

# LEVERAGING AI TO PREDICT AND EXPLAIN DISEASE INCIDENCE FROM CLIMATE DATA

Franka Exner\*, Xufan Lu\*, Shaoming Zhang\*, Enkelejda Kasneci, Yao Rong

School of Computation, Information and Technology  
Technical University of Munich

## ABSTRACT

With the anticipated rise in temperatures due to climate change, organic vectors – known as carriers of certain diseases – may become more active in the coming years. Therefore, it is necessary to forecast the breakout of these diseases based on historical climate data. In this project, we aim to tackle this research challenge by designing a robust predictive model that can predict disease incidence. Specifically, our proposed framework named PEC (Predict and Explain from Climate Data) utilizes time-series climate data complemented by other factors such as land usage (image) or human activity data for a specific area. Moreover, PEC is designed to be conditioned on epidemiological information, enabling context-aware predictions. Beyond this, our framework incorporates post hoc model explanations on predictions, increasing transparency and trustworthiness of the model. In our preliminary experiments, we use Lyme Borreliosis disease as a case study and we are able to effectively demonstrate the correlation between the disease incidence and climate factors, highlighting the potential of using the proposed model in practical deployment.

## 1 INTRODUCTION

Climate change poses direct and indirect threats to human health, highlighting a complex relationship with our environment. This includes the effect of climate on infectious diseases, mortality, and various health outcomes (Rocque et al., 2021). Specifically, climate change impacts the behavior and distribution of vectors like ticks, which transmit diseases such as Lyme Borreliosis. Warmer winters and altered seasonal patterns due to climate change increase tick activity and, therefore, the risk of disease transmission (Gray et al., 2009). Climatic factors significantly influence the life cycle of vectors (ticks), the exact relationship is complicated by numerous factors. Despite evidence of climate change data in tick survival and activity, the precise modeling of disease incidence based on climate data, as well as other factors such as human activity, remains challenging (Vermont, 2023). Therefore, utilizing machine learning models can successfully capture the intrinsic relationship between climate data and disease incidence.

Prior work has demonstrated the potential of AI in predicting the impact of climate change or aiding in climate change mitigation across multiple sectors (Sandalow et al., 2023). For instance, Zeng & Bertsimas (2023) propose to use a multimodal model for multi-year global flood risk prediction. The model extracts embeddings from text-based geographical data and time-series climate data. Moreover, Ning et al. (2023) demonstrate that using graph re-sampling technique and Graph Neural Networks (GNNs) achieves precise global sea surface temperature forecasting. In this study, we address a novel and significant challenge in the sector of *public health*: forecasting disease incidence attributed to climate change.

The objective of this project is to implement the proposed model PEC for predicting disease incidence based on climate data. In addition to training this predictive model, our project will delve into post hoc model explanation techniques, ensuring that high-quality explanations are provided alongside the predictions. This presents several research challenges and novel aspects for our model, including (1) the development of a conditional time-series model and (2) the generation of effective model explanations tailored for the time-series model for regression settings.

---

\*These authors contributed equally to this work

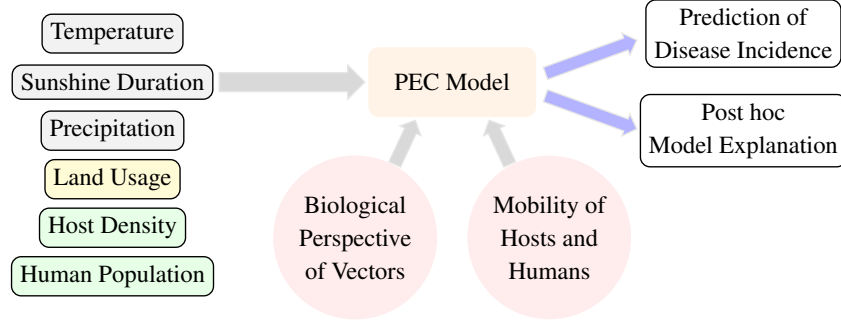


Figure 1: Overview of PEC. Input data (time-series data) is given on the left and the output on the right. Data highlighted in pink is used as the conditional input of the model.

## 2 PEC MODEL DESIGN

### 2.1 PREDICTIVE MODEL

We introduce the PEC model, a novel approach designed to leverage multi-modality data for forecasting with time-series inputs, as illustrated in Figure 1 on the left. The model incorporates climate-related variables, including temperature, sunshine duration, and precipitation, to provide a comprehensive representation of climate conditions. Furthermore, we use land usage information derived from remote sensing images. By analyzing the color of pixels of an area, we can categorize land usage based on colors. This analysis allows us to identify areas with significant vegetation, which are known to be active habitats for vectors (ticks). The third group of input is ecological information of that area such as host density and human population. PEC is predicated on contemporary epidemiological insights, including the “Biological Perspective of Vectors” and the “Mobility of Hosts,” which reflect current information rather than historical data. This allows for more tailored, relevant, and context-specific predictions (specific diseases and locations) from the model. These data will be encoded as a feature embedding to indicate a condition.

