
MODELING THE TEMPORAL POPULATION DISTRIBUTION OF *Ae. aegypti* MOSQUITO USING BIG EARTH OBSERVATION DATA

A PREPRINT

Oladimeji Mudele

Telecommunications and Remote Sensing Laboratory,
Department of Electrical, Computer
and Biomedical Engineering,
University of Pavia, Italy
oladimeji.mudele@universitadipavia.it

Fábio M. Bayer

Department of Statistics
and LACESM,
Federal University of Santa Maria, Brazil
bayer@ufsm.br

Lucas Zanandrez

Ecovec company, Belo Horizonte,
Brazil
lucas@ecovec.com

Alvaro E. Eiras

Laboratory of Technological Innovation
and Entrepreneurship in Vector Control
Department of Parasitology,
Federal University of Minas Gerais, Brazil
alvaro@icb.ufmg.br

Paolo Gamba

Telecommunications and Remote Sensing Laboratory,
Department of Electrical, Computer
and Biomedical Engineering,
University of Pavia, Italy
paolo.gamba@unipv.it

November 26, 2019

ABSTRACT

Over 50% of the world population is at risk of mosquito-borne diseases. Female *Ae. aegypti* mosquito species transmit Zika, Dengue, and Chikungunya. The spread of these diseases correlate positively with the vector population, and this population depends on biotic and abiotic environmental factors including temperature, vegetation condition, humidity and precipitation. To combat virus outbreaks, information about vector population is required. To this aim, Earth observation (EO) data provide fast, efficient and economically viable means to estimate environmental features of interest. In this work, we present a temporal distribution model for adult female *Ae. aegypti* mosquitoes based on the joint use of the Normalized Difference Vegetation Index, the Normalized Difference Water Index, the Land Surface Temperature (both at day and night time), along with the precipitation information, extracted from EO data. The model was applied separately to data obtained during three different vector control and field data collection condition regimes, and used to explain the differences in environmental variable contributions across these regimes. To this aim, a random forest (RF) regression technique and its nonlinear features importance ranking based on mean decrease impurity (MDI) were implemented. To prove the robustness of the proposed model, other machine learning techniques, including support vector regression, decision trees and k-nearest neighbor regression, as well as artificial neural networks, and statistical models such as the linear regression model and generalized linear model were also considered. Our results show that machine learning techniques perform better than linear statistical models for the task at hand, and RF performs best. By ranking the importance of all features based on MDI in RF and selecting the subset comprising the most

informative ones, a more parsimonious but equally effective and explainable model can be obtained. Moreover, the results can be empirically interpreted for use in vector control activities.

Keywords *Ae. aegypti* · Machine learning · Random forest · Remote sensing

1 Introduction

Studying urban ecosystems is a hot topic within various scientific clusters [1–5]. The compounding effects of human activities through urbanization, carbon emission and other biodiversity changes are modifying urban ecosystems, thus creating a need for continuous assessment of the environment to ensure health and quality of life suitability for dwellers [6]. With a repository covering over 40 years of data collected, Earth observation (EO) data from orbiting satellites present a gold mine for environmental change monitoring [7]. Due to recent advances resulting in higher resolution sensors (spectral, spatial and temporal), freely accessible data, and efficient processing algorithms, it is possible to extract static and dynamic phenomena happening on the Earth surface and use this information for various environmental characterization procedures [8–11].

One major area currently experiencing rapid innovation due to EO data is landscape epidemiology [12]. In this field, environmental variables that are proxies to favorable conditions for the spread of disease causing vectors are extracted from EO data and used to model the distribution of these vectors [13, 14]. Major disease carrier vectors with spread and population depending on environmental conditions include mosquito species, ticks and rodents [15–18]. Mapping and predicting vector and diseases spread through EO data can provide reliable information in outbreak early warning, hotspots location, and public health budget optimization [19, 20]. To this aim, there have been seminal works considering how to include major environmental variables obtainable from EO data such as surface temperature, vegetation condition, air quality, precipitation rate, humidity conditions, night-time lights, land cover types, degree of urbanization and velocity of the urban sprawl, among others [21–24].

According to the data provided by the World Health Organization, over 50% of the world population is exposed to the risk of mosquito-borne diseases [18, 25], the ones with the greatest clinical importance being Zika, Chikungunya, Dengue and Yellow fever virus diseases. These diseases are transmitted mainly by the female *Ae. aegypti* mosquito species [16, 19]. While its male counterpart feeds only on fruit nectar, the female *Ae. aegypti* bites humans to feed on blood [26] and poses a serious threat to humans by transmitting dangerous viruses. Areas of the world majorly exposed to risks of arbovirolosis carried by this mosquito species include Latin America, Central Africa, and South-East Asia, with major disease spread outbreaks already recorded in Brazil, Argentina, Colombia, and Venezuela. The *Ae. aegypti* species is known to be adaptable to urban environments due to its inclination to being bred in artificial containers [19, 21, 27, 28]. Spread of the causal diseases by this vector is known to correlate with the local adult vector population [29]. Environmental conditions including precipitation, vegetation conditions, temperature, and humidity have been shown in previous works to significantly influence *Ae. aegypti* mosquito development [21, 27, 30]. After obtaining these environmental variables from satellite data, statistical and machine learning (ML) models of diseases outbreak can be used for vector population prediction [31, 32].

ML is a subdivision of artificial intelligence which deals with the implementation of algorithms to learn complex patterns from machine readable input data for classification and regression purposes [33, 34]. It is a collection of methods that provide multivariate, nonlinear and non-parametric regression or classification of data. ML has proven useful for a very large number of EO system applications such as crop disease detection [35], urban area classification [23], landscape epidemiology [19, 21, 27], and various other bio-geographical information extraction and regression analysis procedures [33, 34].

In landscape epidemiology, ML algorithms can be used to map relationships between proxy environmental conditions and ground collected vector population information for the purpose of prediction, classification and/or explanation of environmental condition effects on vector spread [21, 27]. Previous attempts in modeling the population of female *Ae. aegypti* and other species of mosquitoes have employed ML algorithms including maximum entropy methods, support vector regression (SVR), decision trees regression (DTR), k-nearest neighbor regression (KNN), artificial neural networks (ANN), multilayer perceptron (MLP), among others [19–21, 27].

In spite of the works in this study domain, there are still some gaps. One of such is in the selecting the most informative environmental features subset to obtain the model with best fit and least redundancy. In [27], KNN, SVR, and MLP were used along with classical statistical models in time series modeling of *Ae. aegypti* oviposition activity. To select the best model, feature significance estimated with linear correlation coefficient was used. This metric, however, only considers pairwise linear relationship among candidate features as selection criterion, ignoring nonlinear associations. The same limitation also applies to [28].

In this paper, we present a methodology for modeling the temporal population distribution of female *Ae. aegypti* mosquito based on time series environmental variables obtained from freely accessible EO data products and field collected mosquito population data. To achieve our aim, a random forest (RF) [36, 37] model is considered due to: (i) its robustness against unbalanced data with good performance in complex problems and (ii) its embedded mean decrease impurity (MDI) variable importance measure. MDI is used to rank and extract the most relevant environmental feature subset for modeling the vector population. It considers also nonlinear associations between candidate features. For benchmark purposes, RF is compared with other ML models already used in similar works namely: KNN, SVR, DTR, and MLP, and statistical regression models such as linear regression model (LM) and generalized linear model (GLM) [38]. Finally, we considered the best model, fitted with the most informative features as ranked by MDI, and demonstrated using our data how the effects of the selected features on the vector population can be explained and made into actionable insights for field investigators.

2 Study Area and EO variables

Even though climatic and vegetation condition effects alone do not determine the population of *Ae. aegypti* mosquitoes, entomological and epidemiological studies have described the effects of temperature, humidity, precipitation, and vegetation condition on the population and virus transmission factors of *Ae. aegypti* [39]. In this study, we model the population of mosquitoes as a function of selected remotely sensed environmental variables.

To this aim, we consider spatiotemporal data for the municipality of Vila Velha, located between latitudes $20^{\circ} 19'$ and $20^{\circ} 32'$ south and longitudes $40^{\circ} 16'$ and $40^{\circ} 28'$ west, on the coast of the Espírito Santo State of Brazil. This municipality covers a total area of $209,965 \text{ km}^2$, with an estimated population of 479,664 inhabitants [40]. Figure 1 shows the location of the studied area.

The control of *Ae. aegypti* borne diseases is carried out in Brazil based on the National Dengue Control Program (PNCD) which has been in operation since 2002 [41]. This program addresses important components including: epidemiological surveillance, vector control, environmental sanitation, among others [42]. Plugging into the epidemiological surveillance framework, Ecovec — an expert real-time mosquito surveillance company based in Belo Horizonte, Brazil [43] — began the implementation of the MI-Aedes[®] program in Vila Velha in 2017 to improve the operation of PCND. The MI-Aedes[®] program uses 791 geopositioned adult mosquito traps to obtain weekly *Ae. aegypti* population data across the municipality of Vila Velha. A single mosquito trap is placed at every point, with approximately 250 meters separation from the nearest surrounding traps to limit the effect of spatial autocorrelation. All the utilized traps are of the same type, specifically, the sticky-card oviposition trap, commercially named MosquiTrap[®] [44]. This type of trap is made of a black bucket containing water and a sticky-card to capture any mosquitoes that try to lay eggs on the trap's internal walls. It must be inspected weekly to change the water and collect mosquito population data. Also, a commercial mosquito lure (AtrAedes[®]) was used to enhance attraction of females of *Aedes aegypti* to the trap. An aerial view of the used MosquiTraps[®] location is shown in Figure 1.

