

1 Predictive mapping of mosquito distribution based on environmental and anthropogenic factors in  
2 Taita Hills, Kenya

3  
4 Ruut Uusitalo<sup>1,2,3</sup>, Mika Siljander<sup>1</sup>, C. Lorna Culverwell<sup>2,3,4</sup>, Noah C. Mutai<sup>5</sup>, Kristian M. Forbes<sup>2</sup>,  
5 Olli Vapalahti<sup>2,3,6</sup> Petri K.E. Pellikka<sup>1,7,8</sup>

6  
7 <sup>1</sup>Department of Geosciences and Geography, P.O. Box 64, FI-00014 University of Helsinki,  
8 Finland; mika.siljander@helsinki.fi; petri.pellikka@helsinki.fi

9 <sup>2</sup>Department of Virology, Haartmaninkatu 3, P.O. Box 21, FI-00014 University of Helsinki,  
10 Finland; lorna.culverwell@helsinki.fi; kristian.forbes@helsinki.fi; olli.vapalahti@helsinki.fi

11 <sup>3</sup>Department of Veterinary Biosciences, Agnes Sjöberginkatu 2, P.O. Box 66, FI-00014 University  
12 of Helsinki, Finland

13 <sup>4</sup>Department of Life Sciences, Natural History Museum, Cromwell Road, London, SW5 7BD

14 <sup>5</sup>Department of Mathematics and Informatics, Taita Taveta University, P. O. Box 635-80300, Voi,  
15 Kenya

16 <sup>6</sup>Virology and Immunology, HUSLAB, Helsinki University Hospital

17 <sup>7</sup>Helsinki Institute of Sustainability Science, University of Helsinki, Finland

18 <sup>8</sup>Institute for Atmospheric and Earth System Research, University of Helsinki, Finland

19  
20  
21 Corresponding author:

22 Ruut Uusitalo

23 University of Helsinki, Department of Geosciences and Geography, P.O. Box 64, FI-00014  
24 ruut.uusitalo@helsinki.fi; Tel.: +358 (0)50 5911280

25  
26  
27 Abstract

28 Mosquitoes are vectors for numerous pathogens, which are collectively responsible for millions of  
29 human deaths each year. As such, it is vital to be able to accurately predict their distributions,  
30 particularly in areas where species composition is unknown. Species distribution modeling was used  
31 to determine the relationship between environmental, anthropogenic and distance factors on the  
32 occurrence of two mosquito genera, *Culex* Linnaeus and *Stegomyia* Theobald (syn. *Aedes*), in the  
33 Taita Hills, southeastern Kenya. This study aims to test whether any of the statistical prediction  
34 models produced by the *Biomod2* package in R can reliably estimate the distributions of mosquitoes

35 in these genera in the Taita Hills; and to examine which factors best explain their presence. Mosquito  
36 collections were collected from 122 locations between January–March 2016 along transects  
37 throughout the Taita Hills. Environmental-, anthropogenic- and distance-based geospatial data were  
38 acquired from the Taita Hills geo-database, satellite- and aerial imagery and processed in GIS  
39 software. The Biomod2 package in R, intended for ensemble forecasting of species distributions, was  
40 used to generate predictive models. Slope, human population density, normalized difference  
41 vegetation index, distance to roads and elevation best estimated *Culex* distributions by a generalized  
42 additive model with an area under the curve (AUC) value of 0.791. Mean radiation, human population  
43 density, normalized difference vegetation index, distance to roads and mean temperature resulted in  
44 the highest AUC (0.708) value in a random forest model for *Stegomyia* distributions. We conclude  
45 that in the process towards more detailed species-level maps, with our study results, general  
46 assumptions can be made about the distribution areas of *Culex* and *Stegomyia* mosquitoes in the Taita  
47 Hills and the factors which influence their distribution.

48

49 Keywords: species distribution modeling; vector-borne disease; GIS; predictive mapping; mosquito,  
50 biomod2

51

## 52 1. Introduction

53

54 Environmental and anthropogenic disturbances such as climate change, urbanization and  
55 deforestation are crucial factors in the distribution of pathogen vectors and the emergence of diseases  
56 that they transmit (Crowl *et al.*, 2008). This is already evident as the global rise in human population  
57 densities, and correlated mean temperatures leading to land-use changes, has created new suitable  
58 breeding sites for mosquitoes (Roiz *et al.*, 2011; Campbell *et al.*, 2015). Geographic information  
59 system (GIS) and species distribution modeling (SDM) approaches have been used to understand  
60 these connections and their implications for the spread of invasive species, particularly with regard  
61 to northern regions (Neteler *et al.*, 2011). Other important areas, however, remain unexplored,  
62 including biodiversity rich but increasingly fragmented parts of Africa, where infectious diseases can  
63 be highly prevalent but access to health care is limited (Bhutta *et al.*, 2014). As a result, tools such as  
64 GIS and SDM that can estimate and predict infectious disease risks have the potential to optimize the  
65 use of limited resources and improve public health outcomes.

66

67 Due to their ability to act as vectors for a suite of pathogens, mosquitoes (Diptera, Culicidae) are  
68 among the most economically and socially important taxa on the planet (WHO, 2017). Two notable

69 genera are *Culex* Linnaeus and *Stegomyia* Theobald (following the classification of Reinert *et al.*,  
70 (2009)). *Culex* is a large genus of mosquitoes, with 769 species in 26 subgenera (MTI, 2017),  
71 accounting for 21.6% of all mosquito species worldwide. Species of *Culex* have an almost worldwide  
72 distribution from the tropics to cool temperate regions, but do not extend into extreme northern  
73 latitudes (MTI, 2017). While species from other subgenera are known to vector pathogens, subgenus  
74 *Culex* (*Culex*) contains many of the significant human vector species, including those demonstrated  
75 to transmit West Nile fever virus, Rift Valley fever virus and Japanese Encephalitis virus (MTI, 2017;  
76 WHO, 2017). *Stegomyia* is a moderately sized mosquito genus comprising 128 species, and  
77 distributed in the Afrotropical, Oriental and Australasian regions (MTI, 2017). Following human  
78 dispersal, at least two species, *St. aegypti* and *St. albopicta*, are also present in Neotropical, Nearctic  
79 and Palearctic regions (MTI, 2017; Paupy, *et al.*, 2009). Member species are competent vectors of  
80 yellow fever virus (Huang, 1986), dengue virus serotypes 1-4, Zika virus, and Chikungunya virus,  
81 among others (Huang, 1990; MTI, 2017).

82

83 SDM helps to understand the relationship between wildlife and their environments, to estimate their  
84 distributions where empirical data is missing, and to forecast range expansions in the face of  
85 environmental change. SDM is used in decision-making for a range of global challenges, such as  
86 informing protected areas for wildlife conservation (Guisan *et al.*, 2013), the distribution potential for  
87 invasive species (Václavík & Meentemeyer, 2009), and investigating the effects of climate change  
88 on species that carry pathogens of public health concern such as mosquitoes (Dukes *et al.*, 2009).  
89 Species distribution modeling, therefore, provides a means to understand and predict vector responses  
90 to changing climate patterns in dynamic and fragmented environments.

