# Dengue Fever Spread Predicting Based on Ensemble Learning and Markov Chain Monte Carlo (MCMC) Method

Han Weiyu

October 31, 2023

## Abstract

This study delves into the predictive modeling of Dengue fever transmission, a mosquito-borne disease prevalent in tropical and subtropical regions. Recognizing the potential shifts in disease distribution due to climate change, we employed environmental data collected by various U.S. federal departments to construct a machine learning-based predictive model. After comparing traditional statistical methods with machine learning/deep learning techniques, we opted for a machine learning model for its superior interpretability. To enhance prediction accuracy, we incorporated ensemble learning techniques and combined them with Bayesian methods, specifically the Markov Chain Monte Carlo (MCMC) approach, to ascertain the confidence intervals of our predictions. The model demonstrated commendable performance on the test set, with a Mean Absolute Error (MAE) of 6.39 and a Mean Squared Error (MSE) of 113.50. This research offers an effective methodology for predicting Dengue fever transmission, aiding global public health sectors in better addressing this challenge.

**Keywords:** Machine Learning; Bayesian Method; Epidemicology; Ensemble Learning; Dengue Fever

## 1   Introduction

Dengue fever, a mosquito-borne ailment prevalent in tropical and sub-tropical regions, manifests mild symptoms akin to flu in early stages - fever, rash, and muscle/joint pain. However, severe cases can escalate to intense bleeding, plummeting blood pressure, and potentially, fatality. The disease's transmission dynamics are intertwined with climate variables like temperature and precipitation. Although the climate-disease nexus is intricate, an increasing scientific consensus posits that climate change might induce distributional shifts with notable global public health repercussions. Utilizing environmental data amassed by various U.S. Federal Government entities, including the Centers for Disease Control and Prevention to the National Oceanic and Atmospheric Administration, a prediction model is constructed leveraging robust machine learning instruments.

The predominant methodologies for disease spread prediction encompass traditional statistical approaches and machine learning/deep learning techniques[1]. The former includes standard dynamic epidemiological models such as Susceptible-Infected-Recovered (SIR)[2][3][4], their refined versions[5][6], and time-series prediction models[7]. Some studies juxtapose these traditional methods against machine learning and deep learning, revealing superior predictive accuracy in the latter[2]. Consequently, machine learning and deep learning have gained traction in epidemic forecasting. Notably, deep learning models exhibit proficient time-series data prediction[8][9], while machine learning also demonstrates efficacy in epidemic prediction[10]. Comparative studies further elucidate the relative performance of these methodologies in epidemic forecasting[11]. In this research, machine learning models are selected for their robust interpretability.

Given the inherent limitations of machine learning models, feature engineering emerges as a crucial preliminary step prior to the execution of any fitting and training procedures. In this research, our endeavor is to forecast the incidence of dengue fever at a subsequent time point, predicated on the case data from preceding time points. Hence, the cultivation of more insightful features is imperative to enhance the model's comprehension of the problem at hand. Certain studies suggest that the interaction term of temperature and precipitation could potentially augment the predictive accuracy concerning the spread of dengue fever[12]. In light of this, a plethora of such interaction features are incorporated into the model, aiming to refine its predictive performance substantially.

Following the feature engineering process, to enhance the prediction accuracy for dengue fever incidence, we draw inspiration from the work of Tijana Radivojević et al.[13] and construct an ensemble learning model aimed at forecasting the case count for specific days. Subsequently, to address the over-reliance of the model on the values from previous time orders, we employ a Bayesian approach, specifically the Markov Chain Monte Carlo (MCMC) method, to ascertain the confidence intervals of our predictions. Ultimately, we introduce two evaluation metrics analogous

to precision and recall rates, to rigorously assess the predictive performance of our models.

## 2 Methods

### Data Availability

The dataset utilized in this study was sourced from the DengAI: Predicting Disease Spread competition, as detailed in (14) and further discussed in (15). This dataset encapsulates recorded instances of dengue fever cases in San Juan and Iquitos over a span from 1990 to 2010. Data collection was conducted on a weekly basis, rendering it as time-series data with a total of 1456 records. Any missing values within the dataset have been designated as NaNs (Not a Number). A total of 25 features have been documented, which are elucidated in Table 1 below. The primary focus of our investigation centers on the data pertaining to San Juan, encompassing a total of 936 records.