Another challenge in this project lies in designing the predictive model. Our strategy is to integrate the conditional vector with recurrent neural networks (RNNs) to develop a robust conditional RNN model. Formally, this RNN model can be defined as follows:

$$h_t = f(h_{t-1}, x_t), \quad (1)$$

$$o_t = g(h_t, c), \quad (2)$$

$$\hat{y}_t = \alpha(o_t), \quad (3)$$

where  $f(\cdot)$  and  $g(\cdot)$  denote non-linear transformations, and  $\alpha(\cdot)$  represents an activation function.  $c$  indicates the condition embedding learned from epidemiological data. To integrate the condition embedding  $c$  in the output  $o_t$ , we consider  $g(h_t, c) = \alpha(W_h h_t + W_c c + b)$ . Thus, the prediction  $\hat{y}_t$  is derived from the historical data  $x_{1:t}$  and  $c$ .

### 2.2 MODEL EXPLANATION

Explainability in time-series data has primarily focused on classifying series, aiming to determine how each feature at any timestep influences the output class. These methodologies can be categorized into three classes. Firstly, gradient-based methods, as highlighted by Baehrens et al. (2010); Sundararajan et al. (2017); Smilkov et al. (2017); Bargal et al. (2018), employ gradients to understand feature influence. Secondly, perturbation-based methods (Zeiler & Fergus, 2014; Suresh et al., 2017; Ismail et al., 2020) introduce perturbations to the input data and monitor the resultant impact on the model output, thereby identifying the most critical features for decision-making. Thirdly, attention-based methods (Choi et al., 2016; Alaa & van der Schaar, 2019) leverage training processes to prioritize significant features automatically. In the prior research, only Raman et al. (2023) specifically emphasize post-hoc model-agnostic explanations for probabilistic forecasts. Their approach is rooted in the concept of counterfactual explanation, which identifies the timesteps that are salient for probabilistic forecasts. Thus, we mark the generation of explanations for time-series data in the *regression settings* as another research challenge in this project.

As there is no consensus on whether model explanation technique can improve user’s understanding (Rong et al., 2023), we plan to first adopt the method in Raman et al. (2023) to get promising model explanations based on the automatic evaluation such as using fidelity (Rong et al., 2022; Tomsett et al., 2020) or simulated users (Chen et al., 2022). Upon receiving the results, we will conduct a human subject evaluation on which model assists users in understanding the model decisions.

### 3 PRELIMINARY RESULTS

In this section, we demonstrate our preliminary results on predicting Lyme Borrelisosis disease incidence using climate data. In this experiment, we obtained the climate data from the public local weather service agency and data includes historical temperature, sunshine duration, and precipitation. The disease incidence data was provided by BLIND-FOR-REVIEW.

We first validate the correlation between the climate data and the disease incidence using Pearson correlation  $r$ . The results are demonstrated in Figure 2. From the result, we observe a strong correlation between temperature and incidence with  $r = 0.82$ . Similarly, sunshine duration is also highly correlated to disease transmission ( $r = 0.63$ ). As these factors are attributed directly to warm weather, they are evident indicators of the activities of vectors, and thus link to the disease transmission. On the other hand, precipitation is correlated to the incidence with  $r = 0.39$ .

Given the strong correlation, we deploy a straightforward model to further verify whether the intrinsic relationship can be captured. Concretely, we use Autoregressive Integrated Moving Average (ARIMA) to perform the (short-term) prediction. The process of fitting an ARIMA model involves identifying the optimal values of  $p$  (the number of lag observations),  $d$  (the number of times that the raw observations are differenced) and  $q$  (the size of the moving average window) that best capture the autocorrelation structure of the time series. Figure 3 presents the prediction, where we observe our prediction in orange is close to the ground-truth data in blue. The promising results highlight the possibility of preciser incidence prediction using a more advanced time-series machine learning model with various information input.

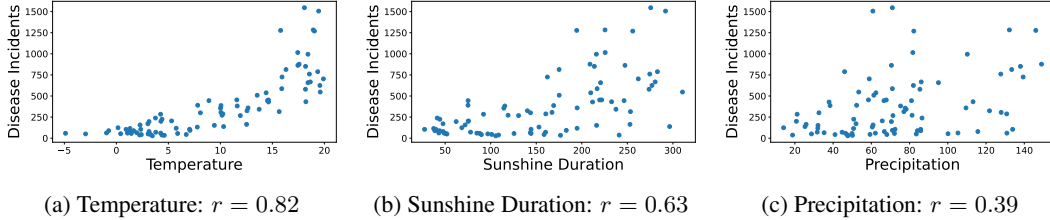


Figure 2: Pearson correlation  $r$  between the disease incidence and various climate data factors.  $r$ -value is shown below each figure.