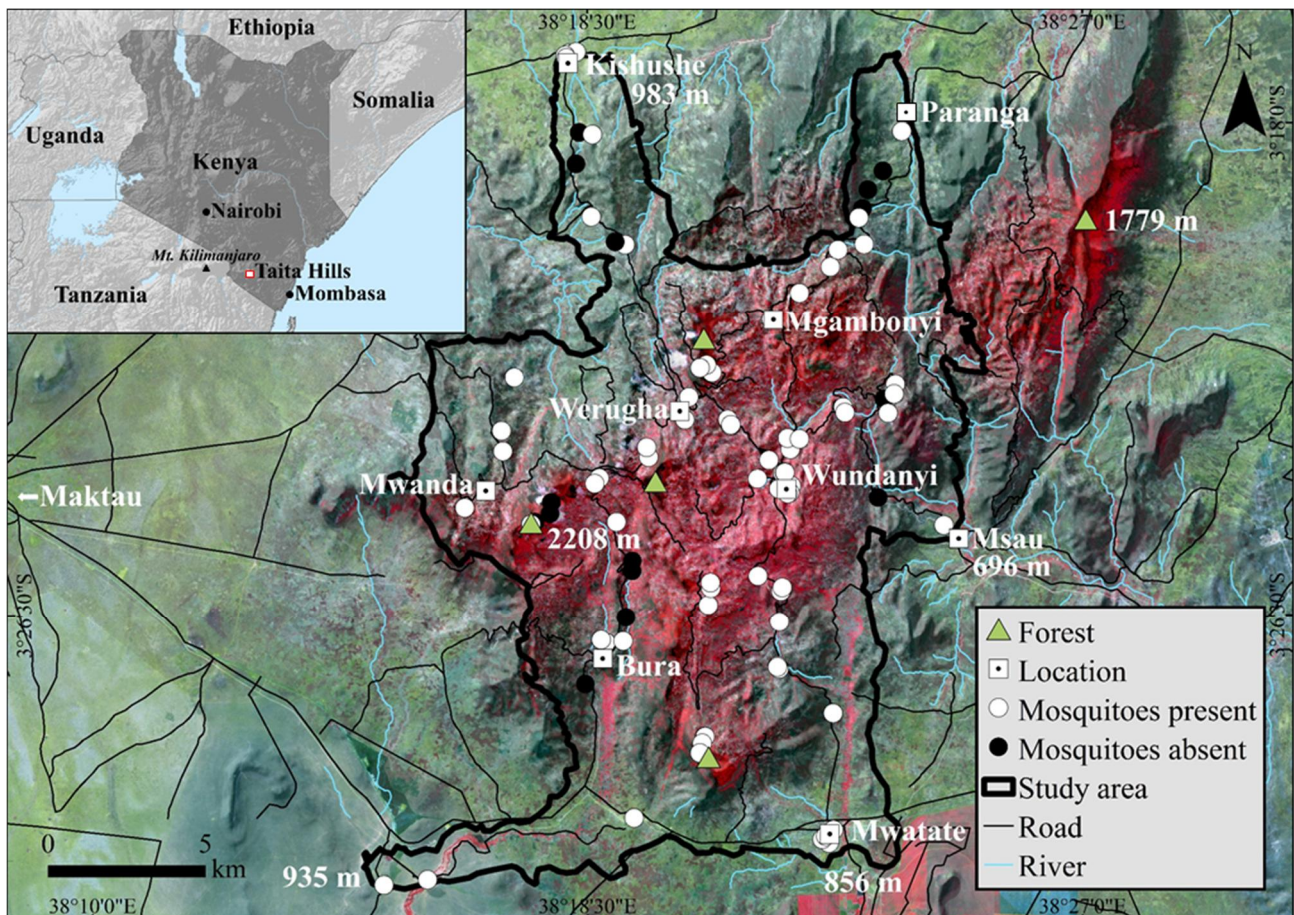
91

92 The purpose of this study is to test whether any of the statistical prediction models produced by the  
93 Biomod2 package in R, can reliably estimate the distributions of mosquitoes in genera *Culex* and  
94 *Stegomyia* in the Taita Hills; to examine which factors best explain their presence; and to create maps  
95 of their predicted occurrence. The target area, Taita Hills in rural south-eastern Kenya, is an  
96 ecologically diverse area with little previous mosquito research, strong variability in rainfall and a  
97 rapidly growing human population. While machine-learning techniques such as maximum entropy  
98 modeling is often used to understand mosquito distributions (Mughini-Gras *et al.*, 2014; Sallam *et*  
99 *al.*, 2016), the employed R package, the Biomod2 (Thuiller *et al.*, 2016), has important benefits for  
100 estimating mosquito distributions. This is due to its ability to provide improved simulations across  
101 initial conditions, test multiple model classes and parameterizations, and include unlimited number  
102 of boundary conditions (Thuiller *et al.*, 2009).

103 2. Methods

104 2.1 Study area and mosquito collections

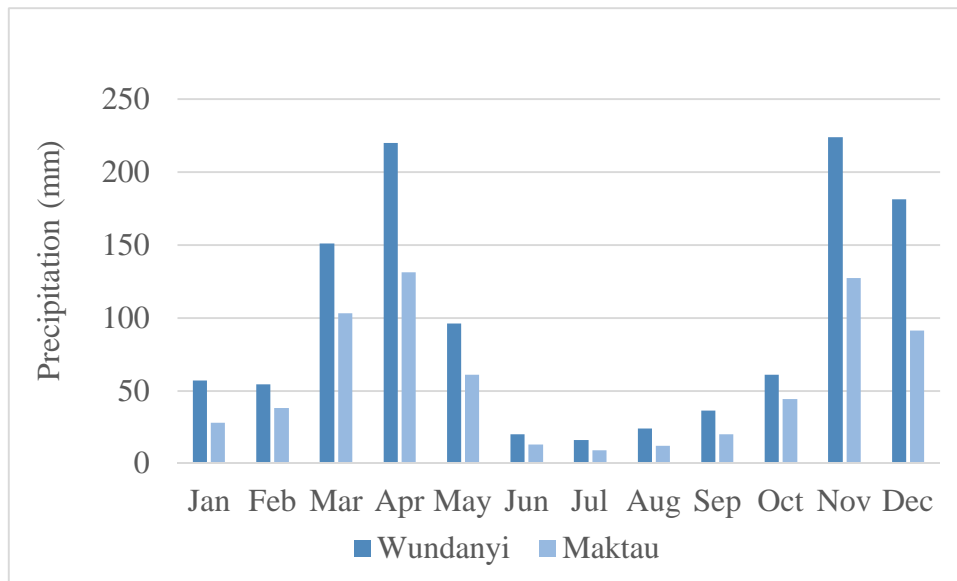
105 The Taita Hills (03°20'S, 38°15'E) in Taita–Taveta county, southeastern Kenya (Fig. 1.) is a range  
106 of peaks varying in altitude from 600 to 2200 m above sea level (a.s.l.) (Erdogan *et al.*, 2011). The  
107 study area of 286 km<sup>2</sup> ranges from the Taita Hills at 1800 m a.s.l. to the surrounding lowlands at  
108 900 m a.s.l., as outlined in Fig..



109  
110 Fig. 1. The Taita Hills and sampling locations in southeastern Kenya. Sampling was conducted  
111 along the main roads, from the surrounding lowlands to mountain peaks.

112  
113 Kenya has two annual wet and dry seasons, with the total rainfall varying across the country (Kaplan  
114 *et al.*, 1976; Fig. 2.). The Taita Hills region receives on average 1330–1910 mm/annually, which  
115 contributes to the formation of many suitable mosquito breeding habitats (Erdogan *et al.*, 2011).

116



117

118 Fig. 2. Monthly average rainfall characteristics for Wundanyi town in the Taita Hills and Maktau  
 119 village in the lowland plains. Weather data was collected between 1982 and 2012, with a resolution  
 120 of 30 arc seconds (Climate-Data.org, 2018).

121

122 Mosquito collections were attempted at over 122 locations across the region between January and  
 123 March 2016. This figure includes ‘sub-locations’, which fell within 500 m of the primary target  
 124 location, and repeat collections from the same locations across the study time (Fig. 1.). Of these 122  
 125 attempts, mosquitoes were present in 107 collections and absent from 15. Mosquito collections started  
 126 after the rainy season on January 25<sup>th</sup>, 2016 and continued until March 16<sup>th</sup>, 2016. Collections did not  
 127 follow transects up to February 12<sup>th</sup>, since the goal at this time was to focus on human-biting species  
 128 around human dwellings. From February 13<sup>th</sup> until the end of the collection period a stratified  
 129 sampling scheme was applied based on the road network. Each main road from the lowlands to the  
 130 highest reachable locations was used to collect mosquitoes on 100 m elevation intervals.

