (2022): 1803-1812.

[6] .      Sarder, Faysal, Sorefa Akter, and Sharmin Akter. "Predicting Dengue Outbreak from Climate Data Using Machine Learning Algorithms." 2022 IEEE International Conference on Data Science and Information System (ICDSIS). IEEE, 2022.

[7] .      Ochida, N., Mangeas, M., Dupont-Rouzeyrol, M. et al. Modeling present and future climate risk of dengue outbreak, a case study in New Caledonia. Environ Health 21, 20 (2022). https://doi.org/10.1186/s12940-022-00829-z.

[8] .          Anuranjan, M. B., et al. "Machine Learning Techniques for Predicting Dengue Outbreak." Innovations in Information and Communication Technologies: Proceedings of ICIICT 2022. Singapore: Springer Nature Singapore, 2022. 45-56.

[9] .          Gupta, G.; Khan, S.; Guleria, V.; Almjally, A.; Alabduallah, B.I.; Siddiqui, T.; Albahlal, B.M.; Alajlan, S.A.; AL-subaie, M. DDPM: A Dengue Disease Prediction and Diagnosis Model Using Sentiment Analysis and Machine Learning Algorithms. Diagnostics 2023, 13, 1093. https://doi.org/10.3390/diagnostics13061093.

[10] .          Rocha, F.P., Giesbrecht, M. Machine learning algorithms for dengue risk assessment: a case study for São Luís do Maranhão. Comp. Appl. Math. 41, 393 (2022).

[11] .          Zheng, Cong-Han, et al. "Predicting Infection Area of Dengue Fever for Next Week Through Multiple Factors." Advances and Trends in Artificial Intelligence. Theory and Practices in Artificial Intelligence: 35th International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems, IEA/AIE 2022, Kitakyushu, Japan, July 19–22, 2022, Proceedings. Cham: Springer International Publishing, 2022.

[12] .          Pacheco, Paolo Ramon DC. "Application of interpretable machine learning in revealing spatial and temporal patterns of dengue disease using climatic factors." (2022).