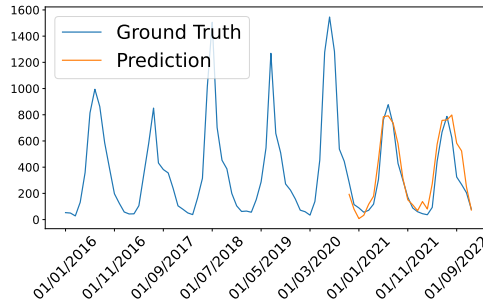


Figure 3: Prediction results from the ARIMA model using climate data. The X-axis refers to time, and the Y-axis refers to disease incidence. The blue line indicates the ground-truth data, while the orange highlights model predictions.

## REFERENCES

- Ahmed M Alaa and Mihaela van der Schaar. Attentive state-space modeling of disease progression. *Advances in neural information processing systems*, 32, 2019.
- David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11:1803–1831, 2010.
- Sarah Adel Bargal, Andrea Zunino, Donghyun Kim, Jianming Zhang, Vittorio Murino, and Stan Sclaroff. Excitation backprop for rnns. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1440–1449, 2018.
- Valerie Chen, Nari Johnson, Nicholay Topin, Gregory Plumb, and Ameet Talwalkar. Use-case-grounded simulations for explanation evaluation. *Advances in Neural Information Processing Systems*, 35:1764–1775, 2022.
- Edward Choi, Mohammad Taha Bahadori, Jimeng Sun, Joshua Kulas, Andy Schuetz, and Walter Stewart. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. *Advances in neural information processing systems*, 29, 2016.
- JS Gray, H Dautel, A Estrada-Peña, O Kahl, E Lindgren, et al. Effects of climate change on ticks and tick-borne diseases in europe. *Interdisciplinary perspectives on infectious diseases*, 2009, 2009.
- Aya Abdelsalam Ismail, Mohamed Gunady, Hector Corrada Bravo, and Soheil Feizi. Benchmarking deep learning interpretability in time series predictions. *Advances in neural information processing systems*, 33:6441–6452, 2020.
- Ding Ning, Varvara Vetrova, and Karin Bryan. Graph-based deep learning for sea surface temperature forecasts. In *ICLR 2023 Workshop on Tackling Climate Change with Machine Learning*, 2023. URL <https://www.climatechange.ai/papers/iclr2023/39>.
- Chirag Raman, Alec Nonnemaker, Amelia Villegas-Morcillo, Hayley Hung, and Marco Loog. Why did this model forecast this future? information-theoretic saliency for counterfactual explanations of probabilistic regression models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Rhea J Rocque, Caroline Beaudoin, Ruth Ndjaboue, Laura Cameron, Louann Poirier-Bergeron, Rose-Alice Poulin-Rheault, Catherine Fallon, Andrea C Tricco, and Holly O Witteman. Health effects of climate change: an overview of systematic reviews. *BMJ open*, 11(6):e046333, 2021.
- Yao Rong, Tobias Leemann, Vadim Borisov, Gjergji Kasneci, and Enkelejda Kasneci. A consistent and efficient evaluation strategy for attribution methods. In *International Conference on Machine Learning*, pp. 18770–18795. PMLR, 2022.
- Yao Rong, Tobias Leemann, Thai-Trang Nguyen, Lisa Fiedler, Peizhu Qian, Vaibhav Unhelkar, Tina Seidel, Gjergji Kasneci, and Enkelejda Kasneci. Towards human-centered explainable ai: A survey of user studies for model explanations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- David B Sandalow, Colin McCormick, Alp Kucukelbir, Julio Friedman, Trishna Nagrani, Zhiyuan Fan, Antoine M Halff, Alexandre d’Aspremont, Ruben Glatt, Elena Méndez Leal, et al. Artificial intelligence for climate change mitigation roadmap. 2023.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- Harini Suresh, Nathan Hunt, Alistair Johnson, Leo Anthony Celi, Peter Szolovits, and Marzyeh Ghassemi. Clinical intervention prediction and understanding using deep networks. *arXiv preprint arXiv:1705.08498*, 2017.

Richard Tomsett, Dan Harborne, Supriyo Chakraborty, Prudhvi Gurram, and Alun Preece. Sanity checks for saliency metrics. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 6021–6029, 2020.

Vermont. Climate change and infectious diseases, 2023. URL <https://www.healthvermont.gov/environment/climate/climate-change-and-infectious-diseases#:~:text=Warmer%20weather%20is%20one%20of,ticks%20are%20active%20each%20year.>

Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.

Cynthia Zeng and Dimitris Bertsimas. Global flood prediction: a multimodal machine learning approach. In *ICLR 2023 Workshop on Tackling Climate Change with Machine Learning*, 2023. URL <https://www.climatechange.ai/papers/iclr2023/5>.