

1 *Artificial Intelligence Predicts and Explains West Nile Virus Risks Across Europe:* 2 *Extraordinary Outbreaks Determined by Climate and Local Factors*

3
4 Albert A Gayle¹

5
6 1. Department of Public Health and Clinical Medicine, Section of Sustainable Health, Umeå University, SE-90187 Umeå, Sweden
7

8 9 **Methods**

10 11 *Data Collection.*

12
13 A literature review was conducted. The purpose of this review was to identify potential covariates
14 of West Nile virus (WNV) emergence. PubMed and Google Scholar were searched. All research
15 published between 1999 and 2019 was identified that included “West Nile” in either the title or
16 abstract. This set of research publications was furthermore filtered. Only those titles or abstracts that
17 referred to geospatial risk associations were included. Search terms for this secondary filter
18 included: “risk”, “model”, “predict”, “estimat”, “associat”, “determinant”, “factor” and “review”.
19 This list was then manually screened to exclude articles that presented neither primary research nor
20 references to the results of such. The determinants presented were then extracted and compiled.
21 Time-lag and other interaction effects were noted separately and in addition to the primary
22 determinants. Determinant specifications with respect to aggregation, normalization, and other
23 processing were also noted. This list was then stratified into eight categories based on class and
24 source of data. Spatial and temporal covariates were included among the eight feature classes.
25

26 The selected feature classes included the following broad categories: “climate”, “environmental”,
27 “economic”, “sociodemographic”, “hosts”, and “vectors”. Public data sources for each feature class
28 were identified and selected based on the availability of data that aligned with determinant
29 specifications. The literature on WNV consistently describes it as “hard to predict”. And this is
30 reflected in broad uncertainties surrounding the precise mechanistic associations surrounding even
31 the most basic, and generally accepted, determinants. In many cases, time-lagged climatic features
32 are included in models to serve as proxy for specific features relevant to vector breeding or disease
33 transmission. For example, what does it mean to associate rain that occurred 4 weeks prior to a
34 reported case today? What is the actual biotic or environmental feature that is driving increased risk
35 today: a) precipitation accumulation, b) soil moisture, c) standing water, or d) increased surface
36 reflectivity indicative of water presence (MNDWI)? In order to avoid such potentially confounding
37 effects, we expanded our inclusion criteria to include all climatic or environmental indicators
38 available from a given selected source – this included some features not present in the literature.
39

40 All data was downloaded and processed using the R platform. Spatial and temporal covariates were
41 included: one for each of the district-level regions in Europe (NUT3) as well as years 2010-2016
42 and 2018. Geometries for each of the NUTS3 regions were obtained using the “Eurostats” package
43 in R. Public sources were sought out for each of the remaining feature classes. See *Supplementary*
44 *Table 1* for a complete summary of all included features, sources, processing specifications, and
45 summary statistics. The “earthEngineGrabR”¹ package in R was used to obtain environmental and
46 climate data from Google Earth Engine. Climatic trend data was obtained directly, in preprocessed
47 form from the International Research Institute for Climate and Society (IRI) at Columbia University
48 (New York, USA)². The “Bioclimatic” feature set is a well-cited set of time-varying geoclimatic
49 metrics applied in the area of ecological niche modeling³. They are derived from standard climatic
50 variables such as precipitation and temperature, over time, for a given region. The “dismo”
51 package⁴ in R was used to calculate these. In addition, indices for intradecadal climatic systems such
52 as the “El Nino South Oscillation” (ENSO) and the “North Atlantic Oscillation” (NAO) were
53 obtained using the “rSOI”⁵ package in R. Environmental and host data was obtained directed from
54 the European Environment Agency (EEA)⁶. This included all Corine Land Cover (CLC) and Natura
55 2000 Bioregions data⁷. Geospatial species distribution data was also obtained from the EEA, by
56 merging the Natura 2000 Habitats and Species tables. This included all species, to allow for

57 potentially unreported hosts. Economic and sociodemographic data was obtained using the
58 “Eurostat” package⁸. The European Center for Disease Control (ECDC) provided mosquito
59 distribution data; this included all reported vectors in the region. West Nile virus case data was also
60 obtained from the ECDC.