Table 1: Dataset features description

| Feature | Description |
| --- | --- |
| city | City abbreviations: sj for San Juan and iq for Iquitos |
| year | The year of record |
| weekofyear | The week of the year |
| week_start_date | Date given in yyyy-mm-dd format |
| station_max_temp_c | Maximum temperature (NOAA's GHCN) |
| station_min_temp_c | Minimum temperature (NOAA's GHCN) |
| station_avg_temp_c | Average temperature (NOAA's GHCN) |
| station_precip_mm | Total precipitation (NOAA's GHCN) |
| station_diur_temp_rng_c | Diurnal temperature range (NOAA's GHCN) |
| precipitation_amt_mm | Total precipitation (PERSIANN satellite) |
| reanalysis_sat_precip_amt_mm | Total precipitation (NOAA's NCEP) |
| reanalysis_dew_point_temp_k | Mean dew point temperature (NOAA's NCEP) |
| reanalysis_air_temp_k | Mean air temperature (NOAA's NCEP) |
| reanalysis_relative_humidity_percent | Mean relative humidity (NOAA's NCEP) |
| reanalysis_specific_humidity_g_per_kg | Mean specific humidity (NOAA's NCEP) |
| reanalysis_precip_amt_kg_per_m2 | Total precipitation (NOAA's NCEP) |
| reanalysis_max_air_temp_k | Maximum air temperature (NOAA's NCEP) |
| reanalysis_min_air_temp_k | Minimum air temperature (NOAA's NCEP) |
| reanalysis_avg_temp_k | Average air temperature (NOAA's NCEP) |
| reanalysis_tdtr_k | Diurnal temperature range (NOAA's NCEP) |
| ndvi_se | Pixel southeast of city centroid (NDVI) |
| ndvi_sw | Pixel southwest of city centroid (NDVI) |
| ndvi_ne | Pixel northeast of city centroid (NDVI) |
| ndvi_nw | Pixel northwest of city centroid (NDVI) |

### Data Pre-processing

Upon examination, we identified a semblance of feature duplication, seemingly encapsulating singular phenomena. To address this probable redundancy, it is imperative to discern the distinctions between the GHCN (Global Historical Climatology Network) and NCEP (National Centers for Environmental Prediction), as the primary divergence arises from these two recording methodologies. Following an investigative pursuit, it was ascertained that GHCN data is procured directly from terrestrial weather stations, while NCEP data is derived through the recalibration of satellite-captured data via specified algorithms. The NCEP data is quantified on a 0.5x0.5 degree scale, encapsulating approximately a 5.5km x 5.5km area. Contrarily, San Juan spans an area of 445 km². In our assessment, on one hand, our objective is to ascertain the influence of ground-level temperature and other variables on the propagation of dengue fever; on the other hand, the recording area of NCEP appears significantly minuscule in comparison to the expanse of San Juan. Consequently, in instances of conflicting data, we have opted to exclude NCEP data in favor of

the more locally pertinent GHCN data.

Having eliminated features potentially causing duplication, our objective transitioned to assessing whether varying measurements of similar concepts could engender redundancy. Specifically, two concepts, precipitation and humidity, were identified as potential culprits. We employed linear models to evaluate the relationship between similar measurements within these two concepts. However, the obtained R-squared values were 0.23 for precipitation and 0.45 for humidity, indicating suboptimal fitting. These results suggest a lack of linear relationship between the respective measurements of these concepts. Consequently, we discerned no necessity to eliminate any of these measurements, thereby retaining them for further analysis.

Following the steps outlined above, our focus transitioned towards gaining a more nuanced understanding of the dataset to adeptly manage the missing data. In line with this, we initiated an examination of the distribution of various features within this dataset, classifying the features with missing data into three distinct groups.

The first group encapsulates features that exhibit an approximately normal distribution and contain only a scant amount of missing data compared to the overall volume of the dataset. Notably, these features encompass reanalysis_relative_humidity_percent, reanalysis_specific_humidity_g_per_kg, station_avg_temp_c, station_diur_temp_rng_c, station_max_temp_c, station_min_temp_c and reanalysis_dew_point_temp_c. In this instance, we opted to use the mean value to impute these null values.

The second group comprises features that have a sparse proportion of missing data yet exhibit significant skewness. These features include reanalysis_precip_amt_kg_per_m2 and station_precip_mm. However, upon plotting how these features vary with time, we discerned the cycle of peak occurrences and noted that the missing data did not seem to occur at these potentially critical peak points. Hence, utilizing the mean value to fill the missing data emerged as a feasible solution.

The third group encompasses features with a relatively large volume of missing data, such as ndvi_ne, ndvi_nw, ndvi_se, ndvi_sw. Specifically, for ndvi_ne, there are 191 missing data points, and nearly half of the records for the year 1994 are absent. Although these features do not exhibit substantial fluctuations over time, using the mean value of the recording window or even the annual mean are not deemed adequate methods for imputation. After decomposing the time-series data into trend, seasonal, and residual components, it became apparent that there is nearly no seasonal cycle present in these four features. Consequently, linear interpolation emerged as a commonly adopted method for addressing missing data in time-series datasets.