131

132 Collection locations and sub-locations were mainly around human dwellings close to roads, but  
 133 collections were also made in forest fragments and in croplands. Adults were collected using  
 134 commercially available Prokopacks (The John W. Hock Company, Gainesville, USA) or by using  
 135 CDC Miniature Light Traps (The John W. Hock Company, Gainesville, USA). Immature life stages  
 136 (larvae, pupae and eggs) were collected from stagnant water using a 1 litre plastic dipper, a fine  
 137 meshed aquarium net or a turkey baster. All water sources were considered, including septic tanks,  
 138 discarded tires, tree cavities, leaf axils and other items filled with rain water, such as discarded  
 139 artificial containers or water tanks. At each collection location, the GPS point was recorded using a  
 140 Garmin Map 64S handset. Collections were initially sorted in a field laboratory and stored in ethanol

141 or RNA-later to preserve the RNA and DNA for future studies. Genus-level identifications were made  
 142 in Finland using suitable identification keys (Service, 1991; Huang, 2000; Huang, 2004).

143  
 144 2.2 Explanatory environmental, anthropogenic and distance variables

145 The selection of explanatory variables to predict mosquito distributions was based on those used in  
 146 existing literature. Environmental and anthropogenic data for the Taita Hills region were obtained  
 147 directly from satellite imagery and airborne digital data, or by deriving them from the satellite  
 148 imagery in ArcGIS (ESRI, Redlands, CA, USA, v.10.3.1) (Table 1). Distance variables were created  
 149 from an existing Taita Hills geo-database. For the distance to buildings calculations, existing airborne  
 150 laser scanning (ALS) data was used. In this data set, buildings were classified from point cloud data  
 151 that were derived from 2014-2015 Taita Hills ALS flights (Adhikari *et al.*, 2017). LAStools software  
 152 (rapidlasso GmbH) was used to classify ground, building and vegetation returns.

153  
 154 Normalized difference vegetation index (NDVI) is a measure of plant greenness, ranging between  
 155 -1 to 1 (Tucker, 1979), and was used to approximate land cover type. Mean precipitation, mean  
 156 temperature and mean relative humidity from January - March was calculated for each collection  
 157 location in ArcGIS. Mean monthly precipitation was derived from long-term precipitation grids  
 158 between 1987 and 2005 (Hutchinson, 1991; Erdogan *et al.*, 2011). A digital elevation model (DEM)  
 159 was used to calculate the elevation, slope and mean monthly solar radiation in ArcGIS, and the values  
 160 were extracted to each mosquito collection location.

161  
 162 Table 1. Value ranges and sources for explanatory variables used in species distribution modeling.

Environmental, anthropogenic or distance factor	Min.	Max.	Avg.	Data source
Distance to houses (m)	0	1270	52	Distance to houses was determined from the building data computed from airborne laser scanning (ALS) data (Adhikari <i>et al.</i> , 2017), modified by digitizing more houses in the study area using QGIS, and calculating Mean Euclidean distance to houses.
Elevation (m)	694	2079	1330	Mean elevation was derived from digital elevation model (DEM). DEM was obtained from scanned Survey of Kenya 1:50 000 scale topographic map from 1991, in which a 20 m planimetric resolution DEM was interpolated from 50 feet interval contours (Clark & Pellikka., 2005).
Distance to roads (m)	0	927	127	Road data was digitized in 2004 from Survey of Kenya 1:50 000 scale topographic map from 1991.

---

				(Broberg <i>et al.</i> , 2004). Modified by calculating Mean Euclidean distance to roads using ArcGIS.
Mean precipitation (mm)	20	113	47	Mean precipitation was obtained from long-term mean precipitation grids, which were interpolated on to a 20 m resolution grid from monthly available meteorological data and surrounding areas between 1987 and 2005, obtained from Kenya Meteorological Department using ANUSPLIN (Australian National University Splines) software (Hutchinson, 1991; Erdogan <i>et al.</i> , 2011).
Mean radiation (kWh/m <sup>2</sup> )	176	228	216	Mean radiation was received on a given surface area in a given time (kWh/m <sup>2</sup> ), was calculated from the DEM using ArcGIS Area Solar Radiation tool (Esri, Redlands, CA, USA).
Mean relative humidity (%)	71	94	77	Mean relative humidity data was created based on data logger observations between April 2013 and May 2014 (Virtanen 2015).
Mean temperature (C°)	15	25	21	Mean temperature data was created based on data logger observations between April 2013 and May 2014 (Virtanen 2015).
Slope (°)	0	43	11	Slope degree was derived from DEM using Slope function in ArcGIS.
Human population density (persons/km <sup>2</sup> )	0	9090	1260	Human population density was estimated from ALS buildings data and from a random non-stratified household survey (n=100) carried out in October 2006, which found an average of 6 persons per dwelling in the area (Siljander <i>et al.</i> , 2011). The data was modified by digitizing more houses in the study area using QGIS.
Normalized Difference Vegetation Index (NDVI)	-0.4	0.2	-0.2	NVDI was derived from a Sentinel-2A MSI Level-1C satellite image from 8 October 2016, downloaded from the Sentinel's Scientific DataHub (ESA, 2015).

---

163

164

### 165 2.3. Data analysis and modeling

166 All the GIS datasets including environmental and other attributes were set to the same spatial extent,  
 167 geographic coordinate system (WGS 1984 UTM Zone 37S) and resolutions (20 m x 20 m). The  
 168 compiled data consisted of presence-absence data and explanatory variables for each sampling  
 169 location. We predicted the distribution of *Culex* and *Stegomyia* mosquitoes using the Biomod2-  
 170 platform (version 3.3–7) in R software (version 3.0.3; R Development Core Team 2013). Prior to  
 171 prediction modelling, multicollinearity of the variables was tested using Pearson correlation

172 coefficients in R, which resulted in the exclusion of four variables. Models were processed mainly  
173 using the default settings of the Biomod2. The following eight predictive modeling techniques were  
174 used: generalized linear models (GLM), generalized additive models (GAM), classification tree  
175 analysis (CTA), artificial neural networks (ANNs), multivariate adaptive regression splines (MARS),  
176 general boosting method (GBM), random forest (RF) and maximum entropy (Maxent). Flexible  
177 discriminant analysis (FDA) and surface range envelope (SRE) were excluded due to their  
178 methodological weaknesses (Hastie *et al.*, 1994).

179

180 The original data set of 122 collection locations was randomly divided into model training (70%, n=  
181 85) and model evaluation sets (30%, n=37) (split-sample approach: Guisan & Zimmermann 2000).  
182 Models were built using the training set, and the models were validated using the model evaluation  
183 set. The area under the curve (AUC) of a receiver operating characteristic (ROC) plot was produced  
184 based on each model to estimate the predictive power of the model (Fielding & Bell, 1997), to assess  
185 the agreement between the presence-absence records and the predictions. The model types with  
186 highest AUC values ( $0.7 < \text{AUC} < 1.0$ ) and statistically significant variables ( $p \leq 0.05$ ) were used to  
187 conduct the predictive modeling (Drew *et al.* 2011). In order to receive the highest possible prediction  
188 accuracy values, we compared the different sets and orders of predictor variables. Depending on their  
189 strength of relationship to the response variable, we selected the best combination of predictors to  
190 include in a model. We compared the variable contributions in the models to define the most powerful  
191 variables, and their relative magnitude. Spatial autocorrelation (SAC) of the predictor variables was  
192 measured using Moran's Index (Moran, 1950).

193

### 194 3. Results

195 A total of 3130 mosquitoes were collected from 107 locations across the Taita Hills region. The  
196 majority (~2600 mosquitoes from 73 locations), belonged to genus *Culex* and included at least three  
197 subgenera: *Culex* (*Culex*), *Cx.* (*Culiciomyia*) and *Cx.* (*Eumelanomyia*). *Culex* mosquitoes were  
198 common across the whole study area, including around human settlements and in forests. *Stegomyia*  
199 mosquitoes were the second most abundant genus in the collections, (~180 individuals from 28  
200 locations). Larvae of both genera were found in water tanks and small ponds, and adults were  
201 collected from houses along the roadsides and in villages. *Stegomyia aegypti* larvae were especially  
202 common in water tanks in the villages of Paranga and Kishushe in the Taita Hills north, and also in



203 car tires in Mwatate village. Both genera were present in lowland and upland areas, including  
204 elevations up to 1900 m.