61
62 To account for time-lag effects, all climatic variables were aggregated into four quarters as done in
63 similarly scoped studies⁹. This excluded the Bioclimatic variables which are defined annually.

64
65 WNV case data was quality-controlled to exclude cases likely to be subject to regional variation in
66 reporting or other biases. Excluded were 1) asymptomatic cases, 2) cases with no syndromatic
67 description, and 3) cases acquired outside the reporting district. This is consistent with previous
68 work, in which all but the most severe neuroinvasive manifestations (WNND) were excluded¹⁰.

69
70 This final data set was then divided into a training set (2010-2016) and test/validation set (2018).
71 The statistical effect of large-scale climatic systems and associated downstream covariation may
72 potentially overlap year-to-year. To mitigate, 2017 was omitted from training.

73
74 *Analysis.*

75
76 The final data set was thus used to construct an explanatory model. This included the full year’s
77 data, including climatic data from periods overlapping and extending beyond the outbreak season
78 (i.e., January - December). This model was therefore not suited for same year prospective early
79 warning, but would serve instead to retrospectively examine the factors driver the 2018 outbreak.

80
81 The official implementation for XGBoost in R¹¹ was used. Cross-validation was used to set the
82 initial iteration limits. A training data set was generated for each using data from years 2010-2016,
83 according to the specification of each model. For cross validation, a hold out subset is specified and
84 then used to evaluate the performance of a model trained with the remaining data. Here, given the
85 cyclical nature of our predictive task, the hold out set was specified to be a single year. For each
86 model, the training set was therefore split according to year and a model was trained using six of the
87 seven years. The seventh, hold out year was then used to test model performance. This process was
88 repeated seven times, using each of the included years once as the hold out set. This process
89 generates several benchmarks, one of which is an optimal number of iterations. For classification
90 tree models, iteration limitations are commonly set so as to avoid over-fitting. A second model was
91 trained using the entire seven years of training data, with the iteration limit set based on the limit
92 assessed by the cross-validation run. An additional criteria was set to stop model training after 10
93 iterations with no reduction in log-loss error. Once this second training round was complete,
94 covariates determined to be insignificant were dropped and the model was rerun. This mode of
95 feature selection is consistent with a common use case for XGBoost¹² and it is furthermore
96 recommended for this class of model¹³. In our case, covariates found to be insignificant with respect
97 to three importance criteria – marginal contribution to predictive output (“gain”), internal predictive
98 prevalence (“frequency”), and case-wise predictive relevance (“cover”) – were removed.

99
100 A final model was therefore generated that included an independently determined subset of features.
101 In a departure from established machine learning norms^{14,15}, model hyperparameters were not
102 subject to automated optimization procedures. Models such as XGBoost are notoriously sensitive to
103 hyperparameter settings¹⁶ – a problem we wished to avoid outright. Inspired by the work of Blagus
104 and Lusa (2017)¹⁷ and Brieman (2001)¹³, hyperparameters were therefore set to ensure internal
105 experimental robustness: 1) 95% was set as the internal discriminatory threshold, 2) each submodel
106 was limited to half of the total feature space, and 3) maximum tree size was set to twice the number
107 of included feature classes (i.e., 16). Consistent with guidance put forth in the seminal paper by
108 Elith J. et al (2008), step size was reduced to 0.03 (10% of the default) to ensure feature interactions
109 were robustly assessed within the ensemble.

110

111 This model was applied to the data set from 2018 to generate a probability of event. This was
112 conducted on an out-of-sample basis, in prospective time: no data used to train the model was used
113 to evaluate the model and all data corresponded to time points preceding the event.

