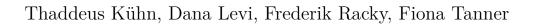
# Temperature sensitivity of dengue in Thailand



Bioinformatics Project group 4 team 4 climate sensitive infectious diseases Summer Term 2023

Supervisor:  $\mathbf{Dr.}$  Stella  $\mathbf{Dafka}$ 

#### The full abstract goes here

## Contents

1	Introduction	2
2	Material and Methods	3
	2.1 Material	3
	2.2 Methods	4
3	Results	6
	3.1 Result for Method A	6
4	Discussion	6

# Abbreviations

**DHF** Dengue Hemmorhagic Fever

**ARIMA** Autoregressive moving average

GAM General additive model

ADF Augumented Dickey Fuller test

## 1 Introduction

The Dengue virus is a vector borne virus, consisting of four serotypes (DENV 1-4). It is transmitted to humans by mosquitoes, more specifically by Aedes aegypti and Aedes albopictus (Phanitchat et al., 2019). The resulting infection goes asymptomatic in many cases, but can also cause Dengue Hemmorhagic Fever (DHF), characterized by symptoms like fever, headache, joint and muscle pain as well as (internal) bleeding and bruising (Gubler, 1998). Dengue is an emerging public health issue, with half of the worlds population at risk of infection. 390 million cases per year are estimated worldwide, with most of them not showing symptoms and thus not being reported. Affected areas range from sub-tropical to tropical regions, with south-east Asia, including Thailand, being one of the most seriously affected regions. Newly affected areas also include Europe (WHO, 2023). The disease control of dengue is challenging due to the absence of an effective treatment or vaccine, which leaves only the treatment of symptoms with painkillers like paracetamol (WHO, 2023). Responsible institutions in affected areas currently focus on

prevention, vector control, case control and prediction of possible future outbreaks (Phanitchat et al., 2019).

An important factor in predicting the epidemiological dynamics of Dengue is the climate: climate fluctuations due to recurring weather phenomenons and climate change are shown to influence *Aedes* biology and infections (Descloux *et al.*, 2012; Phanitchat *et al.*, 2019). In several studies, maximal temperature has an effect on dengue transmission (Descloux *et al.*, 2012), being associated with higher incidence (Phanitchat *et al.*, 2019) between 27 °C and 29,5 °C. Extreme global climate events like el Niño have been shown to affect disease outbreaks like Dengue as well (Anyamba *et al.*, 2019).

In Thailand, a significant climate factor is the monsoon, which can be separated into two seasons: The south-west monsoon between may and october is characterized by higher temperature and high precipitation. The north-east monsoon between november and february is characterized by lower temperature and low precipitation (Kripalani *et al.*, 1995). It has been shown that in northern districts of Thailand, there has been a detectable rise in temperature since the mid 20th century (Masud *et al.*, 2016).

In this analysis, we are going to investigate the correlation between temperature and Dengue in Thailand from 2006 to 2020 to assess the significance of climate change, recurring climate fluctuations and extremes and geographical factors on Dengue infections. It is concentrated on three main points: At first, time periods with DHF incidence will be compared to time points of extreme weather events. It will be examined if provinces with higher temperature also show higher incidences. The observation of trends in temperature and dengue cases over the course of the given time period will be analysed and compared. Based on the findings, two models will be generated to forecast the development of dengue fever cases: Autoregressive moving average (ARIMA) will be used to model the near future. Then, Generalized Additive Model (GAM) will be generated to predict the spatial distribution of incidences over Thailand in the future.

#### 2 Material and Methods

#### 2.1 Material

ERA5 data (climate) The ERA5 database is a global climate database by the European Centre for Medium-Range Weather Forecasts (ECMWF), covering the earth in a 31 km horizontal grid up to 80 km in the atmosphere in the time period from 1950 to present. It was generated from measurements of various climate variables combined with a reanalysis of existing data and past reanalysis data to accurately model and complete the dataset in the given resolution (Hersbach et al., 2020). For this project, monthly temperature data 2m above ground for every province in Thailand in the timeframe of 2006 - 2020 was extracted.

Dengue data The Dengue case numbers are retrieved from annual infectious disease reports published by the Thailand ministry of health. Monthly case numbers of DHF for every province in the timeframe of 2006 - 2020 are used in this project.

The datasets are used with a resolution at province level. As of 2011, Thailand has a total of 77 provinces, but had 76 provinces before 2011, as Bueng Kan was split from Nong Khai in 2011. For better compatibility of the data before and after 2011, the two new provinces are merged into a province equivalent to Nong Khai before 2011.