205

206 The environmental predictors, elevation, mean precipitation, mean temperature and mean relative  
207 humidity, were positively or negatively correlated with each other ( $r \geq 0.9$  or  $r \leq -0.9$ ), as were slope  
208 and mean radiation ( $r \geq -0.75$ ); therefore, only one of each variable group was retained in the final  
209 model. Normalized difference vegetation index, distance to roads, distance to houses, and human  
210 population density did not result in high correlations, and were retained in the model process. The  
211 combination of non-correlating variables, which resulted in highest evaluation values, was run to  
212 obtain reliable estimations for *Culex* and *Stegomyia* (Appendix 1). Mean temperature, mean radiation,  
213 elevation, slope and NDVI were influential factors in the models (Appendix 2). Human population  
214 density had the greatest effect on the distribution of both *Culex* and *Stegomyia* in all models, but its  
215 importance in each model varied (Appendix 2). Moderately explaining variables - NDVI, slope,  
216 distance to roads and elevation - were not ranked consistently in the models. The model that best  
217 explained *Culex* mosquito distribution included the variables slope, population density, NDVI,  
218 distance to roads, and elevation. Mean radiation, NDVI, human population density, distance to roads,  
219 and mean temperature were most influential in explaining the distribution of *Stegomyia* mosquitoes.

220

221 Spatial autocorrelation (SAC) of the predictor variables was measured using Moran's Index (Moran,  
222 1950) (Appendix 3). For *Culex*, human population density and elevation were highly spatially  
223 autocorrelated (Moran's  $I \geq 0.8$ ) for short distances ( $p < 0.05$ ) but not at longer distances. Human  
224 population density and mean temperature were highly spatially autocorrelated at short distances, but  
225 not at longer distances for *Stegomyia*. The models with highest AUC value and significant p-values  
226 for *Culex*, were GAM (AUC=0.791) and MARS (AUC=0.809) (Table 2). Altogether, six of the eight  
227 models provided reliable estimates for *Culex* mosquito distribution (AUC>0.7). For *Stegomyia*, only  
228 two of the eight models resulted in evaluation values considered reliable (AUC>0.7); GBM  
229 (AUC=0.708) and RF models (AUC=0.708).

230

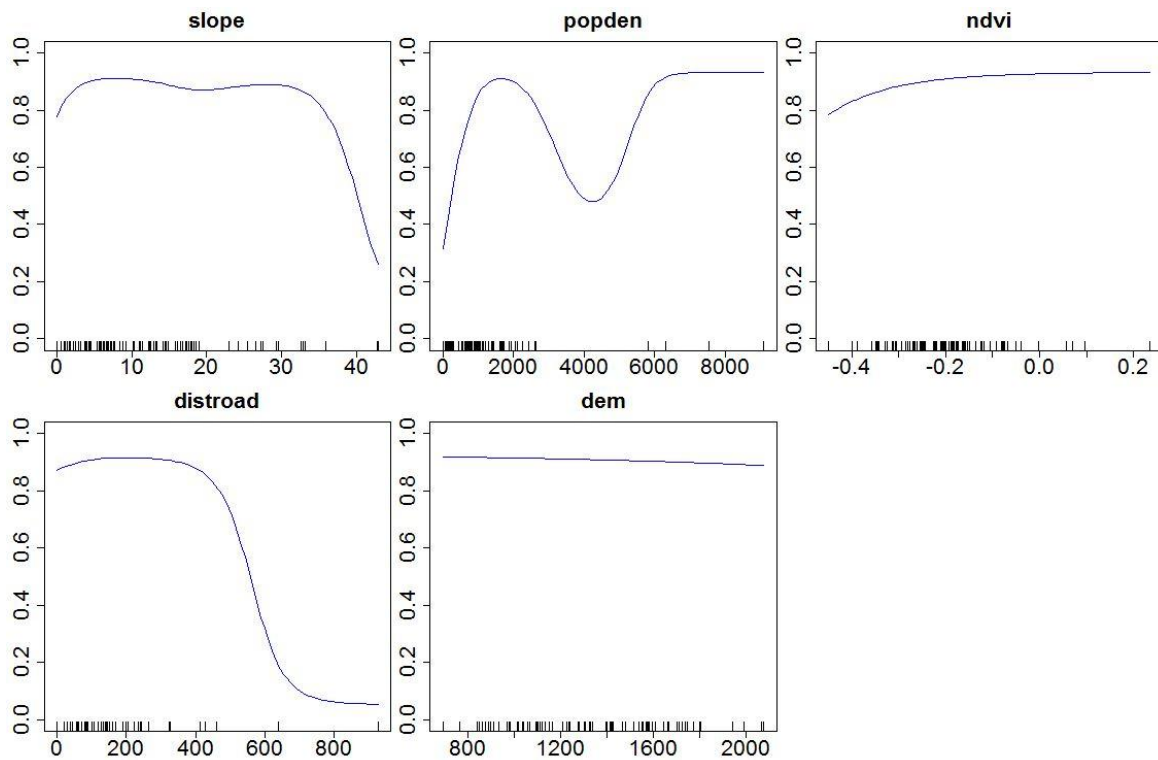
231 Table 2. Area under the curve (AUC) values of all models used to estimate the presence of *Culex*  
232 and *Stegomyia* mosquitoes.

<i>Culex</i>		<i>Stegomyia</i>	
Model	AUC	Model	AUC
GLM	0.730	GLM	-

GAM	0.791	GAM	0.643
GBM	0.750	GBM	0.708
CTA	0.620	CTA	0.616
ANN	0.764	ANN	0.612
MARS	0.806	MARS	-
RF	0.729	RF	0.708
Maxent	0.585	Maxent	0.690

233

234 For the predictive models (described below), we focus on the GAM model for *Culex* and the RF  
 235 model for *Stegomyia* due to their ability to explain the data reliably and produce the best explaining  
 236 variable contributions for this study. According to the GAM model, the probability of *Culex* presence  
 237 was high ( $\geq 80\%$ ) in locations with NDVI- values ranging from -0.4 to 0.2 (Fig. 3.). The likelihood of  
 238 *Culex* presence remained high when roads were within 500 m, and at elevations between 800 and  
 239 2000 m, in locations with moderate slope angles ( $0^\circ$ – $35^\circ$ ), and with population densities of 500–6000  
 240 people/ km<sup>2</sup>.



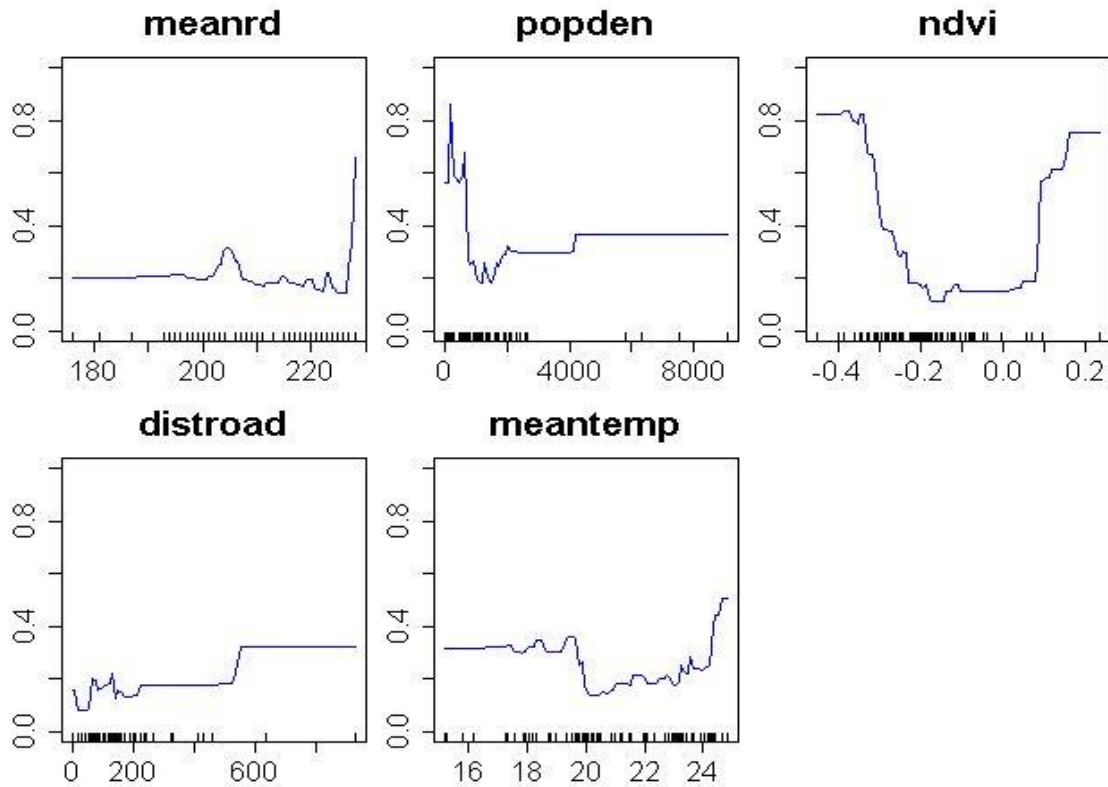
241

242 Fig. 3. Response curves for *Culex* mosquito estimations based on the GAM model. The Y-axis  
 243 represents the probability of presence, and the black marks along the x-axis represent empirical  
 244 observations.