Since the commencement of MI-Aedes[®], vector control actions in Vila Velha has achieved improvements at intervals by using the mosquito population data collected. For example, data collected in 2017 was used to improve control actions from the 1st week of 2018 by directing more control actions to places with higher adult mosquito population. One major advantage of the MI-Aedes[®] program is that it collects weekly data of adult vector population, not immature forms. This ensures that the data highly correlates with human diseases infection cases since only adult female mosquitoes are responsible for the transmission [45].

For the purpose of this work, we used weekly female *Ae. aegypti* MI-Aedes[®] program data from week 15 in 2017 to week 34 in 2019 — a total of 124 weeks. Data collection was carried out during weekly inspections conducted by municipal field workers who are trained by Ecovec for trap inspections, mosquito species identification, and to use an online mobile system to send the data. The inspections occurred on the same day of the week, and same time of the day for each MosquiTrap[®]. The collected and sent data were weekly checked and approved by a trained supervisor and a specialist from Ecovec. To maintain operational standards on the field during the collection tenure, an Ecovec specialist visited the field every six months to check the conditions of the traps and to update staff training.

Our analyses account for changes in data collection conditions such as the earlier reported data-driven improvement of control actions starting from the 1st week of 2018 and an update in the MI-Aedes[®] platform starting from week 42 in 2018 which led to changes in identification number of the mosquito traps and slight changes in their geopositions. In the cases of trap location changes, the traps were moved to new adjacent residences due to monitoring difficulties or because residents specifically requested the changes. All in all, trap relocations beyond a 40 meters radius were rare.

As a result of these discrepancies in data collection conditions, our data and consequent analyses were split in three time batches, presented in Table 1. Since previous works [27, 28] have shown that it is sufficient to work with weekly

Table 1: Mosquito population data Batches used in this study

Batch	Date range	Total weeks	Differentiating condition
1	10/04/2017 - 31/12/2017	38	-
2	02/01/2018 - 05/10/2018	40	Vector control improvement
3	08/10/2018 - 23/08/2019	46	Update to the MI-Aedes [®] platform

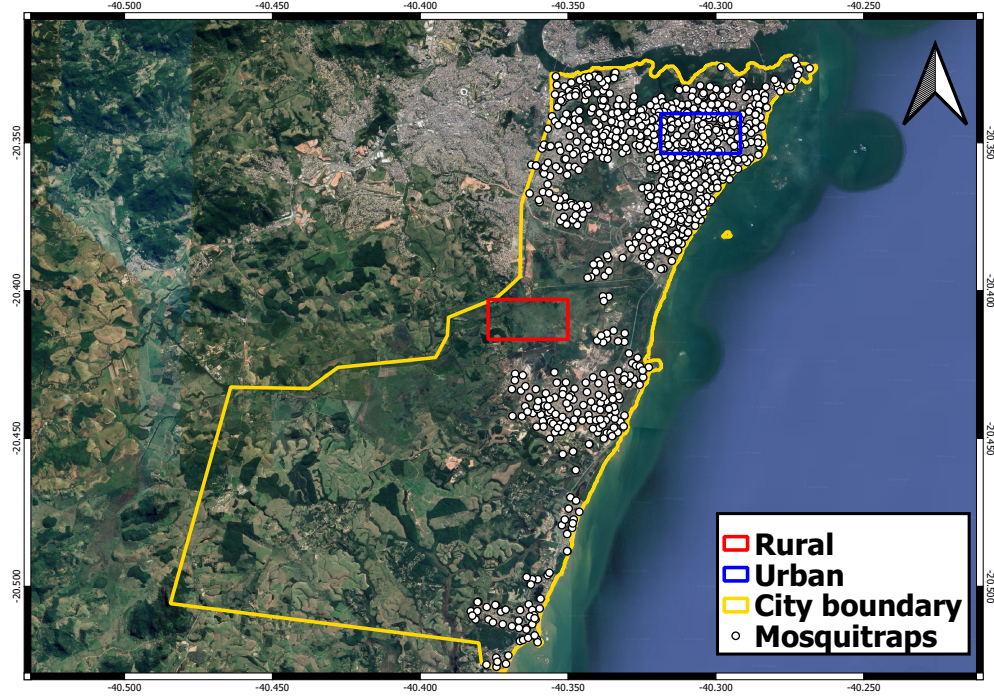


Figure 1: Aerial view of Vila Velha municipality, data collection mosquito trap locations, and urban and rural surface zones selected to extract the environmental variables.

observation samples, we cleaned the *Ae. aegypti* population data to obtain female mosquito population per mosquito trap on a weekly basis. It is important to note two points here: firstly, we expect lower vector population and correlation with environment condition starting from the second Batch in our analyses due to improvements in control actions. Secondly, the three data Batches do not exactly span the same time and seasons of the year.

In studies that associate the development of *Ae. aegypti* with EO derived environmental factors, the Moderate Resolution Imaging Spectroradiometer (MODIS) satellite data products provided by the National Aeronautics and Space Administration (NASA) have been recently used [27, 28]. This is due to its free access, global availability and non-requirement of much further processing to obtain the needed environmental features [46].

For this work, data products providing information corresponding to surface temperature, precipitation, humidity, and vegetation conditions were collected for the weeks matching the vector population data. For vegetation condition information, the Normalized Difference Vegetation Index (NDVI) has been obtained from the MODIS MOD13Q1 product – a 16-day composite product with spatial resolution of 250 meters. The near and mid-infrared bands of the same data product was used to compute the Normalized Difference Water Index (NDWI) using Gao’s definition [47]. This layer provides information about vegetation water content which is a proxy to humidity and surface moisture [27]. In addition, estimates of the minimum and maximum surface temperatures have been obtained from the day and night-time land surface temperature (LST) bands of the 8-day composite MODIS MOD11A2 product at 1 km spatial resolution [48]. Finally, the calibrated precipitation band of the Global Precipitation Measurement (GPM) mission data

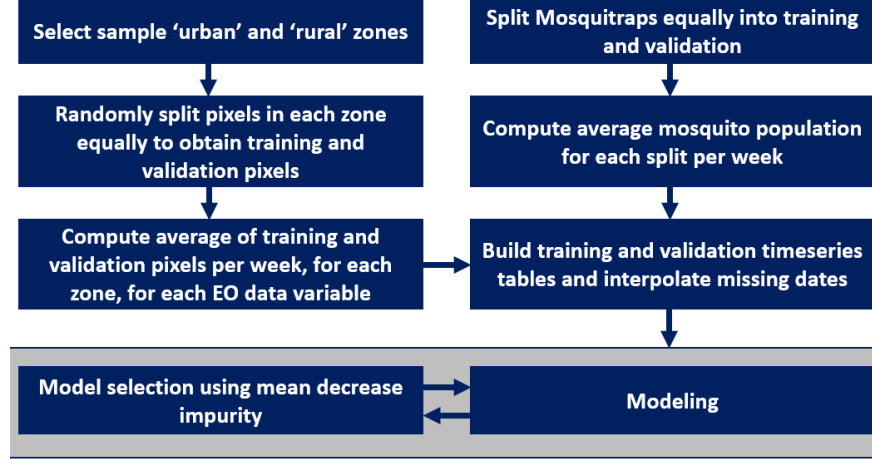


Figure 2: Schematic of the study methodology.

product was used to obtain precipitation information [49], although in a very coarse way, because this data is available daily at 0.1 deg spatial resolution (≈ 11 km). All data were downloaded using the `export` object of the JavaScript application programming interface of google earth engine (GEE). This object has resampling and reprojection methods abstracted into it for easy co-registration of data with different properties. For each download task, we obtained the output data at a resampled spatial resolution of 250 meters because this is the minimum distance separating neighboring mosquito traps. We used the `scale` parameter of the `export` object to set this spatial resolution value. We set the common coordinate reference system projection for all downloaded data as EPSG:4326 using the `crs` parameter.

As previously implemented in [27, 28] and [50], to find the relationship between the measured number of mosquitoes and the model output exploiting EO data, we defined two buffer areas of size 18 km^2 each from which distinct temporal characterization of the considered EO variables are obtained within our study area. One of the zones is in the densely urbanized part of the city (labeled as "urban"), and the other in the peri-urban zone (labeled as "rural") areas. These two areas are shown in Figure 1.

As in typical supervised ML applications, pixels within each selected zone were randomly and equally split into the training and validation sets. The average for training and validation pixels were computed for each image in the stack of the EO data set representing individual environmental variables considered. To obtain features for the mosquito data, we aggregated the training and validation data for each week from all the mosquito traps in the whole study area, not only from the buffer zones as in the case of the EO covariates. We randomly split all the mosquito traps equally into training and validation. We compute the means for each split per week to obtain the response time-series training and validation data. We used a fifth order spline interpolation to obtain values for missing dates. All environmental variables are considered with up to 2 weeks lag to account for non-synchronous effects resulting from vector development life cycle [51]. As in [27], all the variables were standardized (z-score). This was done because variable rescaling is a good practice, especially for training a neural network [52]. A high-level schematic of the whole methodology of our study as described in this section is presented in Figure 2.