114
115 Two sets of binary predictions (presence or absence) were derived using the standard 50%
116 classification threshold and a more discriminating 10% threshold. The choice of classification
117 threshold directly affects predictive performance and several methods exist to optimize¹⁸. However,
118 as with manual or automated tuning of hyperparameters in XGBoost, the effect of such methods on
119 out-of-sample performance is often unpredictable¹⁹. Out-of-sample performance was therefore
120 evaluated at the specified thresholds. In-sample performance was likewise estimated using the
121 2010-2016 training set. For each model, four sets of evaluations were therefore produced: 1) out-of-
122 sample at 50%, 2) out-of-sample at 10%, 3) in-sample at 50%, and 4) in-sample a 10%. Each
123 included the following (reported) metrics: 1) AUC, 2) sensitivity, 3) specificity, and 4) balanced
124 accuracy. AUC is the probability that a randomly selected positive case will be assigned a higher
125 probability than a randomly selected negative case. It indicates how well the model discriminates
126 between positive and negative classes irrespective of classification threshold. The remaining three
127 metrics are threshold-dependent and vary accordingly. Sensitivity, also referred to as true positive
128 rate, measures the probability of detecting event occurrence. Specificity, also referred to true
129 negative rate, measures the probability of correctly identifying non-occurrence. Balanced accuracy
130 is the average of the prior two metrics; it measures the probability of detection assuming negative
131 and positives classes are equally represented in the data set. In the present study, positive classes
132 represent only 0.05% of the total sample. A model that optimizes balanced accuracy irrespective of
133 threshold is therefore preferable.

134
135 The “SHapley Additive exPlanations” (SHAP) engine is an AI-driven deductive framework for
136 imputing local feature effects²⁰. The underlying algorithm is derived from a military intelligence
137 application based on cooperative game theory. This algorithm determines individual contributions
138 for group outcomes and has been mathematically proven to be optimal with respect to his task²¹.
139 The SHAP implementation applies this framework to an XGBoost model to determine the degree to
140 which individual features influence the computed event probability for a given case. The output of
141 this process is a meta-model, consisting of an effect matrix dimensionally identical to the original,
142 nominal value data (matrix): one complete set of feature effect estimations for each case. Each term
143 reflects the local change in log relative risk associated with a given feature with respect to the
144 global, baseline probability of event. This global baseline risk is computed along with the effect
145 matrix and is the same for all cases. This transformation effectively recenters the predictive output
146 at zero, implying a predictive threshold equivalent to the computed baseline probability. Values
147 greater than zero indicate increased relative risk of event and *vice versa*. This transformation
148 provides for several useful analytical modalities that are presented later and in the main text. The
149 inputs for this process are the computed XGBoost model and the original, nominal value data (set).

150
151 The SHAP engine was therefore applied to this model to generate a SHAP meta-model for 2018.
152 Predictive performance was evaluated in a manner identical to that described for the XGBoost
153 model. However, as previously described, threshold selection is implicit to the SHAP generation
154 process and was therefore omitted. The SHAP-implied classification threshold was likewise
155 computed and applied to the original XGBoost model to confirm additive and predictive
156 equivalence. For each of the SHAP model, one set of evaluations was therefore produced.

157
158 The SHAP effect matrix was visually summarized and geospatially mapped. We furthermore
159 conducted means tests ($p < 0.05$) and calculated standardized means differences to determine which
160 effects most strongly discriminated between between a) regions with positive and negative outbreak
161 risk and b) 2018 and the control period (2010-2016) with respect to regions with positive outbreak
162 risk. Supervised clustering is the application of a traditional unsupervised cluster analysis to effect
163 estimates derived from a supervised learning process. This was introduced by Lundberg et al
164 (2020)²⁰ who used it to demonstrate case-wise interaction effects. Previous studies have suggested
165 underlying non-climatic risk factors to be strong predictors of WNV risk, provided suitable climate.

166 We therefore applied a revised methodology to a subset of features corresponding with
167 environmental, host, and vector related covariates. The resulting geospatial clusters and associated
168 outbreak risk statistics were visually summarized and geospatially mapped.

169
170 The “XGBoost”²² package in R was used for all modeling, including generation of the local SHAP
171 meta-models. The “ggplot2”²³ package in R was used for visual summarization and the “tmap”²⁴
172 package in R was used for geospatial mapping. Geospatial polygons were obtained from the
173 Statistical Office of the European Commission using the R package, “eurostat”⁸. All computations
174 as well as as necessary data processing were conducted in R, using the “data.table”²⁵ package.