## Feature Engineering

Prior to conducting any feature engineering, it is imperative to partition the dataset into a training set (80%) and a test set (20%). This segregation is essential to preclude data leakage. The process of feature engineering is conducted solely on the training set and encompasses three distinct steps.

Initially, the creation of lagged variables is undertaken. This is a common practice in the realm of time-series data analysis. The data in time-series is markedly influenced by the values of preceding time points, as the data trajectory is predicated on these prior values. Specifically, the lagged values of 1, 2, 4, 8, and 12 previous time periods for variables such as station_avg_temp_c, station_precip_mm, and total_cases are generated, given their noted significance in previous studies.

Subsequently, interaction features are created encompassing all temperature and precipitation variables. This step is motivated by the hypothesis that these features may exhibit a synergistic effect on disease propagation, as delineated in prior regional studies.

Lastly, the integration of seasonal cycle features into the model is pursued. An analysis of the seasonal trend of case occurrences is conducted, and significant seasonal features discerned within an annual cycle are extracted. The similarities across different cycles are noteworthy; hence, the information gleaned from the initial year's cycle is utilized to construct a new feature. Nonetheless, certain years comprise 53 weeks; for these instances, the data for

the 53rd week is generated by replicating the data of the 1st week. The procedure for extracting seasonal features is illustrated in Figure 1 below.
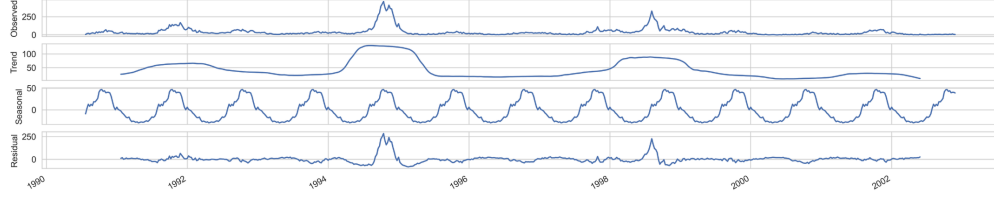


Figure 1: Seasonal trend

After creating these features, we need to select the really important features. A random forest model is trained to help us filter the features. As we set a threshold of the median of all features' importance, finally we get 19 features as follows.

- Week of Year: weekofyear
- NW Normalized Difference Vegetation Index: ndvi_nw
- SE Normalized Difference Vegetation Index: ndvi_se
- Specific Humidity (g per kg): reanalysis_specific_humidity_g_per_kg
- Diurnal Temperature Range (°C): station_diur_temp_rng_c
- Dew Point Temperature (°C): reanalysis_dew_point_temp_c
- 2-Week Lagged Average Temperature (°C): station_avg_temp_c_lag2
- 4-Week Lagged Average Temperature (°C): station_avg_temp_c_lag4
- 8-Week Lagged Average Temperature (°C): station_avg_temp_c_lag8
- 12-Week Lagged Average Temperature (°C): station_avg_temp_c_lag12
- 1-Week Lagged Precipitation (mm): station_precip_mm_lag1
- 4-Week Lagged Precipitation (mm): station_precip_mm_lag4
- 1-Week Lagged Total Cases: total_cases_lag1
- 2-Week Lagged Total Cases: total_cases_lag2
- 4-Week Lagged Total Cases: total_cases_lag4
- 8-Week Lagged Total Cases: total_cases_lag8
- 12-Week Lagged Total Cases: total_cases_lag12
- Interaction of Temperature and Precipitation: station_diur_temp_rng_c*reanalysis_precip_amt_kg_per_m2
- Seasonal Feature: seasonal_feature

## Definition of Prevalent

In order to provide more impactful recommendations, our objective is to concentrate on episodes of higher transmission rates rather than days with minimal case occurrences. The delineation of what constitutes a high transmission rate, however, merits careful consideration. Fortuitously, an examination of the case distribution revealed a pronounced skewness, indicating that a majority of the case occurrences are clustered at lower levels.

To elucidate this further, we devised a plot illustrating the frequency of specific case counts alongside the cumulative proportion of the total case count, both sharing a common x-axis which represents the specific number of cases. This graphical representation yielded an 'elbow point,' a point at which the curve exhibits a noticeable change

in direction. The 'elbow point' on the curve signifies that the specific numbers of cases that occurred less than 5 times collectively account for roughly 60% of the total case count. This pivotal observation led us to designate the instances of cases occurring less than 5 times as prevalent episodes of transmission. Consequently, a threshold of 32 cases was established to define prevalence.