245

246 According to the RF model, the presence of *Stegomyia* was highest in locations with solar radiation  
 247 levels  $\geq 230$  kWh/m<sup>2</sup> (Fig. 4.). They favored temperatures between  $15^\circ\text{C}$  and  $20^\circ\text{C}$ , and over  $23^\circ\text{C}$ ,

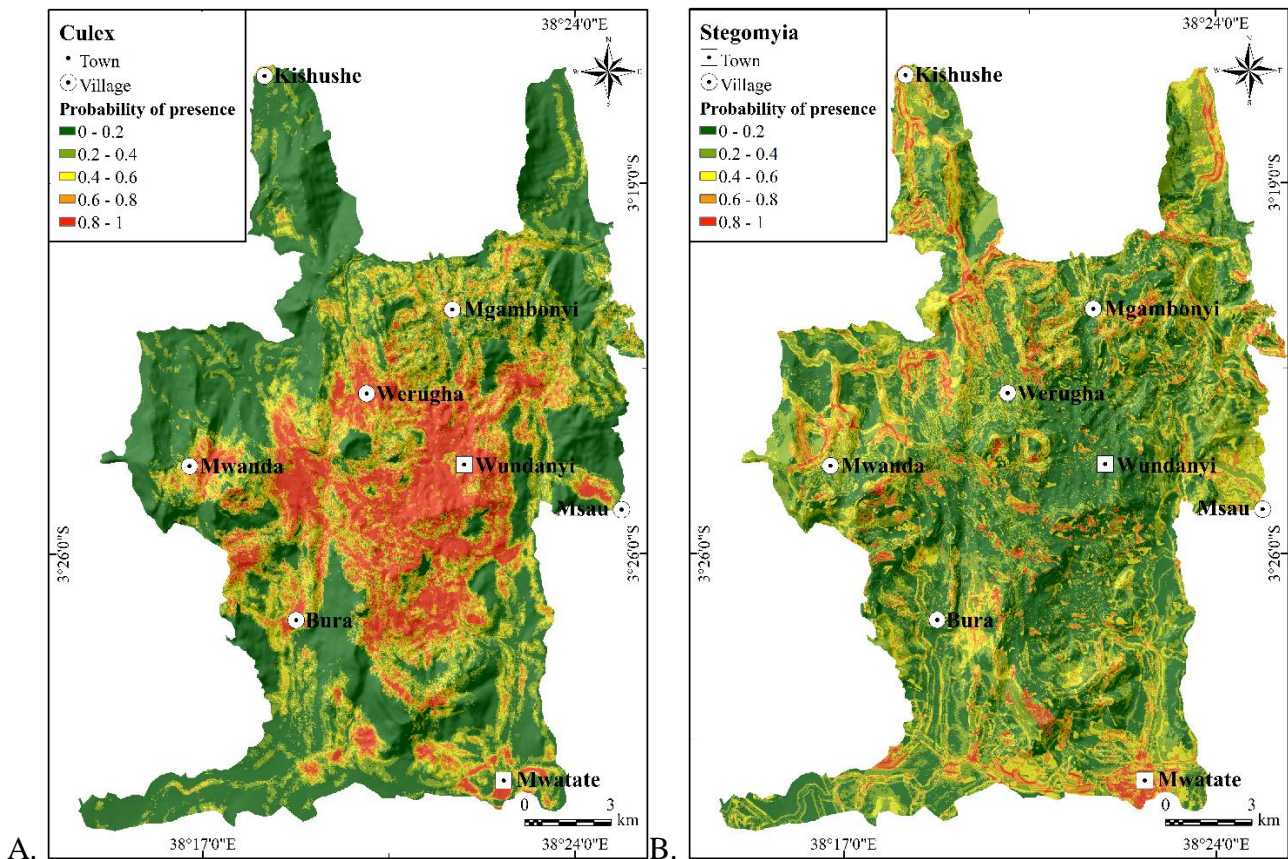
248 and locations with poor ( $NDVI \leq -0.2$ ) or moderate vegetation ( $NDVI \geq 0.2$ ). All distances between 0  
249 and 600 m from roads resulted in moderate probabilities of *Stegomyia* presence.



250  
251 Fig. 4. The response curves of predictors for *Stegomyia* mosquito estimations in RF model. The Y-  
252 axis represent the probability of presence, and the black ticks along the x-axis represent empirical  
253 observations.

254  
255 The prediction model risk maps (Figs. 5A and 5B) show that village areas and forests, including  
256 elevations above 1500 meters in the Taita Hills, were most suitable for *Culex* mosquitoes (80–100%),  
257 whereas the likelihood of their presence was reduced in lowland plateaus, apart from Mwatate village  
258 in the South (0–20%). Conversely, the likelihood of *Stegomyia* mosquitoes was high in villages on  
259 the plateaus and close to roads (Fig. 5.B.). Areas with high *Stegomyia* mosquito probability were  
260 sporadically dispersed.

262



263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

Fig. 5. Predictive risk maps for *Culex* and *Stegomyia* mosquito presence over the study area in the Taita Hills. A. The probability of *Culex* presence was highest in the central part based on the GAM model. B. The probability of *Stegomyia* presence was highest in plateau villages and fragmented areas based on a RF model.

#### 4. Discussion

We show that the Biomod2 ensemble platform in R, was able to predict the distribution of *Culex* and *Stegomyia* mosquitoes in the Taita Hills in Kenya. These two genera were found to overlap in geographical range, but also to be common in different areas based on environmental and anthropogenic factors. *Culex* and *Stegomyia* mosquitoes contain member species that are well-known vectors of significant human pathogens (Lounibos 2002). This approach can be applied to better understand how these vectors, and the infections they carry, vary across space and time, and later, this may be extrapolated to species-level studies and to other understudied areas in rural Africa.

The majority of the models utilized herein produced high AUC values for *Culex* mosquitoes, indicating reliable estimations. We concentrated on the GAM model because it resulted in high AUC values with relevant variable contributions having an automated approach to identify non-linear

280 relationships, and it may result in more accurate predictions than other model types (Yee *et al.*, 1991).  
281 For *Stegomyia* estimations, we focused on the RF model because of its advantages over other models,  
282 including high classification accuracy and the ability to model interactions among predictors (Cutler  
283 *et al.*, 2007). However, some uncertainties exist concerning *Stegomyia* estimations due to a higher  
284 number of absences (94 collections) than presences (28 collections) in the sampling data.  
285 Environmental, anthropogenic and distance variables were spatially autocorrelated to some extent,  
286 which may, among other things, affect the precision of coefficients (Diniz- Filho *et al.*, 2003).

