

Title goes here...

Javier Corvillo^{1, 2*}, Verónica Torralba^{2*}, Diego Campos², Ángel G. Muñoz^{2*}, and Ana Riviére-Cinnamond³

¹Complutense University of Madrid, Department of Earth Science and Astrophysics, Madrid, 28040, Spain

²Barcelona Supercomputing Center, Earth Sciences Department, 08034, Spain

³Pan-American Health Organization, Communicable Diseases and Health Analysis, Panama City, 0843-03441, Panama

*javier.corvillo@bsc.es / veronica.torralba@bsc.es / angel.g.munoz@bsc.es

ABSTRACT

This will be the abstract for the paper...

1 Introduction

Introduction text goes here...

2 Methods

However, the transmission dynamics of these diseases, while partly driven by climate, are not so well understood.

Using data from a number of sources, we explore and analyze the behavior of the climate-component of *Aedes*-borne disease transmission, in order to understand its role on the dynamics of disease outbreaks in the context of a changing climate. Our analysis is composed of three different studies: 1) a timescale decomposition of disease transmissibility values, thereby guiding officials to understand the climatology of outbreaks for budget and resource allocation; 2) a correlation analysis between transmissibility values and different climate variability indices, such as El Niño Southern Oscillation, in order to understand the effects of natural climate patterns onto *Aedes*-borne outbreaks; and 3) a causality analysis to solidify findings obtained through correlation, identifying the most relevant predictors and their potential application under a climate-and-health service framework for forecasting the transmissibility of *Aedes*-borne diseases.

2.1 Analysis 1: Multi-timescale climate decomposition of R_0

With the intent of isolating the human-driven signal from the natural variability R_0 data, a "timescale decomposition" methodology was used to obtain the total variance across different time-scales.

2.1.1 Data

The timescale decomposition analysis was undertaken using R_0 outputs from the *Aedes* Disease Environmental Suitability 2's (AeDES2) monitoring system (citation would go here). The basic reproduction number, or R_0 , is a metric that quantifies the transmissibility of vector-borne diseases, and is defined as the average number of secondary cases generated by a single infected individual in a completely susceptible population. Though R_0 is normally obtained using by combining climate data, vector biology, and human behavior, AeDES2's R_0 values consider climate data as the only factor for disease transmission, as the system is designed to be used as a climate service for vector-borne diseases, giving a first approximation of the transmissibility of these diseases. Therefore, while these climate-based R_0 values are not strictly a tried-and-true metric for epidemiological purposes, its values are still useful for understanding the climatological behavior of vector-borne diseases in the past, ideal for the purposes of this study.

The 1980-2021 monthly-mean period of AeDES2's R_0 values was selected for the analysis, for a total of 504 months or 167 full seasons. Considering that vector borne diseases are extending to previously unaffected areas due to the effects of man-made climate change, AeDES2's coverage has been increased since its inception to contain global outputs, allowing for a comprehensive analysis of the relationship between climate variability indices and R_0 both in current *Aedes* hotspots and emerging regions.

Index Name	Abbreviation	Periodicity	Pattern Type	Source	Detrending Method
Arctic Oscillation	AO	Several weeks to months	Atmospheric	NOAA's Climate Prediction Center (CPC)	Linear
Atlantic Multidecadal Oscillation	AMO	Between 60-80 years	Oceanic	NOAA's Kaplan Extended SST data	Linear
Atlantic 3 Index	ATL3	Several months to a few years	Oceanic	Detrended ORAS5 reanalysis data	-
Indian Ocean Basin	IOB	Several months to a few years	Oceanic	Detrended ORAS5 reanalysis data	-
Indian Ocean Dipole	IOD	Between 2-7 years	Oceanic	Detrended ORAS5 reanalysis data	-
North Atlantic Oscillation	NAO	Several days to decades	Atmospheric	NOAA's CPC	Linear
El Niño 3.4 Index	Niño 3.4	Between 2-7 years	Oceanic	NOAA's CPC	Linear
North Pacific Meridional Mode	NPM	Several months to a few years	Atmospheric	Detrended ORAS5 reanalysis data	-
Pacific Decadal Oscillation	PDO	Between 20-30 years	Oceanic	NOAA's ERSSTv5	Linear
Pacific-North American Pattern	PNA	Several weeks to months	Atmospheric	NOAA's CPC	Linear
Quasi Biannual Oscillation	QBO	~2 years	Atmospheric	NOAA's NCEP/NCAR Reanalysis	Linear
South Atlantic Subtropical Dipole 1	SASD1	Several months to a few years	Oceanic	Detrended ORAS5 reanalysis data	-
Southern Indian Ocean Dipole	SIOD	Several months to a few years	Oceanic	Detrended ORAS5 reanalysis data	-
Southern Oscillation Index	SOI	Between 2-7 years	Pressure-based	NOAA's CPC	Linear
South Pacific Meridional Mode	SPMM	Several months to a few years	Atmospheric	Detrended ORAS5 reanalysis data	-
Tropical North Atlantic	TNA	Several months to a few years	Oceanic	Detrended ORAS5 reanalysis data	-

Table 1. Summary of the climate variability indices used in the analysis used for the correlation and causality studies.

2.1.2 Methodology (Time-based)

As R_0 doesn't follow a clearly defined probability distribution function, the temporal analysis filters a given R_0 time-series of any given grid-point by employing a locally estimated scatterplot smoothing technique (LOESS). In order to obtain the best smoothing parameter for the analysis, a search was performed over a range of values between 1 and 504 months for the fit of the spatial median of the R_0 observational data. The best smoothing parameter for said regression was selected using verification metrics for the goodness of fit of the model: the highest R squared value (RSE), the lowest Akaike Information Criterion (AIC) value, or the lowest GCV value. The GCV value is prioritized over the other two metrics, as it is a more robust measure of the goodness of fit of the model to the data, penalizing overfitting.