Figure 2 below delineates this analysis, providing a visual representation of the frequency and cumulative proportion of cases against the specific number of cases. This visualization not only aids in understanding the distribution of case occurrences but also substantiates our rationale for setting a threshold at 32 cases to demarcate prevalent episodes of transmission. Through such analytical scrutiny, we aim to hone our focus on the most significant transmission episodes, thereby enabling a more effective formulation of recommendations to mitigate the spread of the disease.
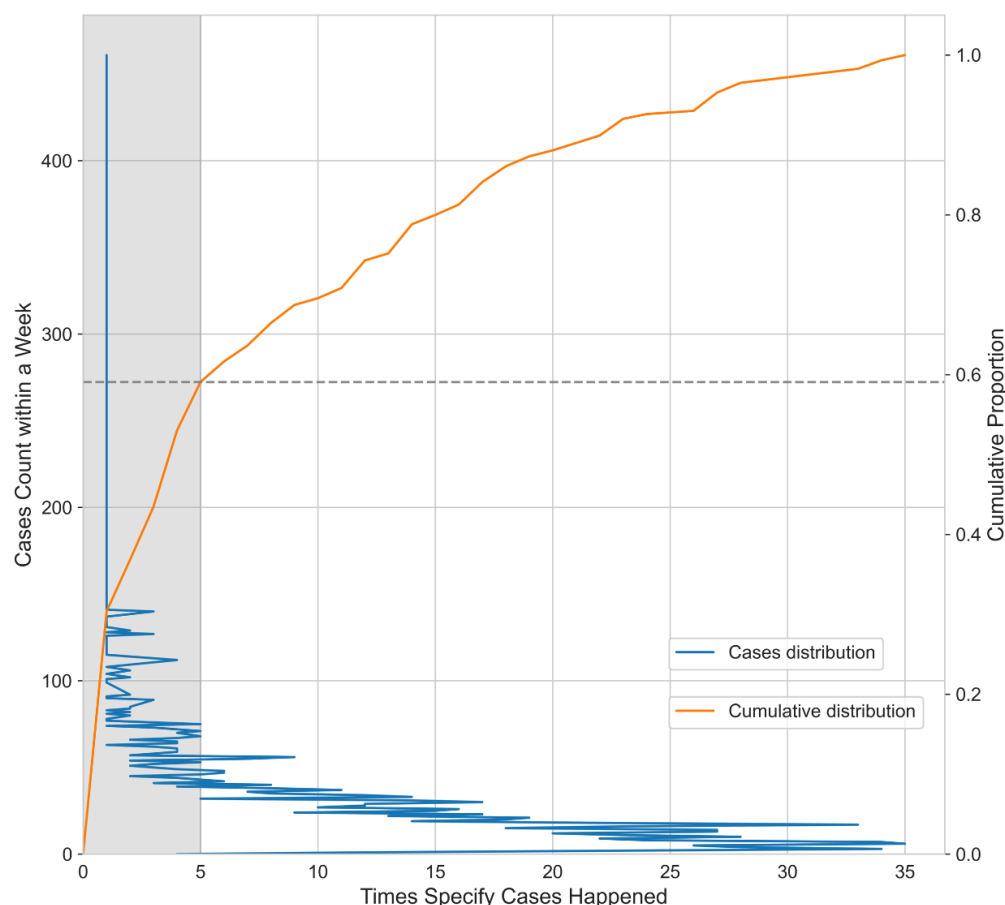


Figure 2: Prevalent threshold

## Model Building

In this investigation, our aim is to construct an ensemble learning framework encompassing seven foundational learners, namely, RandomForestRegressor, MLPRegressor, SVR, KernelRidge, GaussianProcessRegressor, GradientBoostingRegressor, and TPOTRegressor. Subsequent to the establishment of this ensemble model, a Bayesian approach is employed to ascertain the confidence intervals of the predictions. Ultimately, this methodology enables the derivation of a predictive distribution of cases. The architectural schema of the model is delineated in Figure 3 below. This composite framework seeks provide a more nuanced understanding of the case distribution, thereby enhancing the robustness and reliability of the forecasted outcomes.
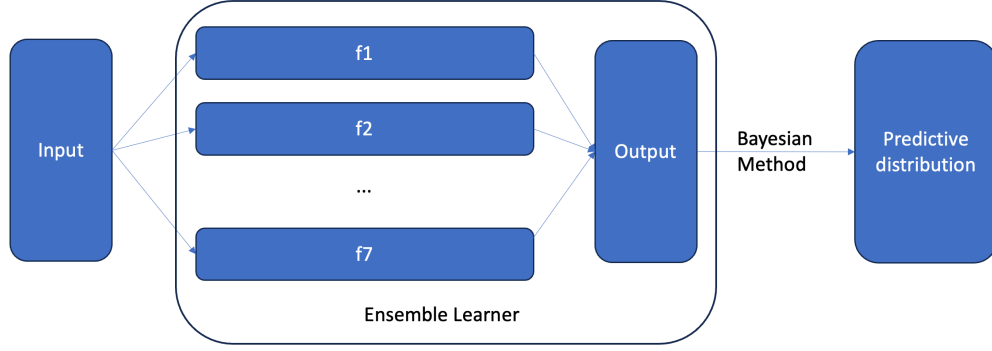
Figure 3: Model structure