287

288 In our study, the spatial resolution of environmental data was 20 meters. This resolution may have  
289 produced some bias to the study results because mosquitoes are also breeding e.g. in tree holes and  
290 water tanks in different microclimates. The lack of very high-resolution environmental data (such as  
291 the resolutions of 2 meters and 5 meters) is a general problem and not limited just to the Taita Hills.  
292 In addition, increasing the resolution of environmental data by deriving them from a high-resolution  
293 environmental data does not always improve the species distribution models (Pradervand *et al.*,  
294 2014). Instead, very high-resolution environmental predictors should be produced by taking into  
295 account more local field measurements such as fine environmental mapping or in-situ measurements  
296 (Pradervand *et al.*, 2014; Lembrechts *et al.*, 2018). The potential solution for the lack of microclimate  
297 data could be the use of unmanned aerial vehicles (UAVs) in order to map the sampling areas in real  
298 time, and to produce higher resolution environmental data of even 3 centimeters (Kimberly *et al.*,  
299 2014; Anderson *et al.*, 2013). Furthermore, there certainly exist other explanatory factors affecting  
300 the mosquito distributions which were not included in the modelling process including landscape  
301 fragmentation indicators. Fragmentation variables, such as distance from forest patch, patch size,  
302 distance from patch edge and the landscape metric of PPU (patches per unit) have resulted significant  
303 statistical relations with mosquito distributions (Richman *et al.*, 2018; Reiter & LaPointe 2007). For  
304 example, the use of FRAGSTATS- a spatial pattern analysis program may add value for our future  
305 mosquito habitat studies (McGarical *et al.*, 2002).

306

307 Consistent with previous research from other global locations (Sallam *et al.*, 2017; Ding *et al.*, 2017),  
308 environmental and anthropogenic variables were important determinants of *Culex* and *Stegomyia*  
309 mosquito distribution in the Taita Hills. For example the response curve for slope (Fig 3.) shows high  
310 probability for *Culex* presence for slopes less than 35°. In the Taita Hills, dwelling units are mainly  
311 absent on slopes greater than 35° (Siljander *et al.*, 2011). High probability for *Culex* presence were  
312 found with high NDVI values caused by strongly reflecting orchard trees and croplands adjacent to  
313 dwelling units in the Taita Hills. Specifically, we found that locations with high population densities,

314 short distances to roads, and gentle slope angles positively influenced *Culex* presence. Residential  
315 and urban areas are often recognized as important for *Culex* distributions (Reiter & LaPointe, 2007;  
316 Conley *et al.*, 2014). Rich and poor vegetation and all elevations were suitable for *Culex* mosquito  
317 presence, strengthening previous findings of its widespread distribution (MTI, 2017). However, we  
318 cannot rule out some effects of potential sample biases (e.g. time of the day and traps involved) which  
319 may have affected the mosquito spectrum captured. Sampling bias may also have been the reason for  
320 an unexpected response curve for population density (*popden*) at values from ca. 4000 to 5000 (Fig.  
321 3).

322

323 Our study suggests, that *Stegomyia* mosquitoes prefer locations with lower human population  
324 densities, higher temperatures and solar radiation, and poor or moderate vegetation. This finding  
325 somewhat contradicts the notion that *St. aegypti* distribution is linked to growing human population  
326 (Fatima *et al.*, 2016), but is consistent with historic records of the variety of *Stegomyia* species  
327 suggesting that most members of the *Stegomyia* genus are forest dwelling, with only some adapted to  
328 breeding close to human habitation (Powell & Tabachnick 2013). Since our collections included  
329 several different *Stegomyia* species - with *St. aegypti* being the only species collected from around  
330 human habitation, and other species being restricted to the forest - the model findings are consistent  
331 with these historical reports (Powell & Tabachnick 2013). *Stegomyia* favored locations with  
332 intermediate and high temperatures, confirming the argument that *Stegomyia* mosquitoes have  
333 temperature-based limits to survival (Brady *et al.*, 2013). Overall, anthropogenic and distance factors  
334 appeared to be more important than environmental drivers for the distribution of both *Culex* and  
335 *Stegomyia* mosquitoes. We note that other variables, not considered in this study, may also affect the  
336 distribution of *Culex* and *Stegomyia* mosquitoes and that the genus-level pooling of the species will  
337 mask species distribution determinants, which will be targeted in future work with larger and more  
338 detailed datasets.

339

## 340 5. Conclusions

341 Our results affirm the utility and reliability of the Biomod2 package in R as a valid modeling method  
342 for species distributions, resulting in insights into vector-ecological interactions. Such work will be  
343 further refined with species-level identification and screening for arboviruses. With our study results,  
344 general assumptions can be made about the distribution of mosquitoes belonging to *Culex* and  
345 *Stegomyia* genera in the Taita Hills and the factors that influence their distribution. High population  
346 densities, short distance to roads and gentle slope angles were associated with the occurrence of *Culex*  
347 mosquitoes. For *Stegomyia* mosquitoes, on the other hand, low human population densities, high

348 temperatures and solar radiations as well as poor or moderate vegetation were suitable factors.  
349 Together, these results have implications for applying the approach for identifying risk areas and  
350 optimizing the use of limited resources for mitigation strategies.

351

#### 352 Acknowledgements

353 We are thankful to Essi Korhonen, Joni Uusitalo, Masika Moses and Peter Mwazi for assistance with  
354 the fieldwork, and Miska Luoto, Juha Aalto, Sakari Äärilä and Ninna Malinen for thoughtful  
355 discussion on the topic. We also acknowledge the logistical support provided by the Taita Research  
356 Station of the University of Helsinki. Funding for this work was provided by the Jane and Aatos  
357 Erkkö Foundation, the Ministry for Foreign Affairs of Finland TAITAGIS project  
358 ([http://www.cimo.fi/programmes/hei\\_ici\\_index/programmes/hei\\_ici/projects/taitagis](http://www.cimo.fi/programmes/hei_ici_index/programmes/hei_ici/projects/taitagis)), and the  
359 Academy of Finland. The work was carried out under a research permit for Taita Research Station  
360 from the National Council for Science and Technology, Kenya (permit NCST/RCD/17/012/33).

361

#### 362 References

- 363 Adhikari, H., Heiskanen, J., Siljander, M., Maeda, E., Heikinheimo V., & Pellikka, P.K.E. (2017).  
364 Determinants of Aboveground Biomass across an Afromontane Landscape Mosaic in Kenya.  
365 *Remote Sensing*, 9(8), 827.
- 366 Anderson, K., & Gaston, K. (2013). Lightweight unmanned aerial vehicles will revolutionize spatial  
367 ecology. *Frontiers in Ecology and the Environment*, 11(3), 138-146.
- 368 Bhutta, Z.A., Sommerfeld, J., Lassi, Z.S., Salam, R.A., & Das, J.K. (2014). Global burden,  
369 distribution, and interventions for infectious diseases of poverty. *Infectious Diseases of Poverty*,  
370 3, 21.
- 371 Broberg, A., Salminen, H., Tolvanen, R., & Ylhäisi, J. (2004). Werugha - village in the heart of the  
372 Taita Hills. In Pellikka, P., Ylhäisi, J., & B. Clark (Eds.), *Taita Hills and Kenya, 2004–seminar,*  
373 *reports and journal of a field excursion to Kenya* (pp.108-113). Expedition reports of the  
374 Department of Geography. University of Helsinki.
- 375 Campbell, L. P., Luther, C., Moo-Llanes, D., Ramsey, J. M., Danis-Lozano, R., & Peterson, A. T.  
376 (2015). Climate change influences on global distributions of dengue and chikungunya virus  
377 vectors. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1665).
- 378 Clark, B., & Pellikka, P. (2005, June). The development of a land use change detection methodology  
379 for mapping the Taita Hills, south-east Kenya. CD Rom publication at the meeting of the  
380 Proceedings of the 31st International Symposium of Remote Sensing of the Environment, St.  
381 Petersburg, Russia.

382 Climate-Data.org (2018). Monthly average rainfall in Maktau. <https://en.climate->  
383 [data.org/location/103844/](https://en.climate-data.org/location/103844/) Accessed 30 April 2018.

384 Conley, A., Fuller, D., Haddad, N., Hassan, A., Gad, A., & Beier, J. (2014). Modeling the distribution  
385 of the West Nile and Rift Valley fever vector *Culex pipiens* in arid and semi-arid regions of the  
386 Middle East and North Africa. *Parasites & Vectors*, 7, (289).

387 Crowl, T., Crist, T., Parmenter, R., Belovsky, G., & Lugo, A. (2008). The spread of invasive species  
388 and infectious disease as drivers of ecosystem change. *Frontiers in Ecology and the Environment*,  
389 6, 238–246.

390 Cutler, D.R., Edwards, T.C.Jr., Beard, K.H., Cutler, A., Hess, K.T., Gibson, J., & Lawler, J.J. (2007).  
391 Random forests for classification in ecology. *Ecology*, 88, 2783–2792.