3 Modeling Procedure

To model the weekly mean number of mosquitoes per mosquito trap (Y), we consider fitting a RF model. For prediction benchmark purpose, ANN, SVR, KNN, DTR, LM, and GLM fitting models were also considered. Finally, a more parsimonious RF model, labeled RF*, considering only the most informative climate variables obtained by using the MDI to rank all predictor covariates was also implemented and compared with the other ones. For prediction performance comparison, the usual quantitative measures between the observed data (y) and the predicted values (\hat{y}) were considered, namely: the linear correlation coefficient (R), the root mean square error (RMSE), and the mean absolute percentage error (MAPE) [53]. We use R and MAPE to measure relative qualities, and RMSE to measure absolute fit of our models. Background details and our implementation (function and model parameters) of RF are presented in Section 3.1. We present the same for the other models in Section 3.2. We implemented all our modeling computation with R programming [54].

3.1 Random Forest Regression

Random forests (RF) [36] are one of the most effective computationally intensive procedures to improve on unstable estimates, especially when it is difficult to find a good model due to problem complexity [37]. It is a predictor consisting of a collection of several randomized regression trees [37]. Let $\mathbf{X} \in \chi \subset \mathbb{R}^p$ an input vector related to p features and $Y \in \mathbb{R}$ the response random variable, the objective is to estimate the regression function $m(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$. Given a training sample $\mathcal{D}_n = ((\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n))$ of independent random variables distributed as the independent prototype pair (\mathbf{X}, Y) , the goal is to use the data set \mathcal{D}_n to construct an estimate $m_n : \chi \rightarrow \mathbb{R}$ of the function m . To this aim, the random forest consists of a collection of M randomized regression trees. For the j th tree in the family, with $j = 1, \dots, M$, the predicted value at each j is denoted by $m_n(\mathbf{x}; \Omega_j, \mathcal{D}_n)$, where $\Omega_1, \dots, \Omega_M$ are random variables independent of \mathcal{D}_n . The variables Ω_j are used to resample the training set prior to the growing of individual trees and to select the directions for splitting. Then, the trees are combined to form the forest estimate given by:

$$\hat{y} = m_{M,n}(\mathbf{x}; \Omega_1, \dots, \Omega_M, \mathcal{D}_n) = \frac{1}{M} \sum_{j=1}^M m_n(\mathbf{x}; \Omega_j, \mathcal{D}_n). \quad (1)$$

In this procedure, we have also to set other tuning parameters, i.e., the number of possible directions for splitting (m_s) at each node of each tree, and the lowest number of examples (m_e) in each cell to cause a split. In this work we set initially $M = 500$, $m_s = \lceil p/3 \rceil$ and $m_e = 5$.

One important characteristic of the random forest is that it can be used to rank the importance of the input variables (i.e., the features we extract from EO data) over the pattern variability of the target variable Y . In this work, we considered the MDI calculated based on the reduction in sum of prediction squared errors averaged over all trees whenever a variable is chosen to split [55].

3.2 Other Predictive Models

As mentioned above, for comparison purposes a number of other regression models have been implemented and their results compared with the ones by the proposed RF procedure. Specifically, they are:

- the Linear Regression Model (LM): the most used predictive model in different fields. With respect to the topic of this paper, the LM has been used in similar works, [28] and [27], to model the *Ae. aegypti* oviposition activity in a northern Argentine city. In this work, LM is implemented using the `lm()` function in R, and an ordinary least square estimator is considered to estimate the regression parameters.
- The Generalized Linear Model (GLM): a class of regression models able to model response variables in the exponential family of distributions [38]. Since the weekly mean value of mosquitoes population (y) is continuous and always in the space of positive real numbers ($y \in \mathbb{R}^+$), the Gaussian assumption considered for inference in the LM can lead to poor results and predictions. In GLM, instead, the gamma distribution is considered to model y , together with the log link function in the GLM regression structure.
- Artificial Neural Networks (ANN): one of the most used ML methods for geoscience problems. More recently, deep learning methods – corresponding to ANN architectures with several hidden layers [34] – became widely applied in many EO data processing problems [33, 56]. The ANN is a nonlinear data modeling technique used to model complex relationships between sets of input and output variables [57]. In this work, a multilayer perceptron with three layers, each with three neurons, is used. Each of the neuron is activated with a logistic activation function, and the resilient backpropagation with weight backtracking algorithm [58] was considered for training the neural network.
- Support Vector Regression (SVR): an approach based on support vectors with radial basis function kernel with tuning parameter γ set as $1/(\text{number of features})$, tuned via the the epsilon-regression method.
- k-Nearest Neighbor Regression (KNN): a regression based on the 4-nearest neighbors according to Euclidean distance and uniform weights for local interpolation.
- Decision Trees Regression (DTR): a regression model in the form of a tree structure [59]. It is a simple but powerful prediction method [60].

4 Results and Discussion

All the EO data variables used in this work were considered both in an urban surface zone (U) and in a peri-urban, or rural zone (R). We thus have the following variables: NDVI-U, NDWI-U, TempD-U, TempN-U, and Prec-U, as well as

Table 2: Descriptive measures of the EO-based environmental variables of interest, computed separately for each Batch

Variable	Mean	Median	SD	Min	Max
Batch 1 (2017)					
y	0.1910	0.1587	0.1446	0.0319	0.7873
NDVI-U	0.2269	0.2260	0.0261	0.1702	0.2841
NDVI-R	0.7280	0.7395	0.0448	0.5826	0.7983
NDWI-U	0.0024	0.0035	0.0350	-0.0994	0.0748
NDWI-R	0.5407	0.5492	0.0473	0.3952	0.6640
TempD-U	31.4166	31.0544	4.1390	22.5969	39.3821
TempD-R	27.7208	27.7761	2.9338	20.5000	33.8600
TempN-U	20.0180	20.6767	2.4379	13.4942	23.5909
TempN-R	19.0909	19.3264	1.8757	13.6278	23.1461
Prec-U	2.9118	0.2000	7.5662	0.0000	43.2500
Prec-R	2.9039	0.5375	6.7846	0.0000	37.1750
Batch 2 (2018)					
y	0.1267	0.0923	0.0956	0.0531	0.5057
NDVI-U	0.2500	0.2474	0.0179	0.2086	0.2996
NDVI-R	0.7475	0.7554	0.0261	0.6802	0.7922
NDWI-U	0.0030	-0.0021	0.0365	-0.0591	0.1065
NDWI-R	0.5923	0.5965	0.0518	0.4762	0.6925
TempD-U	31.8615	31.5774	3.4781	26.4761	38.8394
TempD-R	27.5138	26.9434	2.6868	23.7850	33.8161
TempN-U	21.2936	21.0190	2.2088	16.3886	25.3319
TempN-R	20.1959	19.6536	2.0855	16.5067	23.7806
Prec-U	2.4707	0.0654	4.9280	0.0000	23.0803
Prec-R	1.7527	0.0990	3.4678	0.0000	16.0950
Batch 3 (2018-2019)					
y	0.1134	0.1075	0.0315	0.0620	0.1922
NDVI-U	0.2557	0.2551	0.0188	0.2117	0.3075
NDVI-R	0.7210	0.7265	0.0331	0.5843	0.7829
NDWI-U	-0.0137	-0.0200	0.0324	-0.0762	0.0661
NDWI-R	0.5350	0.5450	0.0538	0.4031	0.6272
TempD-U	33.6604	33.8508	4.3575	18.6561	39.0947
TempD-R	28.9434	29.4046	2.3435	24.7489	33.1411
TempN-U	22.0600	22.4421	2.1145	17.3431	25.2092
TempN-R	20.6418	21.1400	2.6758	13.2544	24.1517
Prec-U	1.6132	0.1596	2.8811	0.0000	11.3604
Prec-R	1.3545	0.0419	2.4574	0.0000	9.5575

NDVI-R, NDWI-R, TempD-R, TempN-R, and Prec-R, where TempD is the daytime temperature and TempN is the night-time temperature. Table 2 reports the mean, median, standard deviation (SD), minimum (Min) and maximum (Max) of the EO data variables together with the female mosquito population (y). These measures were computed for the three Batches of time considered in our study (cf. Table 1) in a separate way, to evaluate the influence of the vector control program improvement implemented in 2018 and 2019, affecting Batches 2 and 3. We observed that the average value of the number of female mosquitoes in Batch 1 is 0.1910, while the values in Batches 2 and 3 are 0.1267 and 0.1134, respectively. This represents about 40% decrease in the mean number of mosquitoes from Batch 1 compared to subsequent Batches considered and highlights the efficacy of the data-driven vector control program improvement which has been implemented. Furthermore, the nonparametric Kruskal-Wallis test [61] rejects the null hypothesis of three populations being equally distributed (p -value = 0.0092), confirming that the decrease in the mean number of mosquitoes from Batches 1 to 2 and 3 is statistically significant. In addition to the Kruskal-Wallis test, the Nemenyi test shows that the population from Batch 1 differs from Batches 2 and 3, but population distributions of mosquitoes of Batches 2 and 3 do not differ statistically. Also, we can note that the standard deviation of the number of mosquitoes decreases from 0.1446 in Batch 1 to 0.0956 and 0.0315 in Batches 2 and 3, respectively. The boxplot in Figure 3 shows the differences among the population of female mosquitoes in the different Batches. Particularly, we can see that the maximum value in Batch 3 is four times less than the maximum number in Batch 1.

The environmental characteristics, as presented in Table 2, slightly differ across the three considered batch periods due to differences in observation time and seasons. Also, it can be seen that there was more precipitation in 2017 (Batch 1) than 2018 and 2019 (Batches 2 and 3). Regarding the urban and rural zones, the average land surface temperature in urban zone, TempD-U and TempN-U, respectively, is, as expected, higher and more variable than in the rural zone. Moreover, as equally expected, the opposite is the case for both NDVI and NDWI.