# 3 Results

## Model Training and Tuning

The training dataset was partitioned into a training set and a validation set to facilitate the model training and validation process. The model was trained utilizing the training set, followed by prediction exercises on the validation set. As illustrated in Figure 4, it can be deduced that the performance of the ensemble model was adversely impacted by two base learners: SVR and GaussianProcessRegressor. Consequently, these two models were excluded from the ensemble learning model configuration.
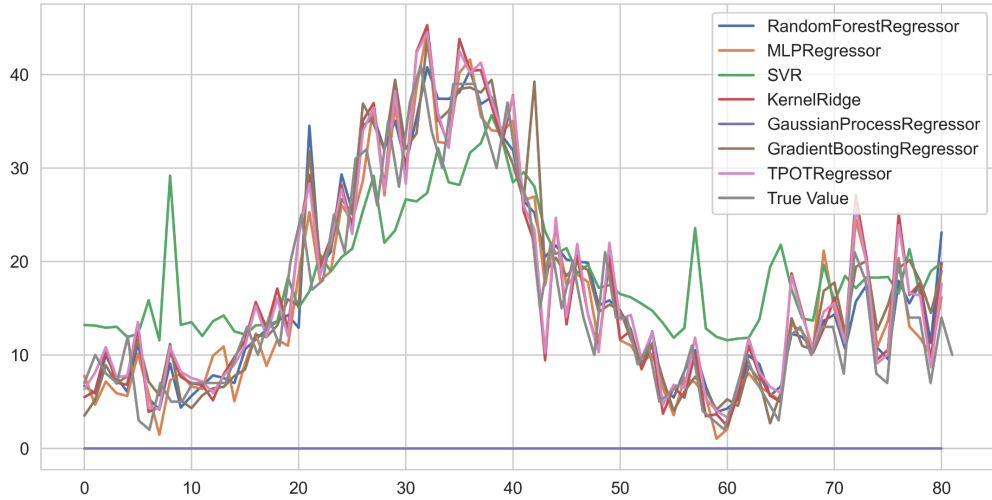


Figure 4: Performance of different base learners

Furthermore, a notable observation from Figure 5 elucidates that the model exhibited a significant reliance on the preceding time order's cases, with a tendency to minimally adjust the prior time order's cases for prediction. This behavior deviates from the desired model characteristics. To bolster the robustness of our model and mitigate its sensitivity, the confidence interval of the prediction was incorporated to expand the decision boundary of the prediction.
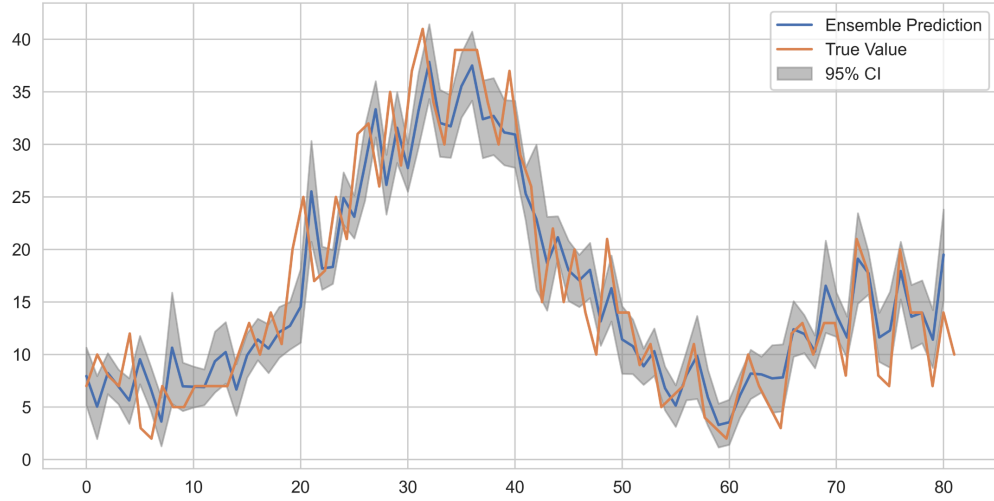
Figure 5: Performance of different base learners