392 Ding, F., Fu, J., Jiang, D., Hao, M., & Lin, G. (2017). Mapping the spatial distribution of *Aedes*  
393 *aegypti* and *Aedes albopictus*. *Acta Tropica*, 178, 155–162.

394 Diniz- Filho, J., Bini, L., & Hawkins, B. (2003). Spatial autocorrelation and red herrings in  
395 geographical ecology. *Global Ecology & Biogeography*, 12, 53–64.

396 Drew, C., Wiersma, Y., & Huettmann, F. (2011). *Predictive species and habitat modeling in*  
397 *landscape ecology*. 319p. The United States.

398 Dukes, J., Pontius, J., Orwig, D., Garnas, J., Rodgers, V., Brazee, N., Cooke, B., Theoharides, K.,  
399 Stange, E., Harrington, R., Ehrenfeld, J., Gurevitch, J., Lerdaun, M., Stinson, K., Wick, R. & M.  
400 Ayres (2009). Responses of insect pests, pathogens, and invasive plant species to climate change  
401 in the forests of northeastern North America: What can we predict? *Canadian Journal of Forest*  
402 *Research*, 39, 231–248.

403 Environment Systems Research Institute (ESRI). (1991). ARC/INFO User's guide. Cell-based  
404 modelling with GRID. Analysis, display and management (Redlands, CA: ESRI). European  
405 Spatial Agency (ESA). (2015). Sentinel-2 user handbook. *ESA Standard Document*, 64.

406 Erdogan, E.H., Pellikka, P., & Clark, B. (2011). Modelling the impact of land-cover change on  
407 potential soil loss in the Taita Hills, Kenya, between 1987 and 2003 using remote-sensing and  
408 geospatial data. *International Journal of Remote Sensing*, 32(21), 5919–5945.

409 Fatima, S., Atif, S., Rasheed, S., Zaidi, F., & Hussain, E. (2016). Species distribution modeling of  
410 *Stegomyia aegypti* in two dengue-endemic regions of Pakistan. *Tropical Medicine and*  
411 *International Health*, 21, 427–436.

412 Fielding, A., & Bell, J. (1997). A review of methods for the assessment of prediction errors in  
413 conservation presence/ absence models. *Environmental Conservation*, 24, 38–49.



- 414 Fornace, K., Drakeley, C., William, T., Espino, F. & Cox, J. (2014). Mapping infectious disease  
415 landscapes: unmanned aerial vehicles and epidemiology. *Trends in Parasitology*, 30(11), 514-  
416 519.
- 417 Franklin, J., & Miller, J. (2010). Statistical methods – modern regression. In: Franklin, J.(ed.)  
418 *Mapping species distribution: spatial inference and prediction*. Cambridge: Cambridge  
419 University Press, (pp. 340).
- 420 Guisan, A., & Zimmermann, N. (2000). Predictive habitat distribution models in ecology.  
421 *Ecological Modeling*, 135(2–3), 147–186.
- 422 Guisan, A., Tingley, R., Baumgartner, J., Naujokaitis-Lewis, I., Sutcliffe, P., Tulloch, A., Regan, T.,  
423 Brotons, L., McDonald-Madden, E., Mantyka-Pringle, C., Martin, T., Rhodes, J., Maggini, R.,  
424 Setterfield, S., Elith, J., Schwartz, M., Wintle, B., Broennimann, O., Austin, M., Ferrier, S.,  
425 Kearney, M., Possingham, H., & Buckley, Y. (2013). Predicting species distributions for  
426 conservation decisions. *Ecology Letters*, 16(12), 1424–1435.
- 427 Hastie, T., Tibshirani, R., & Buja, A. (1994). Flexible discriminant analysis by optimal scoring.  
428 *Journal of the American Statistical Association*, 89(428), 1255–1270.
- 429 Huang, Y-M. (1986). *Aedes (Stegomyia) bromeliae* the yellow fever virus vector in East Africa.  
430 *Journal of Medical Entomology*, 23(2), 196–200.
- 431 Huang, Y-M. (1990). The subgenus *Stegomyia* of *Aedes* in the Afrotropical Region I. The Africanus  
432 group of species (Diptera: Culicidae). *Contributions of the American Entomological Institute*,  
433 26(1), 1–90.
- 434 Huang, Y-M. (2001). A pictorial key for the identification of the subfamilies of Culicidae, genera of  
435 Culicinae, and subgenera of *Aedes* of the Afrotropical Region (Diptera: Culicidae). *Proceedings*  
436 *of the Entomological Society of Washington*, 103, 1–53.
- 437 Huang, Y-M. (2004). The subgenus *Stegomyia* of *Aedes* in the Afrotropical Region with keys to the  
438 species (Diptera: Culicidae). *Zootaxa*, 700(1), 1–120.
- 439 Hutchinson, M.F. (1991). The application of thin plate splines to continent-wide data assimilation. In  
440 Data Assimilation Systems. J.D. Jasper (Ed.). Melbourne: Bureau of Meteorology. *BMRC*  
441 *Research Report*, 27, 104–113.
- 442 Kaplan, I., Dobert, M.K., Marvin, B.J., McLaughlin, J.L., & Whitaker, D.P. (1976). *Area Handbook*  
443 *for Kenya*. Foreign Area Studies (FAS). The American University, Washington D.C., U.S.A, (pp.  
444 472).
- 445 Lembrechts, J., Nijs, I., & Lenoir, J. (2018). Incorporating microclimate into species distribution  
446 models. *Ecography*.

- 447 Lounibos, L.P. (2002). Invasions by insect vectors of human disease. *Annual Review of Entomology*,  
448 47, 233–266.
- 449 McGarigal, K., Cushman, S., Neel, M., & Ene, E. (2002). FRAGSTATS: Spatial pattern analysis  
450 program for categorical maps. Computer software program produced at the University of  
451 Massachusetts, Amherst. <https://www.umass.edu/landeco/research/fragstats/fragstats.html/>  
452 Accessed 7 November 2018.
- 453 Moran, P (1950). Notes on continuous stochastic phenomena. *Biometrika*, 37, 17–23.
- 454 Mosquito taxonomic inventory (MTI). (2017). *Anopheles; Culex; Stegomyia*. [http://mosquito-](http://mosquito-taxonomic-inventory.info/simpletaxonomy/term/8577/)  
455 [taxonomic-inventory.info/simpletaxonomy/term/8577/](http://mosquito-taxonomic-inventory.info/simpletaxonomy/term/8577/) Accessed 5 April 2017.
- 456 Mughini-Gras, L., Mulatti, P., Severini, F., Boccolini, D., Romi, R., Gioia, B., Khoury, C., Bianchi,  
457 R., Montarsi, F., Patregnani, T., Bonfanti, L., Rezza, G., Capelli, G., & Busani, L. (2014).  
458 Ecological niche modelling of potential West Nile virus vector mosquito species and their  
459 geographical association with equine epizootics in Italy. *EcoHealth*, 11, 120.
- 460 Neteler, M., Roiz, D., Rocchini, D., Castellani, C., & Rizzoli, A. (2011). Terra and Aqua satellites  
461 track tiger mosquito invasion: modeling the potential distribution of *Aedes albopictus* in north-  
462 eastern Italy. *International Journal of Health Geographics*, 10 (49).
- 463 Paupy, C., Delatte, H., Bagny, L., Corbel, V. & Fontenille, D. (2009). *Aedes albopictus*, an arbovirus  
464 vector: from the darkness to the light. *Microbes and Infection*, 11, 1177–1185.
- 465 Pellikka, P. K. E., Heikinheimo, V., Hietanen, J., Schäfer, E., Siljander, M. & Heiskanen, J. (2018).  
466 Impact of land cover change on aboveground carbon stocks in Afromontane landscape in Kenya.  
467 *Applied Geography*, 94, 178–189.
- 468 Powell, J. R., & Tabachnick, W. J. (2013). History of domestication and spread of *Aedes aegypti* - A  
469 Review. *Memórias Do Instituto Oswaldo Cruz*, 108(1), 11–17.
- 470 Pradervand, J., Dubuis, A., Pellissier, L., Guisan, A., & Randin, C. (2014). Very high resolution  
471 environmental predictors in species distribution models: Moving beyond topography? *Progress*  
472 *in Physical Geography*, 38, 79-96.
- 473 R Core Team. (2013). R: A language and environment for statistical computing. *R Foundation for*  
474 *Statistical Computing*. Vienna, Austria. URL <http://www.R-project.org/>.
- 475 Reinert, J. F., Harbach, R. E., & Kitching, I. J. (2009). Phylogeny and classification of tribe Aedini  
476 (Diptera: Culicidae). *Zoological Journal of the Linnean Society*, 157, 700–794.
- 477 Reiter, M., & LaPointe, D. (2007). Landscape factors influencing the spatial distribution and  
478 abundance of mosquito vector *Culex quinquefasciatus* (Diptera: Culicidae) in a mixed  
479 residential-agricultural community in Hawai ‘i. *Journal of Medical Entomology*, 44(5), 861–868.