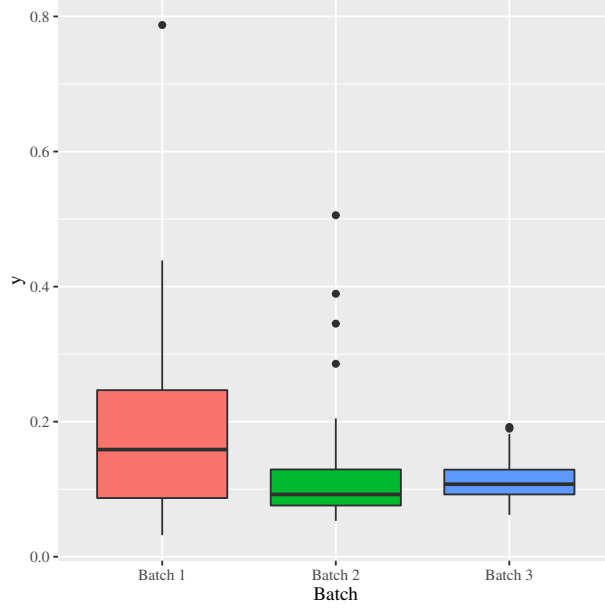


Figure 3: Boxplot comparing the distribution of the number of female mosquitoes in the three considered periods.

The linear correlation between each environmental variable, including lagged ones, and the female mosquito population, y , is presented in Table 3. This table shows, first of all, that none of the considered covariates show very strong linear correlation (above 0.7) with y . However, in all Batches, the non-lagged daytime temperature variables show the highest correlation with y : TempD-R and TempD-U — in that order — in both Batches 1 and 2, and only the former in Batch 3. In general, more EO variables show higher correlation with the vector population in Batch 2, with five non-lagged variables showing correlation above 0.4, compared to the other Batches.

In addition, we can note that the lagged NDWI-R and NDVI-R and the precipitation variables show greater correlation with y than with non-lagged ones. This observation can be explained by the aquatic cycle of *Ae. aegypti* from egg to adult being approximately 7-9 days [51], thus making the transition dependent on non-synchronous environmental effects, as also highlighted in [62]. We observe that the linear relationships between the targets and the predictors could be insufficient to describe the temporal variance in the mosquito population, justifying the need for non-linear approaches such as the considered ML methods.

Considering the results for the validation dataset, the quantitative measures of the prediction obtained from each considered model are shown in Table 4 for all Batches. In Batch 1, the measures show that GLM produces the worst results in terms of R and RMSE, while LM the worst in terms of MAPE. The best performances are reached in R by RF, in RMSE by RF*, and in MAPE by ANN. The outstanding performance of RF, and its more parsimonious variant, RF*, is highlighted by the correlation measure of 0.90 and 0.86, and RMSE of 0.0369 and 0.0396, respectively. Also, in Batch 2, we obtain the best result in terms of RMSE with RF*, followed by the fully fitted RF. In this Batch, however, best results considering R and MAPE are obtained using KNN. In Batch 3, RF* produced the best model as measured by R, RMSE and MAPE. Considering RMSE, RF* produces equal or better quality in all Batches. In general, from the whole results, though the statistical models in certain cases produced sufficient or comparable measured qualities, there are much more cases of better performances by machine learning models. For example, LM produces comparable quality with respect to RF in all Batches of our analyses if we look only at the explained variance which is measured by R. When we consider the RMSE, however, we see that that RF performs better than LM in terms of absolute fit. LM's high correlation measure of 0.8062 and 0.8535 in Batches 1 and 2 respectively are as a result of better correlation between more of their EO covariates and y compared to what is obtained in Batch 3. The much lower R measure for the all models on Batch 3 data could be as a result of time cumulative effect of data driven improvement in the vector control actions which led to less variation in the mosquito population due to environmental changes. This same cause may be used to explain the lower correlation among more EO covariates and y in the Batch 3 with respect to Batches 1 and 2, as earlier presented in Table 3.

Figure 4 shows a scatterplot of the actual and predicted female mosquito population values. In this plot, GLM, LM, and DTR are not considered due to their low RMSE figures. In all the Batches, the EO-based environmental features show significant effects on the vector population, less so in Batch 3. The environmental effects in Batches 2 and 3

Table 3: Correlation matrix between the average number of mosquitoes (y) and each covariate among the pool of EO-based variables in Table 2

Covariate	Lag 0	Lag 1	Lag 2
Batch 1			
NDVI-U	-0.3097	-0.0926	-0.1127
NDVI-R	0.0030	0.0958	-0.0142
NDWI-U	-0.3989	-0.2592	-0.1480
NDWI-R	0.0624	0.1566	0.1342
TempD-U	0.4046	0.3086	0.2236
TempD-R	0.4157	0.3327	0.1532
TempN-U	-0.2959	0.0895	0.1403
TempN-R	0.1105	0.1701	0.0328
Prec-U	0.0084	-0.0349	-0.0553
Prec-R	0.0261	-0.0368	-0.0335
Batch 2			
NDVI-U	-0.1401	-0.2117	-0.3098
NDVI-R	-0.5025	-0.5852	-0.4079
NDWI-U	-0.2454	-0.2772	-0.2118
NDWI-R	-0.2702	-0.3954	-0.3392
TempD-U	0.5810	0.4700	0.3565
TempD-R	0.6418	0.4962	0.4988
TempN-U	0.4835	0.4150	0.3578
TempN-R	0.4323	0.3584	0.4254
Prec-U	-0.0185	-0.1644	0.0360
Prec-R	0.0053	-0.1595	-0.0232
Batch 3			
NDVI-U	-0.1527	-0.0546	0.0913
NDVI-R	0.0439	0.0807	0.3296
NDWI-U	0.1716	0.0320	-0.0873
NDWI-R	0.0380	0.1005	0.1480
TempD-U	-0.1082	-0.0414	-0.0750
TempD-R	-0.2703	-0.1654	-0.2530
TempN-U	-0.1004	-0.0566	-0.0871
TempN-R	-0.1352	-0.1698	-0.0909
Prec-U	-0.0969	-0.1376	0.0578
Prec-R	-0.0384	-0.1142	0.0589

show pointers with which the vector control program can be further improved. A plausible approach to consider is weighting the intensity of the program implementation by the most relevant environmental conditions across the year. This can ensure a higher control intensity during periods of favorable climate for vector development, thus reducing risks exposure during these periods and helping to optimize vector control resources allocation. As opposed to current common practice which favors the application of control actions mostly in rainy season, these actions should also be performed in dry season and periods of low vector population. Implementing control actions during the dry season has been shown in [63] to suppress the population of eggs laid during this season and consequently reduce the hatched population in the following rainy season.

In Table 5, we present a summary of the observed and fitted data: mean, median, minimum (Min), maximum (Max), first (Q1) and third (Q3) quartiles. In Batch 1, LM exaggerates the minimum value, producing a negative minimum value which has no physical meaning in relation to mosquito population. Considering that this same model produced a good quality measure of R, we deduce that LM, regardless of the quality score, is not a good model for predicting mosquito population because its distribution assumption of the response variable is inappropriate for this study use case. This problem, however, does not occur with GLM, due to its gamma distribution assumption which assumes positive values for y . GLM produces a good estimated minimum value in all Batches, but still exhibits low prediction performance, which can be seen in its predicted mean and Q3 which are above the observed values for all Batches. Among the ML methods, ANN produces better prediction in relation to minimum and maximum values. RF and RF*, though both overestimate the minimum in all Batches produced good estimated mean values and showed good balance. The results in Figure 4, Table 4 and Table 5 show that the best models for predicting the vector population are RF* and RF, followed by SVR.

Figure 5 shows the mosquito population data (validation set) along with the predicted values by the best models. It is seen that the SVR and KNN models underestimate the seasonal spikes in Batches 1 and 2. The model performances are

Table 4: Quality measures of predictions in the validation dataset

Model	R	RMSE	MAPE
Batch 1			
GLM	0.6274	0.0722	60.6487
LM	0.8062	0.0480	75.0815
ANN	0.8420	0.0447	34.1122
SVR	0.8270	0.0471	36.1444
KNN	0.7456	0.0478	64.1194
DTR	0.6334	0.0547	67.9013
RF	0.9066	0.0369	56.9834
RF* (8 features)	0.8618	0.0396	61.8109
Batch 2			
GLM	0.6366	0.0490	43.0603
LM	0.8535	0.0309	44.3038
ANN	0.7202	0.0360	44.3968
SVR	0.8651	0.0252	30.9242
KNN	0.8839	0.0223	30.7137
DTR	0.7925	0.0281	32.5429
RF	0.8530	0.0255	35.4397
RF* (8 features)	0.8594	0.0247	33.4977
Batch 3			
GLM	0.5707	0.0180	23.8297
LM	0.5697	0.0168	23.3971
ANN	0.4048	0.0192	24.3962
SVR	0.4096	0.0174	23.8113
KNN	0.3662	0.0180	27.1013
DTR	0.3466	0.0179	25.7221
RF	0.5981	0.0156	21.4197
RF* (8 features)	0.5885	0.0156	21.5879