In essence, the definition of 'prevalent', as elucidated earlier, underpins our objective to issue warnings whenever the 95% confidence interval of the succeeding time order's dengue fever cases prediction encompasses 32 cases, as depicted in Figure 6.
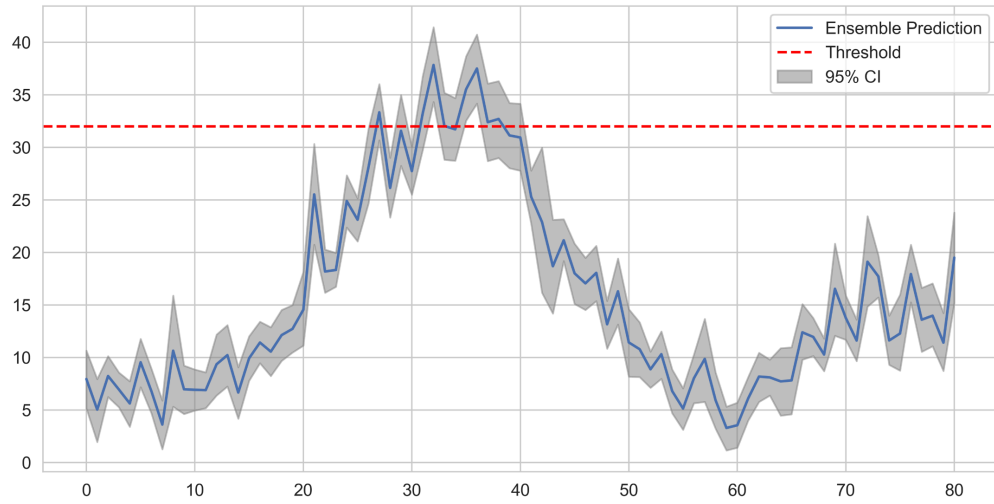


Figure 6: Prediction with confidence interval

## Prediction Result

Following the model training and tuning phases, the training and validation sets were amalgamated for a retraining exercise on the model. Subsequently, the performance of the model was evaluated on the test set. The outcomes of this evaluation are presented in Figure 7. The model yielded a Mean Absolute Error (MAE) of 6.39 and a Mean Square Error (MSE) of 113.50.
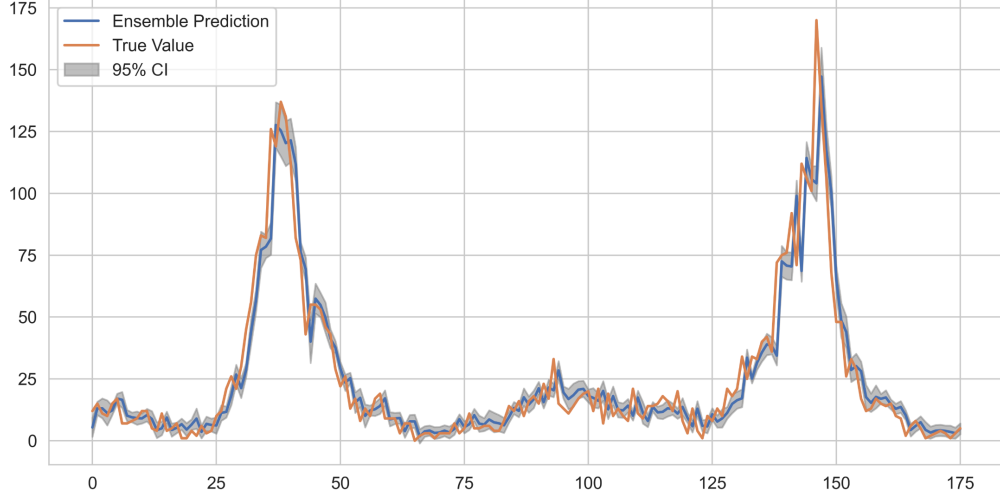
Figure 7: Prediction on test set

## Other Measurements

As articulated earlier, our objective is to trigger alerts whenever the 95% confidence interval of the predicted dengue fever cases for the subsequent time period encompasses 32 cases. It is imperative to establish metrics to evaluate the efficacy of these warnings. Issuing excessive warnings that do not correspond to actual prevalent cases could lead to an unwarranted allocation of public resources. Conversely, a model that fails to identify potential prevalence is deemed ineffective. Therefore, we have devised two metrics to assess the performance of the model in this regard, delineated as follows:

$$P1 = \frac{\text{Cases exceeding 32 within all warnings}}{\text{Total warnings issued}}$$

$$R2 = \frac{\text{Warnings issued for all cases exceeding 32}}{\text{Total cases exceeding 32}}$$

The performance of the model on the test set yielded a $P1$ score of $0.88$ and an $R2$ score of $0.9$. These results suggest that the model demonstrates a commendable efficacy in issuing advance warnings, approximately a week prior, to potential surges in dengue fever cases.

