Predicting Dengue Outbreaks with Explainable Machine Learning

Robson Aleixo¹, Fabio Kon¹, Rudi Rocha², Marcela Santos Camargo², **Raphael Y. de Camargo**³

¹Department of Computer Science, University of São Paulo, Brazil ²São Paulo School of Business Adm., Institute for Health Policy Studies (IEPS) ³Center for Mathematics, Computing and Cognition, Federal University of ABC

Al4Health, May, 2022

InterSCity

http://interscity.org

Introduction

- Dengue disease: Mosquito-borne tropical disease
- CDC estimates hundreds of millions infections
 - Tens of thousands deaths per year
- State of Rio de Janeiro, in Brazil
 - Located in a tropical area and most cities are frequently under high-risk state of alert
- Prevention of Dengue outbreaks
 - Reduce the mosquito population
 - Work-intensive: requires agents to visit locations







https://pt.wikipedia.org







- Develop a Machine Learning model for predicting Dengue outbreaks up to 3 months in advance
 - Use of multiple features, such as environmental, epidemiological, demographic, and spatial
 - Provide explainable predictions for health agents
- Evaluate the model using 6 years of data from Rio de Janeiro metropolis

Related Work

- Extensive literature on predicting Dengue cases and correlate with climatic and socioeconomical variables
 - Most works on tropical countries: Thailand, Indonesia, Malaysia, and Latin America
 - Used different ML techniques, from Linear Models to Neural Networks
- A limitation of existing work is that they provide little explainability of predictions
 - Linear models provide limited interpretability
- A second limitation is the use of different datasets, which are unavailable, and metrics
 - We will make all model code and data publicly available
 - We evaluate and compare several regression and classification metrics

Data and Feature Extraction

- Obtained from government agencies, such as
 - National Health Notification Information System (SINAN)
 - National Register of Health Facilities (CNES)
 - National Weather Institute (INMET)
- Dengue cases: last 3 months, prevalence, neighbors
- Environmental
- Mosquito infestation in the city and similar
- Demographic density and health coverage

Feature
cases m-1
cases m-2
cases m-3
dengue_prevalence
neighbor_cases
precipitation (mm)
temperature (°C)
air_humidity (%)
liraa
chikungunya
zika
demographic density
num_health_unit

Regression Method

- Boosted-tree regression method: CatBoost
 - Series of simple trees, where each tries to improve the prediction of previous trees
 - Captures non-linear relationships between features
- Baseline: Seasonal Autoregressive Integrated Moving Average (SARIMA)
- Used data from 2015 to optimized hyperparameters



Predictions

- Predictions for each district, using data from 2016 to 2020
- CatBoost: 5-fold validation, using 4 years for training and 1 for testing
 - 3-month multistep-ahead predictions
- SARIMA model: separate model for each year
 - Used the four previous years to adjust the model
- Evaluated Regression and Classification errors
 - Regression: R², MAE, MAPE, RMSE
 - Classification: precision and recall
 - Severe and Mild outbreaks





Outbreak Predictions 3 months in advance

- Three months in advance provides enough time for health authorities to act
- Predicted no outbreaks 97% of the time
- Recall of 59% (58%) for mild (severe) outbreaks
- When predicting an outbreak as mild
 - 57% chance of a mild or severe outbreak
- When predicted as severe
 - 86% chance of a mild or severe outbreak
- SARIMA model had more modest performance
 - Precision of 40% and 38% for mild and severe outbreak predictions



Probability of Model Prediction for each Outbreak Group

Number of Model Predictions for each Outbreak Group



Seasonal Effects

- Dengue has a seasonal pattern
 - High-season (January to May)
 - Varying outbreak degrees
- Our model captures well the seasonality
 - Results for one month are closer to real values, but three-months also worked well
 - It tended to overestimate when there were fewer cases overall
- SARIMA also captures the seasonality, although it tended to overshoot predictions



Evaluation per District

- MAE error unevenly distributed among districts
- Higher for districts which had very large number of cases, such as Bangu and Realengo
- F-score there were comparable to other districts, as the outbreaks were detected
- Harder to detect outbreaks in districts with fewer cases
- Remain mostly in the borderline of the outbreak threshold



Explaining the Results

- SHAP (SHapley Additive exPlanations) values
 - Method to explain the contribution of each feature on each prediction
 - Has many desirable features: consistency, symmetry, additiviness, etc.

	Real: 969 Predicted:1060		
-	1000	0 1000	
cases_m-1	(44	2)	
cases_m-3		(37)	
temperature (°C)		(25)	
dengue_prevalence		(1)	
cases_m-2		(93)	
air_humidity (%)		(70)	
zika		(14)	
chikungunya		(0)	
neighbor_cases		(169)	
precipitation (mm)		(64)	
demograph c dens ty		(70)	
liraa		(1)	
num_health_unit		(34)	
-	1000 Model ou	0 1000 tput value	

	Real: 79 Predicted:72		
	- 50	0	50
precipitation (mm)		(24	18)
cases_m-1		(12	2)
neighbor_cases			(116
cases_m-2			(14)
dengue_prevalence			(0.13)
temperature (°C)			(27)
num_health_unit		(30)
cases_m-3		(4)
air_humidity (%)		('	70)
chikungunya		(()
demographicidensity		(2	2)
liraa		(0	.9)
zika		()	D)
	- 50 Model	outpu	50 It value

Explanations for the Complete Model

- Summary plot:
 - Features that most contributed to all predictions and the output direction
- Dependence plots:
 - Effects on predictions for feature values and dependence with another feature



Conclusions

- Developed an explainable machine learning model to predict Dengue outbreaks
 - Predictions three-month in advance provide health agents enough time to act
 - Combined with explanations for the predictions for better interpretation
 - Can be combined with other information available at health agencies
- Next steps
 - Enhancing the model with better features, such as sorotypes and other indicators
 - Extraction of better time series features and inclusion of data from more regions
 - Provide open-source visualization tool for health authorities



raphael.camargo@ufabc.edu.br

http://interscity.org