Table 5: Summary of the observed and fitted data

	Min	Q1	Median	Mean	Q3	Max
Batch 1						
y	0.0093	0.0474	0.0803	0.0933	0.1204	0.3848
GLM	0.0063	0.0397	0.0745	0.1043	0.1611	0.4340
LM	-0.0724	0.0413	0.0793	0.0946	0.1330	0.3565
ANN	0.0157	0.0546	0.0752	0.0995	0.1312	0.4309
SVR	0.0257	0.0610	0.0822	0.0838	0.1090	0.1496
KNN	0.0325	0.0621	0.0928	0.0984	0.1357	0.1967
DTR	0.0392	0.0392	0.0741	0.0986	0.1411	0.1715
RF	0.0362	0.0703	0.0918	0.0985	0.1224	0.2512
RF*	0.0461	0.0646	0.0858	0.0956	0.1187	0.2510
Batch 2						
y	0.0228	0.0351	0.0461	0.0619	0.0638	0.2402
GLM	0.0208	0.0384	0.0695	0.0890	0.0976	0.5542
LM	0.0140	0.0376	0.0716	0.0795	0.0977	0.3030
ANN	0.0247	0.0370	0.0488	0.0694	0.0857	0.2653
SVR	0.0309	0.0456	0.0498	0.0614	0.0703	0.1351
KNN	0.0310	0.0450	0.0487	0.0595	0.0566	0.1688
DTR	0.0438	0.0438	0.0438	0.0618	0.0641	0.1439
RF	0.0370	0.0454	0.0500	0.0642	0.0701	0.1769
RF*	0.0367	0.0436	0.0507	0.0639	0.0678	0.1856
Batch 3						
y	0.0303	0.0417	0.0531	0.0563	0.0657	0.1176
GLM	0.0303	0.0440	0.0576	0.0606	0.0702	0.1020
LM	0.0251	0.0453	0.0598	0.0596	0.0693	0.0888
ANN	0.0380	0.0419	0.0599	0.0574	0.0659	0.0893
SVR	0.0409	0.0515	0.0568	0.0566	0.0626	0.0700
KNN	0.0389	0.0544	0.0591	0.0592	0.0657	0.0777
DTR	0.0478	0.0478	0.0584	0.0576	0.0673	0.0673
RF	0.0405	0.0504	0.0558	0.0565	0.0626	0.0743
RF*	0.0405	0.0491	0.0547	0.0562	0.0636	0.0760

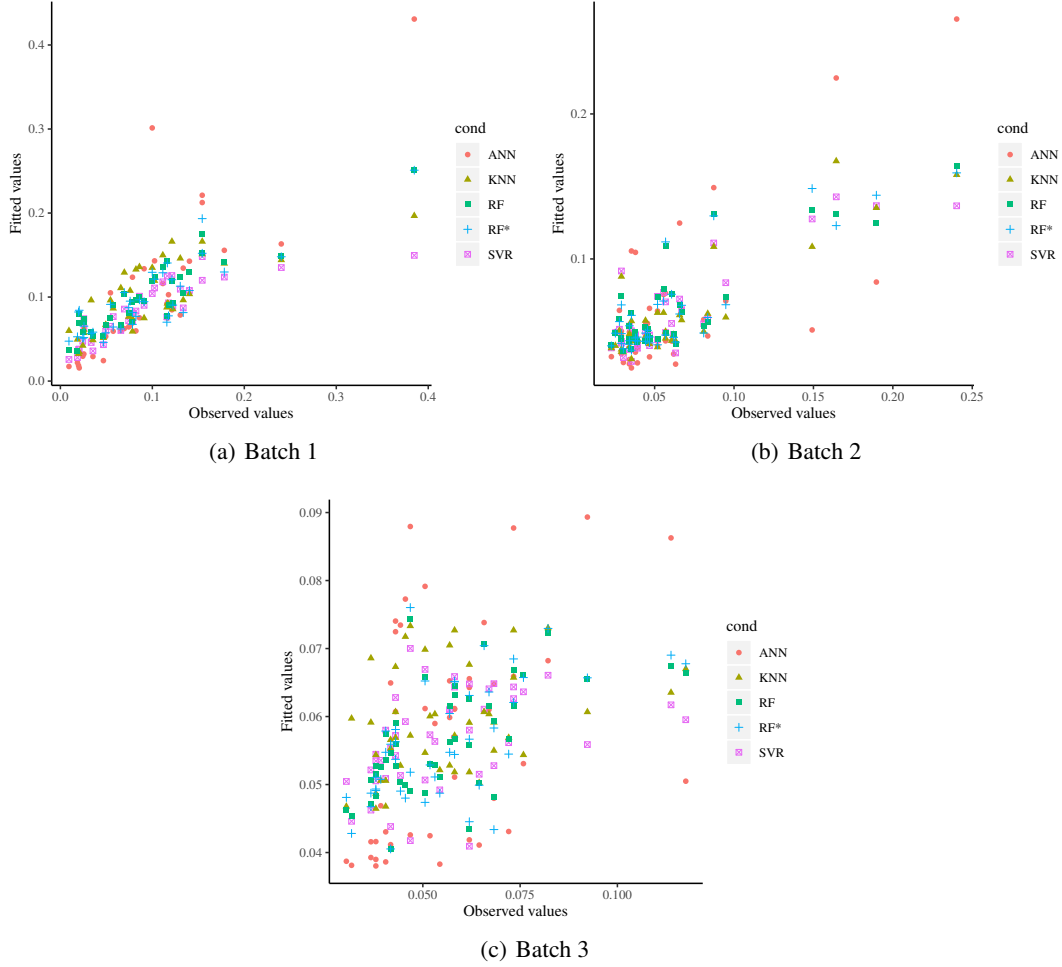


Figure 4: Scatterplot of observed and predicted values by means of different predictors: ANN, KNN, RF, RF* and SVR (see the text for more explanation).

reduced in Batches 2 and 3 to their overestimation of the very low mosquito population value points resulting from the improved control activities in the Batch periods.

In further analysis of RF*, here, we discuss the selected most informative features subset across the considered observation time periods. Figure 6 shows the average value of the MDI for each variable considering 50 RF replicas. For Batch 1, the eight selected most informative EO variables are TempN-R, TempD-U, TempN-R1, TempN-U2, TempD-R, TempN-U, NDVI-U and NDWI-U. For Batch 2, the selected variables features are: TempD-R1, NDVI-R1, TempD-R, NDWI-R2, NDVI-R, NDWI-R1, TempD-R2 and NDWI-R2. For Batch 3, they are: TempD-R, NDVI-U, TempN-U1, Prec-U1, TempN-R, TempD-U2, NDVI-U1, NDVI-R2. The improved control actions in Batches 2 and 3 are mostly responsible for the differences in important variables across all Batches. Only TempD-R is commonly selected in all the Batches. In addition, we can see that temperature features extracted from the rural surface zone provide the highest quality of information to our models in all Batches: TempN-R, TempD-R1 and TempD-R for Batches 1, 2 and 3, respectively. This is consistent with previous studies which show that non-artificial surface characterizations of the environment are more informative for predicting *Ae. aegypti* vector population [27, 28], and that the weekly (or daily) temperature is the most important environmental condition affecting the development of *Ae. aegypti* [51]. The non-synchronous effects of temperature is also seen in the selection of both TempN-R and TempN-R1 in Batch 1, and TempD-R, TempD-R1 and TempD-R2 in Batch 2. Of all eight selected subset variables, six of them lagged in both Batches 2 and 3, compared to only two in Batch 1. This shows that most of the environmental effects on vector population during the improved vector control regimes come from non-synchronous compounding effects. Six temperature variables are selected in Batch 1. This number reduced to three in both Batches 2 and 3. We deduce from this that the effect of non-temperature environmental variables increases during the improved control actions regime.

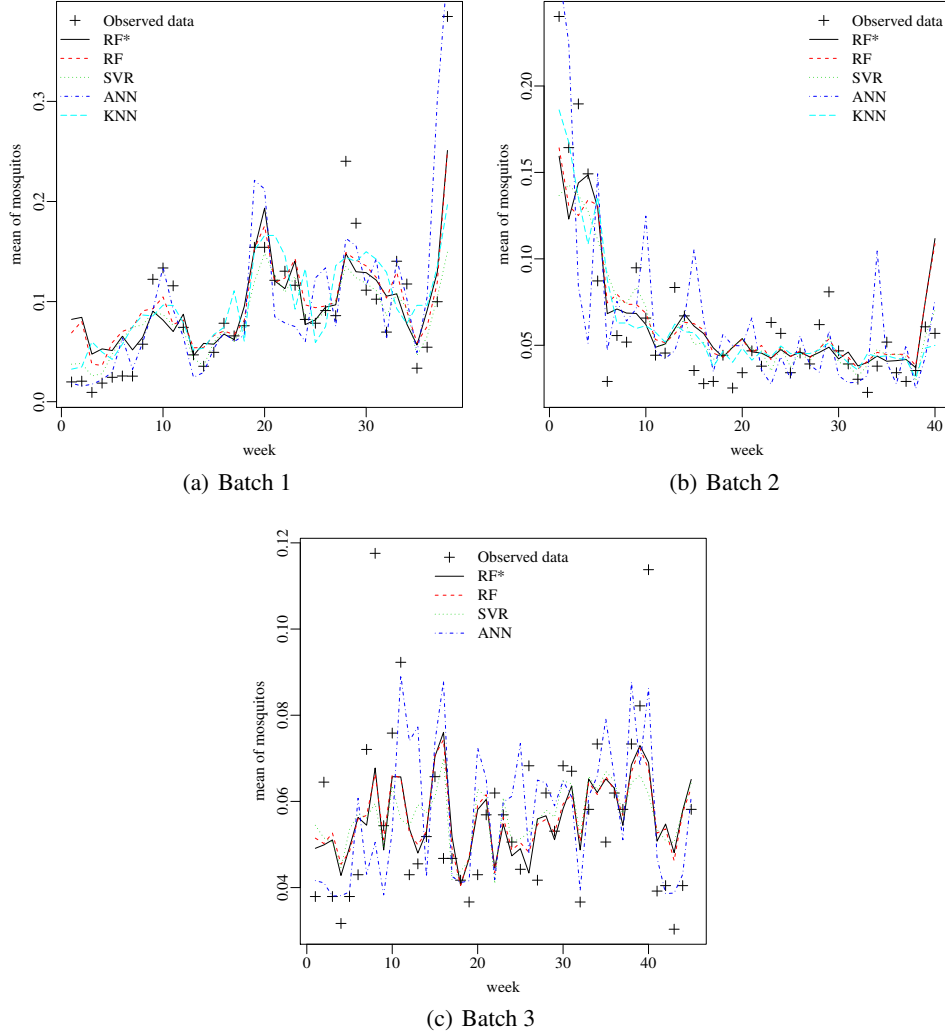


Figure 5: Time series of the actual (points) versus predicted (lines) values, using only a subset of the prediction algorithms in Fig. 4: ANN, SVR, RF and RF*.