## 4 Discussion

### Contributions

The endeavor to model and predict the spread of Dengue Fever in this study has resulted in several innovative strides that contribute a novel perspective to the existing body of knowledge in epidemiological modeling. The highlights of the innovations are elucidated as follows:

**Employment of Ensemble Learning Models:** Unlike conventional approaches that often rely on singular machine learning models, this study leverages the power of ensemble learning to amalgamate the strengths and mitigate the weaknesses inherent in individual models. By orchestrating an ensemble of diverse learning algorithms, the resultant model demonstrates a robust capability to capture the underlying patterns in the data with a heightened level of accuracy and reliability.

**Integration of Bayesian Methodology:** To tackle the challenge of overly narrow prediction intervals which could unduly constrain the interpretability and applicability of the predictions, a Bayesian framework is ingeniously integrated into the model. This Bayesian infusion not only broadens the prediction intervals to a more realistic range but also endows the model with a measure of tolerance, thereby enhancing the model's capability to provide actionable

insights amidst uncertainties inherent in epidemiological data.

**Introduction of Customized Evaluation Metrics $P1$ and $R2$:** To finely gauge the efficacy of the model, especially in the critical task of issuing timely warnings, bespoke evaluation metrics $P1$ and $R2$ are crafted. The model's satisfactory performance on these metrics, with a $P1$ score of 0.88 and $R2$ score of 0.9 on the test set, underscores its proficiency in discerning potential outbreaks, thereby aligning well with the primary objectives of the study.

**Achievement of Exceptional Predictive Accuracy:** The model exhibits stellar predictive accuracy, with a Mean Absolute Error (MAE) of a mere 6.39 on the test set, a level of precision that many extant models struggle to attain. This refined level of accuracy is instrumental in bolstering the model's utility as a reliable tool for anticipating dengue fever outbreaks, thus potentially aiding in the allocation and mobilization of public health resources in a timely and efficient manner.

These innovative facets significantly bolster the potential impact and practical applicability of the study. They also serve as a testament to the judicious blend of machine learning techniques, statistical methodologies, and domain-specific insights employed in this inquiry. The promising results beckon further exploration and refinement in future research endeavors, with the aim of advancing the frontier of predictive modeling in public health epidemiology.

## Limitations

In this inquiry into the dynamics of Dengue Fever spread, the employed methodologies and data analytics pipelines have demonstrated substantial promise. Nevertheless, several limitations and avenues for future exploration have emerged which merit attention to enhance the robustness and predictive accuracy of the models.

**Advanced Model Exploration:** The adoption of more sophisticated machine learning and statistical models could potentially unveil intricate patterns within the data that simpler models might overlook. Investigating the merits of models that have shown proficiency in handling time-series data or those designed to tackle epidemiological data could be of significant value.

**Environmental Variable Augmentation:** Our current model harnesses a selected set of environmental variables. However, the interplay between climatic conditions and vector-borne diseases like Dengue Fever is complex. Introducing a more comprehensive set of environmental variables, potentially encompassing variables like wind speed, water storage levels, and urbanization metrics, could provide a more nuanced understanding of the transmission dynamics.

**Bayesian Approach Refinement:** The Bayesian framework employed for deriving prediction intervals could be refined further. A more nuanced utilization of prior information, and exploring alternative Bayesian computation techniques, might yield more accurate and informative prediction intervals, thus improving the model's utility for real-world decision-making scenarios.

**Hyperparameter Optimization:** The performance of machine learning models is often significantly impacted by the choice of hyperparameters. A more rigorous approach towards hyperparameter tuning, possibly through grid search or Bayesian optimization, might lead to models that better fit the underlying data structure, thereby enhancing predictive performance.

**Lag_1 Data Dependency Reduction:** Our model exhibits a pronounced dependency on the lag_1 data. This dependency might be masking other informative patterns in the data. Exploring models that reduce this reliance, or alternatively, investigating the merits of omitting the lag_1 data, could potentially lead to a more balanced and generalized model.

**Extension to Iquitos:** The current investigation is primarily centered around data from San Juan. Extending the analysis to cover dengue cases in Iquitos could provide a more holistic understanding of the disease spread dynamics, and test the model's generalizability across varying geographical and climatic conditions.