Climate forecast For GAM modeling, temperature data from the CORDEX climate model is used. It includes the temperature 2m above ground for June to August (south west monsoon) of the years 2021 -2040 at a 22 km grid (Copernicus Climate Change Service, 2019).

Population data We used ... population data of every Thailand province from 2006 - 2020 (Quelle?).

Spatial data To associate our data with the different provinces and visualize it, we use spatial data of Thailand's provinces. The two main data types in spatial data are vector-data and raster-data. Vector-data consists of a list of points with their exact location, which can then form lines or polygons. Raster-data assigns a value to every square of a raster. In this case, maps consisting of polygons for each province are used (Pebesma and Bivand, 2023). With the sf package, objects associating our data for each province with its coordinates and polygons are created, describing its shape. The function geom\_sf of ggplot2 are used for mapping.

#### 2.2 Methods

Descriptive analysis Linear regression is a method to describe a linear relationship between variables, using the minimum sum of squares between regression line and data points to identify possible trends (Schneider *et al.*, 2010).

ARIMA To predict the development of the dengue cases with an ARIMA model, a time series with the total dengue cases of Thailand was created. A time series can be decomposed in the components trend, seasonality and random. A requirement for fitting an ARIMA model is a stationary time-series. This is obtained, when the mean value doesn't change over time, the variance doesn't increase and the seasonality effect is minimal (Prabhakaran, 2017). Two tests were used to test for stationarity of the data. The Augmented Dickey-Fuller (ADF) test examines whether the time series has a unit root, indicating non-stationarity. In the ADF test,

the null hypothesis assumes the presence of a unit root, implying non-stationarity, while the alternative hypothesis suggests stationarity. The Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test is also a unit root test but focuses on the presence of a deterministic trend in the series. In contrast to ADF-test, for the KPSS-test the null hypothesis assumes stationarity. If the time series is initially found to be non-stationary, the differences between consecutive observations can be calculated, and the stationarity tests can be applied again (Hyndman and Athanasopoulos, 2018). ARIMA models combine an autoregressive model AR(p) and a moving average model MA(q). The autoregressive model computes the current value from previous values and the error term:

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_n y_{t-n} + \epsilon_t$$

 $\epsilon_t$  = white noise

 $1, \ldots, \phi_p = \text{parameters}$ 

 $y_{t-1}, \ldots, y_{t-p} = \text{lagged values}$ 

For the moving average the current value consists of the mean value of the time series and weighted current and past error terms:

$$y_t = c + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_t - 2 + \dots + \theta_q \epsilon_t - q$$

$$\theta_1, \ldots, \theta_q = \text{parameters}$$

I(t) is the number of times differencing was performed to make the time series stationary (Venkat, 2018). To find the optimal values for p and q, the Autocorrelation function (ACF) and partial Autocorrelation function (pACF) were evaluated. The ACF plot shows the correlations of a time-series with lags of itself. The pACF calculates the relationship between a time series and its lag, excluding the influence of linear dependencies among other lags (Prabhakaran, 2017). A second evaluation tool is the auto.arima function. The function automatically fits the best ARIMA model by minimizing the Akaike's Information Criterion (AIC), which is a criterion for the quality of the model. The auto.arima function can also consider seasonal models. Optimal models will yield uncorrelated residuals with zero mean and constant variance. This can be evaluated by plotting the ACF of the residuals and by performing a portmontreau test, for example the Ljung-Box test (Hyndman and Athanasopoulos, 2018).

GAM GAM provides insight into the shape and direction of the relationship between temperature and dengue cases. Additionally, it was used to forecast the future dengue case development based on the temperature development prediction. A GAM is a flexible extension of Generalized Linear Models (GLMs) that allows for the modelling of non-linear relationships between the response variable and predictor variables. A linear model can be described as follows:

 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p + \epsilon$  GAMs are now a nonparametric form of regression where the linear predictors  $(\beta_i x_i)$  of the regression are replaced by smooth functions of the explanatory variables,  $f(x_i)$ . The model can be defined as:  $y_i = f(x_i) + \in_i$  where  $y_i$  is the response variable,  $x_i$  is the predictor variable, and f is the smooth function (Wood, 2006). It is called an additive model, because all the  $f(x_i)$  functions and therefore their predictor variables contribute individually to the response variable and are added up. The advantage is that the different smooth functions can capture complex relationships by flexibly fitting curves to the data. There are two different ways to interpolate the functions of the predictor values. Finding a polynomial with a specific degree that passes through all the data points is a polynomial interpolation approach. Because high-degree polynomials may result into wide oscillation or overfitting, piece-wise interpolation is sometimes more accurate. Here, the data is divided into smaller intervals, each one is described by an individual function. Defining the number of knots determines into how many segments the model is divided. All polynomials together are called splines with a degree of k. They connect the knots one by one. The spline can be differentiated k-1 times. A smaller number of knots makes the response smoother, while a higher k value results in a curve that closely follows the individual data points. Choosing from various types of splines, in the GAM smoothing splines were used, which try to fit the data closely as well as maintaining smoothness (Peri, 2021). The obtained model includes the relationship between incidence of dengue cases and temperature. Such GAMs can be used to forecast the development of the response variable (incidence of dengue fever) based on the predictor variables (temperature). The predicted values can than be plotted on a map of Thailand. The prediction is based on the average temperature for the months June to August for the south west monsoon over the time period of 2021 until 2040 in Thailand.

## 3 Results

## 3.1 Result for Method A

Results go here as subsections.

## 4 Discussion

Your Discussion goes here!

## References

Anyamba, A., Chretien, J.-P., Britch, S. C., Soebiyanto, R. P., Small, J. L., Jepsen, R., Forshey, B. M., Sanchez, J. L., Smith, R. D., Harris, R., Tucker, C. J., Karesh, W. B., and Linthicum, K. J. (2019). Global Disease Outbreaks Associated with the 2015–2016 El Niño Event. Scientific Reports 9, 1930.

Copernicus Climate Change Service (2019). CORDEX regional climate model data on single levels.

Descloux, E., Mangeas, M., Menkes, C. E., Lengaigne, M., Leroy, A., Tehei, T., Guillaumot, L., Teurlai, M., Gourinat, A. C., Benzler, J., Pfannstiel, A., Grangeon, J. P., Degallier, N., and de Lamballerie, X. (2012). Climate-Based Models for Understanding and Forecasting Dengue Epidemics. PLOS Neglected Tropical Diseases 6, e1470, doi: 10.1371/JOURNAL. PNTD.0001470.

Gubler, D. J. (1998). Dengue and dengue hemorrhagic fever. Clinical Microbiology Reviews 11, 480–496.

Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., Chiara, G. D., Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L., Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., de Rosnay, P., Rozum, I., Vamborg, F., Villaume, S., and Thépaut, J. N. (2020). The ERA5 global reanalysis. Quarterly Journal of the Royal Meteorological Society 146, 1999–2049, doi: 10.1002/QJ.3803.

Hyndman, R., and Athanasopoulos, G. (2018). Forecasting: Principles and Practice (2nd ed). OTexts: Melbourne, Australia. https://otexts.com/fpp2/.

Kripalani, R. H., Singh, S. V., Panchawagh, N., and Brikshavana, M. (1995). Variability of the summer monsoon rainfall over Thailand—comparison with features over India. International Journal of Climatology 15, 657–672.

Masud, M. B., Soni, P., Shrestha, S., and Tripathi, N. K. (2016). Changes in climate extremes over North Thailand, 1960–2099. downloads.hindawi.com, 1960–2099.

Pebesma, E., and Bivand, R. (2023). Spatial Data Science: With Applications in R (Chapman and Hall/CRC).

Peri, S. P. (2021). GAMs and Smoothing Splines. Towards AI.

Phanitchat, T., Zhao, B., Haque, U., Pientong, C., Ekalaksananan, T., Aromseree, S., Thaewnongiew, K., Fustec, B., Bangs, M. J., Alexander, N., and Overgaard, H. J. (2019). Spatial and temporal patterns of dengue incidence in northeastern Thailand 2006–2016. BMC Infectious Diseases 19, 743.

Prabhakaran, S. (2017). Time Series Analysis With R. r-statistics.co .

Schneider, A., Hommel, G., and Blettner, M. (2010). Linear Regression Analysis: Part 14 of a Series on Evaluation of Scientific Publications. Deutsches Ärzteblatt International 107, 776.

Venkat, A. (2018). Time Series Analysis for Epidemiological Data.

WHO, W. H. O. (2023). Dengue and severe dengue, https://www.who.int/news-room/fact-sheets/detail/dengue-and-severe-dengue.

Wood, S. (2006). Generalized Additive Models: An Introduction With R, vol. 66.