Figure 7 presents the relationship between the mosquito population y and the standardized values of environmental features selected for RF*. It is seen that the relationships vary across the different features (x). For example, Figure 7(a) shows that in Batch 1, higher values of all six selected temperature variables, between their mean and two standard deviations above their mean, result in larger numbers of mosquitoes. Also, for Batches 2 and 3, as shown in Figures 7(b) and 7(c), even with the improved anthropological interference due to better control actions, when the three selected most informative temperature variables in each case synchronously rise above one standard deviation from their mean, there is a rapid increase in the vector population. These results are in accordance with laboratory studies presented in [39, 64]. These works show that at higher temperatures below 40 °C, the development life cycle of *Ae. aegypti* is accelerated, thus increasing the vector population. In addition, the extrinsic incubation period (EIP) of *Ae. aegypti* – the time interval between virus agent acquisition by the mosquito vector and the moment it is able to transmit it to humans – reduces at high temperature [30, 65]. The results in Figure 7(b) also show that the improved control program in starting from Batch 2 is effective except in cases of synchronous extreme values of important temperature variables. This observation shows pointers that can be used to further improve the efficacy of the control actions in subsequent years.

Regarding non-temperature environmental effects, the relation curves for Batch 3, Figure 7(c), also show that at about -0.5 standard deviation below the mean of Prec-U1, there is a spike in the vector population. As reported in [27, 39], in regions where primary larval sites are in rain-filled containers, rainfall has been shown to positively correlate with larval

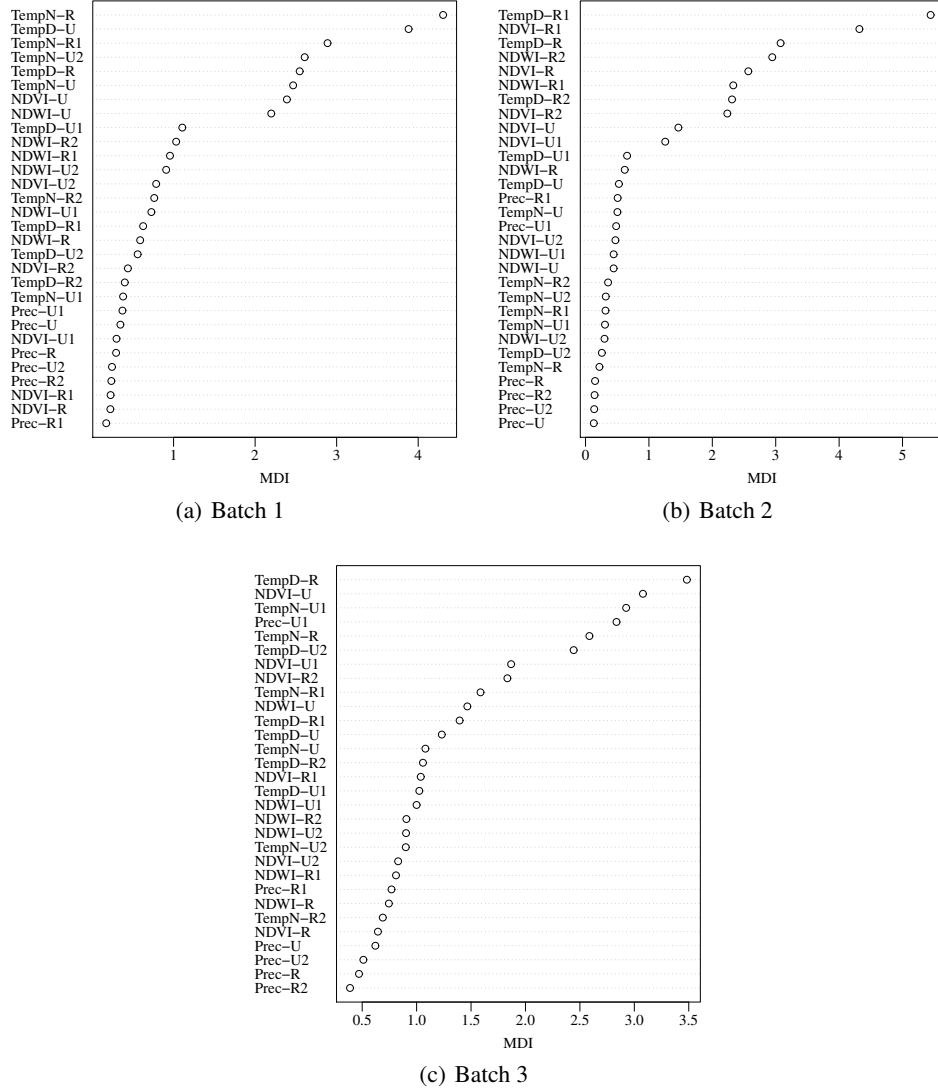


Figure 6: Average values of MDI considering 50 RF replicas.

and adult mosquito abundance. Further studies have corroborated the importance of humidity and vegetation conditions in abundance and reproduction of mosquito species, with longer dry season and lower relative humidity resulting in higher egg mortality [66]. Fully developed tree canopies by providing shade can reduce evaporation of hatching water, and can also increase near ground humidity, thus increasing density of *Ae. aegypti* mosquito larvae [21, 67]. Stronger positive effects of vegetation conditions can be observed through the effects of NDVI-R2 and NDVI-U1 in Batch 3.

Finally, and in general, this research shows, based on our study data, that the random forest model performs better than all the other models for predicting the mosquito population from EO-based variables. By selecting the most informative features using the MDI measure, it is possible to obtain a much less complex model with comparable (or even better) results, labeled in this work as RF*. Additionally, with the MDI ranking it is possible to understand and explain the rationale for the model variables and analyse the effects of each EO variable on the vector population. This better understanding is gained by considering the relation graphs based on RF*, which provide a visualization of the nonlinear relationships among the variables.

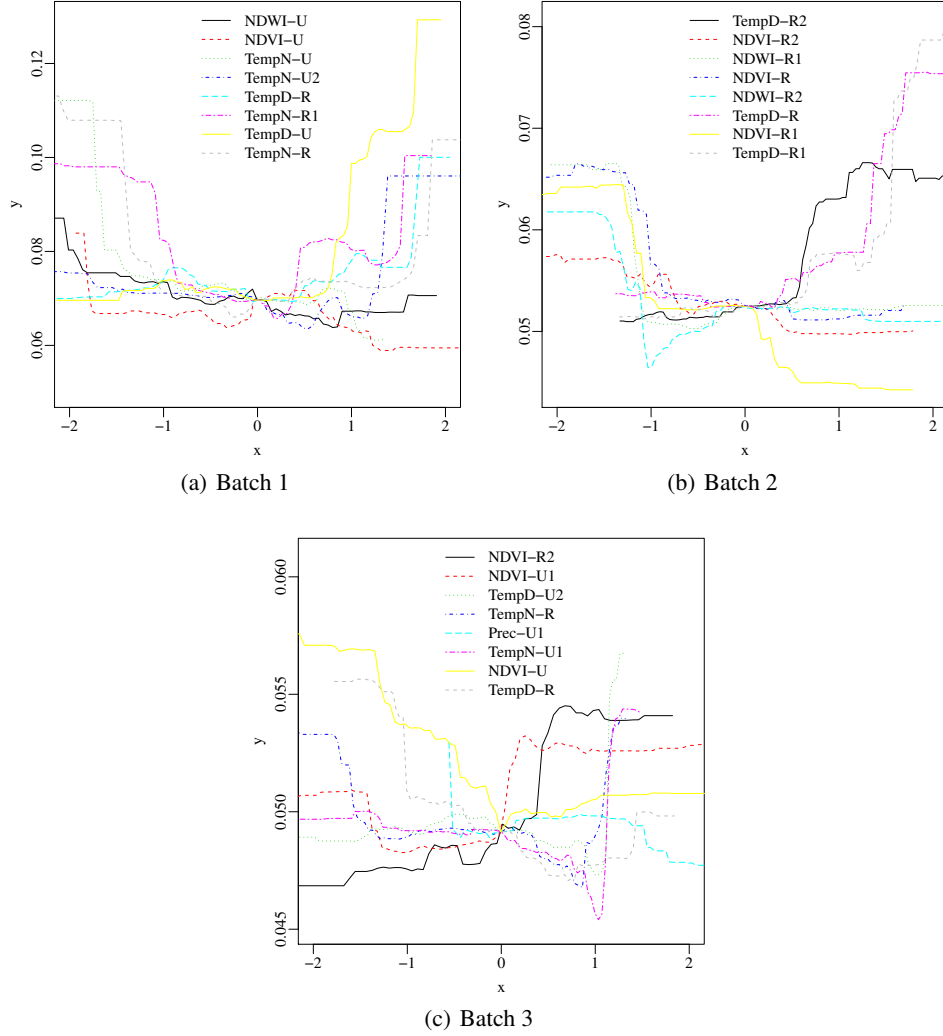


Figure 7: The relation between the mean of mosquitoes (y) and each covariate (x).