Once the ideal smoothing parameter is found for the R_0 data, the R_0 time-series for each gridpoint is separated into four components: a long-term trend signal (understood to be the trend caused by anthropogenic climate change), an inter-annual signal (year to year), a decadal signal (10-30 years), and lastly, a remainder signal which contains other signals of the data (i.e., inter-annual and inter-decadal variability, among others). Variance maps for each of these four components capture the overall direction of the data over time, as well as the climatological variability of R_0 in any given grid-point.

Once variance maps obtained, Strongest Seasonal Signal Regions (SSSRs) are identified as regions with a significant percentage of variance explained by the seasonal component of the data, to be used in the following parts of the analysis. The boundaries for the selection of SSSR regions are defined by the current Intergovernmental Panel on Climate Change set of reference regions for subcontinental analysis of climate model data (Iturbide et al., 2020).

2.2 Analysis 2: Correlation studies between R_0 and climate variability indices

After analyzing the R_0 signal and its variability through timescale decomposition, we assess the impact of several climate variability indices on global R_0 values over the chosen 1980-2021 monthly-mean period.

2.2.1 Data

Correlation studies are performed over both global and SSSR regions, using the same R_0 data as in the previous analysis. A total of 16 climate variability indices have been used for the correlation analysis. Their respective sources, as well as detrending methods for each, are listed and summarized in Table 1.

2.2.2 Methodology

The correlation analysis was performed using the Pearson correlation coefficient, which quantifies the linear relationship between two variables. The Pearson correlation coefficient is defined as:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

where x_i and y_i are the values of the two variables, \bar{x} and \bar{y} are their respective means, and n is the number of observations. The correlation coefficient ranges from -1 to 1, where -1 indicates a perfect negative correlation, 0 indicates no correlation, and 1 indicates a perfect positive correlation. For computation of statistical significance in correlation, a Monte Carlo method was used, with a p-value threshold of 0.05.

In order to avoid spurious correlation outputs, anthropogenic signal from the natural variability of the data is isolated by using the long-term trend component obtained in the timescale decomposition analysis. For the climate variability indices,

we apply the detrending methods listed in Table 1. The correlation analysis is performed over the different seasons over the 1980-2022 period.

2.3 Analysis 3: Causality studies between R_0 and climate variability indices. Outlining of predictors for disease outbreaks

Causal-based patterns can be identified after this analysis, which, as opposed to correlation, allow for a more robust foundation for the understanding of the underlying mechanisms between climate variability and R_0 patterns. In discarding potentially spurious results obtained through correlation, this causality analysis can be used to outline the most relevant predictors for disease outbreaks. These predictors, in turn, can be used for the refining and building of AeDES2's prediction system for improving the accuracy and skill of the ensemble forecasts respect its predecessor.

2.3.1 Data

The datasets that were used for the causality analysis are the same as those employed in assessing the impact of climate variability indices on R_0 values across the globe (Section 2.2.1).

2.3.2 Methodology

Causality analysis between R_0 and climate variability indices was performed by using Liang-Kleeman's proposed methodology for computing information flow between two entities of a dynamical system (Liang and Kleeman, 2005), quantifying the amount of information that one time series (the climate variability indices) can provide about another time series (R_0 patterns). This formalism is based on the concept of transfer entropy, which allows to compute causality as:

$$T_{2 \rightarrow 1} = \frac{r}{1 - r^2} (r'_{2,\partial 1} - r'_{1,\partial 1}) \quad (2)$$

where $T_{2 \rightarrow 1}$ is the rate of entropy transfer from time series 2 to time series 1, r is the correlation coefficient between the two time series, and $r'_{2,\partial 1}$ and $r'_{1,\partial 1}$ are the partial correlation coefficients of time series 2 and 1 with respect to each other. While normalizing the transfer entropy can help to streamline the analysis, it is not advised for the purposes of this study, as higher correlation values in the denominator of the equation will naturally lead to very high values of transfer entropy that can influence in the normalization process.

Much like for the correlation analysis, the causality analysis is performed over the different seasons, regions, and time period, with the detrending methods listed in Table 1. For the causality analysis, a p-value threshold of 0.05 was used for statistical significance, which, following Liang-Kleeman's causality formalism, has been computed using Fisher's information matrix.

3 Results

Results text goes here...

3.1 Analysis 1: Multi-timescale climate decomposition of R_0

3.2 Analysis 2: Correlation studies between R_0 and climate variability indices

3.3 Analysis 3: Causality studies between R_0 and climate variability indices. Outlining of predictors for disease outbreaks

4 Discussion

Discussion text goes here...

5 Conclusion

Conclusion text goes here...

Figure references

Author contributions statement

Á.M., V.T. and D.C. conceived the methodology to be undertaken in this manuscript. Data sources, code and figures were obtained and developed from the ground up by J.C, who also analysed the results. All authors have reviewed the manuscript.

Code and data availability statement

Code for the generation of R_0 values employed for this study, computed using AeDES2's monitoring system, is available under request, and its values can be visualized in an operational, in-development Shiny App (link) for any region and grid-point. Additionally, the necessary datasets, functions and scripts to generate the maps and plots for this manuscript and supplementary material are available under the following GitHub repository: https://github.com/jacorvillo/monitoring_system_analysis

Competing interests

The authors declare no competing interests.