175 176 *Limitations*

177
178 Several *a priori* assumptions and modeling choices were required. Such choices naturally influence
179 the final model and output in various ways, not limited to overall predictive performance and
180 imputed associations with respect to individual features. Most notably, the final feature set size for
181 any given model is in no way deterministic and could very well be a function of model
182 hyperparameters, namely “tree-depth” and “eta”. These two hyperparameters effectively control the
183 degree to which the model explores causal alternatives and complexity of causation, and it is
184 plausible to presume some effect. Furthermore, the data set used to derive this model was in no way
185 comprehensive or statistically refined and outputs are sure to change given additional fine-tuning in
186 this regard. That having been said, exploring the performance of these methods with data of
187 uncertain or uneven quality was an implicit goal of this work.

188
189

1. Simplify the acquisition of remote sensing data. <https://jesjehle.github.io/earthEngineGrabR/>.
2. IRI Climate and Society Map Room. [index.html](#).
3. Kriticos, D. J., Jarošik, V. & Ota, N. Extending the suite of BIOCLIM variables: a proposed registry system and case study using principal components analysis. *Methods Ecol. Evol.* **5**, 956–960 (2014).
4. Hijmans, R. J., Phillips, S. & Elith, J. L. and J. *dismo: Species Distribution Modeling*. (2017).
5. Import Various Northern and Southern Hemisphere Climate Indices. <https://boshek.github.io/rsoi/>.
6. European Environment Agency's home page — European Environment Agency. <https://www.eea.europa.eu/>.
7. Natura 2000 - Environment - European Commission. https://ec.europa.eu/environment/nature/natura2000/index_en.htm.
8. Lahti [aut, L. *et al. eurostat: Tools for Eurostat Open Data*. (2020).
9. Keyel, A. C. *et al.* Seasonal temperatures and hydrological conditions improve the prediction of West Nile virus infection rates in Culex mosquitoes and human case counts in New York and Connecticut. *PLOS ONE* **14**, e0217854 (2019).
10. Laperriere, V., Brugger, K. & Rubel, F. Simulation of the seasonal cycles of bird, equine and human West Nile virus cases. *Prev. Vet. Med.* **98**, 99–110 (2011).
11. Chen, T. & Guestrin, C. Xgboost: A scalable tree boosting system. in *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* 785–794 (2016).
12. He, X. *et al.* Practical Lessons from Predicting Clicks on Ads at Facebook. in *Proceedings of 20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining - ADKDD'14* 1–9 (ACM Press, 2014). doi:10.1145/2648584.2648589.
13. Breiman, L. Random forests. *Mach. Learn.* **45**, 5–32 (2001).
14. Wang, Y. & Ni, X. S. A XGBoost risk model via feature selection and Bayesian hyper-parameter optimization. *ArXiv190108433 Cs Stat* (2019).
15. Xia, Y., Liu, C., Li, Y. & Liu, N. A boosted decision tree approach using Bayesian hyper-parameter optimization for credit scoring. *Expert Syst. Appl.* **78**, 225–241 (2017).
16. Mayr, A., Binder, H., Gefeller, O. & Schmid, M. The evolution of boosting algorithms. From machine learning to statistical modelling. *Methods Inf. Med.* **53**, 419 (2014).
17. Blagus, R. & Lusa, L. Gradient boosting for high-dimensional prediction of rare events. *Comput. Stat. Data Anal.* **113**, 19–37 (2017).
18. Cortes, C. & Mohri, M. AUC optimization vs. error rate minimization. in *Advances in neural information processing systems* 313–320 (2004).
19. Yang, K., Yu, R., Wang, X., Quddus, M. & Xue, L. How to determine an optimal threshold to classify real-time crash-prone traffic conditions? *Accid. Anal. Prev.* **117**, 250–261 (2018).

20. Lundberg, S. M. *et al.* From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* **2**, 56–67 (2020).
21. Lundberg, S. M., Erion, G. G. & Lee, S.-I. Consistent Individualized Feature Attribution for Tree Ensembles. *ArXiv180203888 Cs Stat* (2019).
22. XGBoost. <https://xgboost.ai/>.
23. Wickham, H. *et al.* *ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics*. (2020).
24. Tennekes, M. *et al.* *tmap: Thematic Maps*. (2020).
25. Dowle, M. *et al.* *data.table: Extension of 'data.frame'*. (2019).