5 Conclusions

In this paper, a procedure for modeling the temporal population distribution of female *Ae. aegypti* mosquito species based effects environmental factors estimated from freely available EO data has been presented. The procedure was implemented using MODIS and GPM EO data products to obtain information about temperature, precipitation, moisture and vegetation conditions. The mosquito vector population data was collected on the field from April 2017 to August 2019 (124 weeks) using 791 mosquito traps. We divided the field collected data in time into three operational Batches based on the difference in data collection conditions as obtained from the field staff. From January 2018, there were data-driven improvements to vector control program, affecting data Batches 2 and 3 in our study.

Different studies in the recent past have been conducted in this domain due to the rise in availability of more freely available earth observation data and off-the-shelf machine learning analyses tools [27, 28]. These studies, being prediction focused, do not consider explainable modeling using non-linear machine learning techniques like RF.

For this reason, here, we have presented a procedure for explainable modeling of female *Ae. aegypti* based on EO data estimated environmental factors based on RF. A quantitative measure of the variable importance (MDI)— wrapped in RF — was used to extract the most informative environmental features, and obtain a less complex but still accurate predictive RF model, labeled RF* in this study. To prove their robustness, the resulting models were compared to other machine learning models including SVR, ANN, KNN and DTR, as well as statistical models, such as LM and GLM.

Our results show that our proposed RF-based approach is capable of mapping the complex relationship among the EO variables and vector population. Furthermore, the features selected thanks to the MDI value ranking can be empirically interpreted and provide hints about the relationships among vector population and these environmental conditions from an operational point of view.

Acknowledgment

The first author gratefully acknowledges the financial support by the European Commission through Horizon 2020 research and innovation programme, grant agreement No 734541 - Project "EOXPOSURE". FMB thanks to CNPq/MCTI, Brazil. AEE thanks CAPES/MEC, CNPq/MCTI and Decit/SCTIE/MS (grant No 440358/2016-7). All the authors also thank Ecovec for providing the mosquito trap data and Marcelaine Marculano (Vigilância Ambiental) for the entomological data.

References

- [1] Catherine Linard, Andrew J. Tatem, and Marius Gilbert. Modelling spatial patterns of urban growth in Africa. *Applied Geography*, 44:23–32, 2013.
- [2] Jean-Paul Donnay, Mike J Barnsley, and Paul A Longley. *Remote Sensing and Urban Analysis: GISDATA 9*. CRC Press, 2000.
- [3] Hasi Bagan and Yoshiki Yamagata. Landsat analysis of urban growth: How Tokyo became the world's largest megacity during the last 40years. *Remote Sensing of Environment*, 127:210–222, 2012.
- [4] Paolo E. Gamba Qihao Weng, Dale Quattrochi. *Urban Remote Sensing, Second Edition*. CRC Press, 2018.
- [5] Qihao Weng. *Global Urban Monitoring and Assessment through Earth Observation*. CRC Press, 2014.
- [6] Aidan While and Mark Whitehead. Cities, urbanisation and climate change. *Urban Studies*, 50(7):1325–1331, 2013.
- [7] Brian L. Markham and Dennis L. Helder. Forty-year calibrated record of Earth-reflected radiance from Landsat: A review. *Remote Sensing of Environment*, 122:30–40, 2012.
- [8] Nirav N Patel, Emanuele Angiuli, Paolo Gamba, Andrea Gaughan, Gianni Lisini, Forrest R Stevens, Andrew J Tatem, and Giovanna Trianni. Multitemporal settlement and population mapping from Landsat using Google Earth Engine. *International Journal of Applied Earth Observation and Geoinformation*, 35:199–208, 2015.
- [9] Zhe Jiang and Shashi Shekhar. *Spatial Big Data Science*. Springer, 2017.
- [10] Andrea Marinoni and Paolo Gamba. Unsupervised data driven feature extraction by means of mutual information maximization. *IEEE Transactions on Computational Imaging*, 3(2):243–253, 2017.
- [11] Agnes Turyahikayo. Assessing land use induced disturbance to vegetation cover in the upper Molopo catchment, South Africa, using Landsat images. In *2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE, 2015.
- [12] R. Ostfeld, G. Glass, and F Keesing. Spatial epidemiology: an emerging (or re-emerging) discipline. *Trends in Ecology & Evolution*, 20(6):328–336, 2005.
- [13] William K. Reisen. Landscape epidemiology of vector-borne diseases. *Annual Review of Entomology*, 55(1):461–483, 2010.
- [14] Sean G. Young, Jason A. Tullis, and Jackson Cothren. A remote sensing and GIS-assisted landscape epidemiology approach to West Nile virus. *Applied Geography*, 45:241–249, 2013.
- [15] Jean-Pierre Rossi, Ibrahima Kadaouré, Martin Godefroid, and Gauthier Dobigny. Landscape epidemiology in urban environments: The example of rodent-borne trypanosoma in Niamey, Niger. *Infection, Genetics and Evolution*, 63:307–315, 2018.
- [16] Tammi L Johnson, Ubydul Haque, Andrew J Monaghan, Lars Eisen, Micah B Hahn, Mary H Hayden, Harry M Savage, Janet McAllister, John-Paul Mutebi, and Rebecca J Eisen. Modeling the environmental suitability for *Aedes (stegomyia) aegypti* and *Aedes (stegomyia) albopictus* (diptera: Culicidae) in the contiguous United States. *Journal of Medical Entomology*, 54(6):1605–1614, 2017.
- [17] Kimberly M. Fornace, Chris J. Drakeley, Timothy William, Fe Espino, and Jonathan Cox. Mapping infectious disease landscapes: unmanned aerial vehicles and epidemiology. *Trends in Parasitology*, 30(11):514–519, 2014.

- [18] Amanda Jamison, Elaina Tuttle, Ryan Jensen, Greg Bierly, and Rusty Gonser. Spatial ecology, landscapes, and the geography of vector-borne disease: A multi-disciplinary review. *Applied Geography*, 63:418–426, 2015.
- [19] Camilo Rotela, Laura Lopez, María Frías Céspedes, Gabriela Barbas, Andrés Lighezzolo, Ximena Porcasi, Mario A Lanfri, Carlos M Scavuzzo, and David E Gorla. Analytical report of the 2016 dengue outbreak in Córdoba city, Argentina. *Geospatial Health*, 2017.
- [20] Ximena Porcasi, Camilo H. Rotela, María V. Intorini, Nicolás Frutos, Sofía Lanfri, Gonzalo Peralta, Estefanía A. De Elia, Mario A. Lanfri, and Carlos M. Scavuzzo. An operative dengue risk stratification system in Argentina based on geospatial technology. *Geospatial Health*, 6(3):31, 2012.
- [21] Jane P Messina, Moritz UG Kraemer, Oliver J Brady, David M Pigott, Freya M Shearer, Daniel J Weiss, Nick Golding, Corrine W Ruktanonchai, Peter W Gething, Emily Cohn, John S Brownstein, Kamran Khan, Andrew J Tatem, Thomas Jaenisch, Christopher JL Murray, Fatima Marinho, Thomas W Scott, and Simon I Hay. Mapping global environmental suitability for Zika virus. *eLife*, 5, 2016.
- [22] Christopher N.H. Doll, Jan-Peter Muller, and Jeremy G. Morley. Mapping regional economic activity from night-time light satellite imagery. *Ecological Economics*, 57(1):75–92, 2006.
- [23] Emanuele Angiuli and Giovanna Trianni. Urban mapping in Landsat images based on normalized difference spectral vector. *IEEE Geoscience and Remote Sensing Letters*, 11(3):661–665, 2014.
- [24] Annelise Tran, Assane Gueye Fall, Biram Biteye, Mamadou Ciss, Geoffrey Gimonneau, Mathieu Castets, Momar Talla Seck, and Véronique Chevalier. Spatial modeling of mosquito vectors for Rift Valley fever virus in Northern Senegal: Integrating satellite-derived meteorological estimates in population dynamics models. *Remote Sensing*, 11(9), 2019.
- [25] World Health Organization. Mosquito-borne diseases, 2016.
- [26] Thomas W. Scott, John D. Edman, Pattamaporn Kittayapong, Jonathan F. Day, and Amara Naksathit. A fitness advantage for aedes aegypti and the viruses it transmits when females feed only on human blood. *The American Journal of Tropical Medicine and Hygiene*, 57(2):235–239, aug 1997.
- [27] Juan M. Scavuzzo, Francisco Trucco, Manuel Espinosa, Carolina B. Tauro, Marcelo Abril, Carlos M. Scavuzzo, and Alejandro C. Frery. Modeling Dengue vector population using remotely sensed data and machine learning. *Acta Tropica*, 185:167 – 175, 2018.
- [28] A. German, M.O. Espinosa, M. Abril, and C.M. Scavuzzo. Exploring satellite based temporal forecast modelling of Aedes aegypti oviposition from an operational perspective. *Remote Sensing Applications: Society and Environment*, 11:231 – 240, 2018.
- [29] Fernando Jose Antonio, Andreia Silva Itami, Sergio de Picoli, Jorge Juarez Vieira Teixeira, and Renio dos Santos Mendes. Spatial patterns of Dengue cases in Brazil. *PloS one*, 12(7):e0180715, 2017.
- [30] Zhuanzhuan Liu, Zhenhong Zhang, Zetian Lai, Tengfei Zhou, Zhirong Jia, Jinbao Gu, Kun Wu, and Xiao-Guang Chen. Temperature increase enhances Aedes albopictus competence to transmit Dengue virus. *Frontiers in Microbiology*, 8, 2017.
- [31] SCOTT A Ritchie. The production of Aedes aegypti by a weekly ovitrap survey [Dengue, Caribbean]. *Mosquito News*, 1984.
- [32] Alvaro E Eiras, Marcelo C Resende, José L Acebal, and Kelly S Paixão. New cost-benefit of Brazilian technology for vector surveillance using trapping system. In *From Local to Global Impact of Mosquitoes*. IntechOpen, 2018.
- [33] David J. Lary, Amir H. Alavi, Amir H. Gandomi, and Annette L. Walker. Machine learning in geosciences and remote sensing. *Geoscience Frontiers*, 7(1):3–10, 2016.
- [34] A. Karpatne, I. Ebert-Uphoff, S. Ravela, H. A. Babaie, and V. Kumar. Machine learning for the geosciences: Challenges and opportunities. *IEEE Transactions on Knowledge and Data Engineering*, 31(8):1544–1554, 2019.
- [35] Qing Yao, Zexin Guan, Yingfeng Zhou, Jian Tang, Yang Hu, and Baojun Yang. Application of support vector machine for detecting Rice diseases using shape and color texture features. In *2009 International Conference on Engineering Computation*. IEEE, 2009.
- [36] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [37] Gérard Biau and Erwan Scornet. A random forest guided tour. *TEST*, 25(2):197–227, 2016.
- [38] P. McCullagh and J.A. Nelder. *Generalized Linear Models*. Chapman and Hall, 2nd edition, 1989.
- [39] Cassie C. Jansen and Nigel W. Beebe. The Dengue vector Aedes aegypti: what comes next. *Microbes and Infection*, 12(4):272–279, 2010.