These limitations not only offer a transparent reflection on the current study but also pave the way for future research endeavors. Addressing these identified areas could significantly augment the predictive prowess of the models and contribute to a deeper comprehension of Dengue Fever spread dynamics, a crucial step towards better public health preparedness and response strategies.

# 5 Code availability

The project, titled "Dengue Fever Spread Predicting Based on Ensemble Learning and Markov Chain Monte Carlo (MCMC) Method," has been developed under Python version 3.9.7. It is generously made available under a license that permits usage in any form, fostering an open environment for both academic and commercial endeavors. The official repository, which houses the project, is a treasure trove of resources for those interested in delving into the intricacies of the implemented algorithms and methodologies. It is hosted on GitHub and can be accessed via the following link: https://github.com/HironyHan/ELM_prediction. Users are encouraged to explore, utilize, and contribute to the repository, thus aiding in the continual enhancement and broadening the scope and efficacy of the project. This collaborative spirit underpins the ethos of open dissemination of knowledge and mutual advancement in tackling real-world issues such as Dengue Fever spread prediction through innovative computational frameworks.

# References

[1] S. M. Shakeel et al., "Covid-19 prediction models: a systematic literature review," Osong public health and research perspectives, vol. 12, no. 4, pp. 215–229, 2021.

[2] S. F. Ardabili, A. Mosavi, P. Ghamisi, F. Ferdinand, A. R. Varkonyi-Koczy, U. Reuter, T. Rabczuk, and P. M. Atkinson, "Covid-19 outbreak prediction with machine learning," Algorithms, vol. 13, no. 10, 2020. [Online]. Available: https://www.mdpi.com/1999-4893/13/10/249

[3] M. Siwiak et al., "From the index case to global spread: the global mobility based modelling of the covid-19 pandemic implies higher infection rate and lower detection ratio than current estimates," PeerJ, vol. 8, p. e9548, 2020.

[4] M. A. Achterberg et al., "Comparing the accuracy of several network-based covid-19 prediction algorithms," International journal of forecasting, vol. 38, no. 2, pp. 489–504, 2022.

[5] Z. Yang et al., "Modified seir and ai prediction of the epidemics trend of covid-19 in china under public health interventions," Journal of thoracic disease, vol. 12, no. 3, pp. 165–174, 2020.

[6] K. Chatterjee et al., "Healthcare impact of covid-19 epidemic in india: A stochastic mathematical model," Medical journal, Armed Forces India, vol. 76, no. 2, pp. 147–155, 2020.

[7] S. Bhandari, A. Tak, J. Gupta, B. Patel, J. Shukla, A. S. Shaktawat, S. Singhal, A. Saini, S. Kakkar, A. Dube, S. Dia, M. Dia, and T. Wehner, "Evolving trajectories of covid-19 curves in india: Prediction using autoregressive integrated moving average modeling," 2020. [Online]. Available: https://europepmc.org/article/PPR/PPR184651

[8] A. Tomar and N. Gupta, "Prediction for the spread of covid-19 in india and effectiveness of preventive measures," The Science of the total environment, vol. 728, p. 138762, 2020.

[9] P. Arora et al., "Prediction and analysis of covid-19 positive cases using deep learning models: A descriptive case study of india," Chaos, solitons, and fractals, vol. 139, p. 110017, 2020.

[10] R. Sujath et al., "A machine learning forecasting model for covid-19 pandemic in india," Stochastic environmental research and risk assessment : research journal, vol. 34, no. 7, pp. 959–972, 2020.

[11] L. J. Muhammad et al., "Supervised machine learning models for prediction of covid-19 infection using epidemiology dataset," SN computer science, vol. 2, no. 1, p. 11, 2021.

[12] Y. Li, Q. Dou, Y. Lu, H. Xiang, X. Yu, and S. Liu, "Effects of ambient temperature and precipitation on the risk of dengue fever: A systematic review and updated meta-analysis," Environmental Research, vol. 191, p. 110043, 2020. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0013935120309403

[13] T. Radivojević, Z. Costello, K. Workman, and H. Garcia Martin, "A machine learning automated recommendation tool for synthetic biology," Nature Communications, vol. 11, no. 1, p. 4879, 2020.

[14] DrivenData, "Dengai: Predicting disease spread," https://www.drivendata.org/competitions/44/dengai-predicting-disease-spread/, 2016, retrieved [10 31 2023].

[15] P. Bull, I. Slavitt, and G. Lipstein, "Harnessing the power of the crowd to increase capacity for data science in the social sector," CoRR, vol. abs/1606.07781, 2016. [Online]. Available: http://arxiv.org/abs/1606.07781