- [40] Alexandre Rosa dos Santos, Felício Santos de Oliveira, Aderbal Gomes da Silva, José Marinaldo Gleriani, Wantuelfer Gonçalves, Giselle Lemos Moreira, Felipe Gimenès Silva, Elvis Ricardo Figueira Branco, Marks Melo Moura, Rosane Gomes da Silva, Ronie Silva Juvanhól, Kaíse Barbosa de Souza, Carlos Antonio Alvares Soares Ribeiro, Vagner Tebaldi de Queiroz, Adilson Vidal Costa, Alexandre Simões Lorenzon, Getulio Fonseca Domingues, Gustavo Eduardo Marcatti, Nero Lemos Martins de Castro, Rafael Tassinari Resende, Duberli Elera Gonzales, Lucas Arthur de Almeida Telles, Thaisa Ribeiro Teixeira, Gleissy Mary Amaral Dino Alves dos Santos, and Pedro Henrique Santos Mota. Spatial and temporal distribution of urban heat islands. *Science of The Total Environment*, 605-606:946–956, 2017.
- [41] Elisângela Martins da Silva Costa, Rivaldo Venâncio da Cunha, and Edgar Aparecido da Costa. National dengue control program implementation evaluation in two border municipalities in Mato Grosso do Sul State, Brazil, 2016. *Epidemiologia e Serviços de Saúde*, 27(4), 2018.
- [42] Helena Araújo, Danilo Carvalho, Rafaella Ioshino, André Costa-da Silva, and Margareth Capurro. Aedes aegypti control strategies in Brazil: incorporation of new technologies to overcome the persistence of dengue epidemics. *Insects*, 6(2):576–594, 2015.
- [43] Ecovec Company. About ecovec.
- [44] Álvaro Eduardo Eiras and Marcelo Carvalho Resende. Preliminary evaluation of the "dengue-mi" technology for aedes aegypti monitoring and control. *Cadernos de Saúde Pública*, 25:S45–S58, 2009.
- [45] Diogo Portella Ornelas de Melo, Luciano Rios Scherrer, and Álvaro Eduardo Eiras. Dengue fever occurrence and vector detection by larval survey, ovitrap and MosquiTRAP: A space-time clusters analysis. *PLoS ONE*, 7(7):e42125, jul 2012.
- [46] C.O Justice, J.R.G Townshend, E.F Vermote, E Masuoka, R.E Wolfe, N Saleous, D.P Roy, and J.T Morisette. An overview of MODIS land data processing and product status. *Remote Sensing of Environment*, 83(1-2):3–15, 2002.
- [47] Bo cai Gao. NDWI—a normalized difference water index for remote sensing of vegetation liquid water from space. *Remote Sensing of Environment*, 58(3):257–266, 1996.
- [48] Z. Wan, Y. Zhang, Q. Zhang, and Z.-L. Li. Quality assessment and validation of the MODIS global land surface temperature. *International Journal of Remote Sensing*, 25(1):261–274, 2004.
- [49] Gail Skofronick-Jackson, Walter A Petersen, Wesley Berg, Chris Kidd, Erich F Stocker, Dalia B Kirschbaum, Ramesh Kakar, Scott A Braun, George J Huffman, Toshio Iguchi, et al. The global precipitation measurement (GPM) mission for science and society. *Bulletin of the American Meteorological Society*, 98(8):1679–1695, 2017.
- [50] Elizabet L. Estallo, Mario A. Lamfri, Carlos M. Scavuzzo, Francisco F. Ludueña Almeida, María V. Introini, Mario Zaidenberg, and Walter R. Almirón. Models for predicting Aedes aegypti larval indices based on satellite images and climatic variables. *Journal of the American Mosquito Control Association*, 24(3):368–376, 2008.
- [51] Yoon Cheong, Katrin Burkart, Pedro Leitão, and Tobia Lakes. Assessing weather effects on Dengue disease in Malaysia. *International Journal of Environmental Research and Public Health*, 10(12):6319–6334, 2013.
- [52] Christopher M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1996.
- [53] Claude Sammut and Geoffrey I. Webb. *Encyclopedia of Machine Learning and Data Mining*. Springer, 2017.
- [54] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2018.
- [55] Gilles Louppe, Louis Wehenkel, Antonio Sutera, and Pierre Geurts. Understanding variable importances in forests of randomized trees. In *Advances in Neural Information Processing Systems*, pages 431–439, 2013.
- [56] L. Zhang, L. Zhang, and B. Du. Deep learning for remote sensing data: A technical tutorial on the state of the art. *IEEE Geoscience and Remote Sensing Magazine*, 4(2):22–40, 2016.
- [57] David John Lary. *Geoscience and Remote Sensing New Achievements*, chapter Artificial Intelligence in Geoscience and Remote Sensing, pages 105–128. IntechOpen, 2010.
- [58] M. Riedmiller. Rprop - description and implementation details. Technical report, University of Karlsruhe, 1994.
- [59] Leo Breiman, Jerome Friedman, Charles J. Stone, and R.A. Olshen. *Classification and Regression Trees*. Chapman and Hall/CRC, 1984.
- [60] Martin Krzywinski and Naomi Altman. Classification and regression trees. *Nature Methods*, 14:757–758, 2017.
- [61] Myles Hollander, Douglas A. Wolfe, and Eric Chicken. *Nonparametric Statistical Methods*. John Wiley and Sons, third edition, 2013.

- [62] Szu-Chieh Chen, Chung-Min Liao, Chia-Pin Chio, Hsiao-Han Chou, Shu-Han You, and Yi-Hsien Cheng. Lagged temperature effect with mosquito transmission potential explains Dengue variability in southern Taiwan: insights from a statistical analysis. *Science of the Total Environment*, 408(19):4069–4075, 2010.
- [63] LS Barsante, KS Paixão, KH Laass, RTN Cardoso, AE Eiras, and JL Acebal. A model to predict the population size of the dengue fever vector based on rainfall data. *arXiv preprint arXiv:1409.7942*, 2014.
- [64] Kiyoshi Kamumura, Ines Tomoco Matsuse, Hanako Takahashi, Jun Komukai, Takayo Fukuda, Kayo Suzuki, Miho Aratani, Yoshikazu Shirai, and Motoyoshi Mogi. Effect of temperature on the development of *Aedes aegypti* and *Aedes albopictus*. *Medical Entomology and Zoology*, 53(1):53–58, 2002.
- [65] Nils Benjamin Tjaden, Stephanie Margarete Thomas, Dominik Fischer, and Carl Beierkuhnlein. Extrinsic incubation period of Dengue: Knowledge, backlog, and applications of temperature dependence. *PLoS Neglected Tropical Diseases*, 7(6):e2207, 2013.
- [66] M Trpiš. Dry season survival of *Aedes aegypti* eggs in various breeding sites in the Dar es Salaam area, Tanzania. *Bulletin of the World Health Organization*, 47(3):433, 1972.
- [67] Douglas O Fuller, A Troyo, and John C Beier. El nino southern oscillation and vegetation dynamics as predictors of Dengue fever cases in Costa Rica. *Environmental Research Letters*, 4(1):014011, 2009.