480 Richman, R., Diallo, D., Diallo, M., Sall, A., Faye, O., Diagne, C., Dia, I., Weaver, S., Hanley, K., &  
481 Buenemann, M. (2018). Ecological niche modeling of *Aedes* mosquito vectors of chikungunya  
482 virus in southeastern Senegal. *Parasites & vectors*, *11*(1), 255.

483 Roiz, D., Neteler, M., Castellani, C., Arnoldi, D., & Rizzoli, A. (2011). Climatic factors driving  
484 invasion of the tiger mosquito (*Aedes albopictus*) into new areas of Trentino, Northern Italy.  
485 *PLoS ONE*, *6*(4), e1480.

486 Sallam, M., Xue, R., Pereira, R., & Koehler, P. (2016). Ecological niche modeling of mosquito  
487 vectors of West Nile virus in St. John's County, Florida, USA. *Parasites & Vectors*, *9*, 371.

488 Sallam, M.F., Michaels, S.R., Riegel, C., Pereira, R.M., Zipperer, W., Lockaby, B.G., & Koehler,  
489 P.G. Spatio-Temporal Distribution of Vector-Host Contact (VHC) Ratios and Ecological Niche  
490 Modeling of the West Nile Virus Mosquito Vector, *Culex quinquefasciatus*, in the City of New  
491 Orleans, LA, USA. *International Journal of Environmental Research and Public Health*, *14*(8),  
492 892.

493 Service, M. (1991). *Mosquito ecology: Field sampling methods*. (2<sup>nd</sup> edition). 988 p. Elsevier.

494 Siljander, M., Clark, B., & Pellikka, P. (2011). A predictive modelling technique for human  
495 population distribution and abundance estimation using remote sensing and geospatial data in a  
496 rural mountainous area in Kenya. *International Journal of Remote Sensing*, *32*(21), 5997–6023.

497 Thuiller, W., Georges, D., Engler, R., & Breiner, F. (2016). Package “biomod2”. Ensemble platform  
498 for species distribution modeling. Version 3.3-7. [https://cran.r-](https://cran.r-project.org/web/packages/biomod2/biomod2.pdf/)  
499 [project.org/web/packages/biomod2/ biomod2.pdf/](https://cran.r-project.org/web/packages/biomod2/biomod2.pdf/) Accessed 5 November 2016.

500 Thuiller, W., Lafourcade, B., Engler, R., & Araújo, M. (2009). BIOMOD - A platform for ensemble  
501 forecasting of species distributions. *Ecography*, *32*, 369–373.

502 Tucker, C. (1979). Red and photographic infrared linear combinations for monitoring vegetation.  
503 *Remote Sensing of Environment*, *8*, 127–150. University of Helsinki (2003-2009). The Taita  
504 Hills. <http://www.helsinki.fi/science/taita/taitahills.html/> Accessed 22 April 2016.

505 Virtanen, E. (2015). *Fine-resolution climate grids for species studies in data-poor regions*. Master's  
506 Thesis. University of Helsinki.

507 The World Health Organization (WHO). (2017). Vector-borne diseases. October 2017. Fact sheets.  
508 <http://www.who.int/news-room/fact-sheets/detail/vector-borne-diseases/> Accessed 6 June 2018.

509 Yee, T., & Mitchell, N. (1991). Generalized additive models in plant ecology. *Journal of Vegetation*  
510 *Science*, *2*, 587–602.

511  
512  
513

514 Appendix 1. Multicollinearity of the variables included in the models. All statistical significance  
 515 values (p-values), are marked as follows: statistically highly significant \*\*\* = <0.001, statistically  
 516 significant \*\* = <0.01, statistically significant \* = <0.05 and ns = not statistically significant.  
 517

<i>Culex</i>					
Variable	Slope	NDVI	Population density	Distance to roads	Elevation
Slope		0.115ns	-0.294**	0.172ns	0.279**
NDVI	0.115ns		-0.226*	0.391***	0.557***
Population density	-0.294**	-0.226*		-0.26**	-0.313***
Distance to roads	0.172ns	0.391***	-0.26**		0.143ns
Elevation	0.279**	0.557***	-0.313***	0.143ns	
<i>Stegomyia</i>					
Variable	Mean radiation	NDVI	Population density	Distance to roads	Mean temperature
Mean radiation		-0.04ns	0.312***	-0.12ns	0.10ns
NDVI	-0.04ns		-0.226*	0.391***	-0.543***
Population density	0.313***	-0.226*		-0.26**	0.286**
Distance to roads	-0.12ns	0.391***	-0.26**		-0.12ns
Mean temperature	0.10ns	-0.543***	0.286**	-0.12ns	

518  
 519  
 520  
 521  
 522  
 523  
 524  
 525  
 526  
 527  
 528  
 529  
 530  
 531  
 532  
 533  
 534

535 Appendix 2. Variable importance presented in each model for *Culex* and *Stegomyia* estimates.

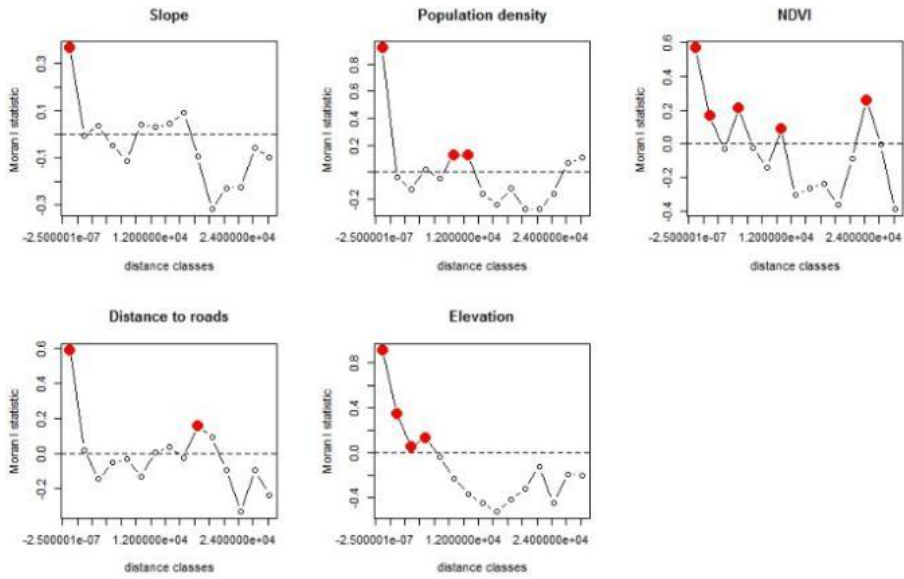
536

<i>Culex</i>								
Variable	GLM	GAM	GBM	CTA	ANN	MARS	RF	Maxent
Slope	0	0.147	0.078	0	0.147	0.385	0.200	0.432
Population density	0.658	0.764	0.678	1	0.773	0.702	0.429	0.513
NDVI	0.261	0.165	0.088	0	0	0	0.114	0.382
Distance to roads	0	0.179	0.039	0	0.092	0	0.064	0.453
Elevation	0	0.028	0.006	0	0.097	0	0.030	0.436
<i>Stegomyia</i>								
Variable	GLM	GAM	GBM	CTA	ANN	MARS	RF	Maxent
Mean radiation	-	0.041	0	0	0.289	-	0.059	0.176
Population density	-	0.575	0.680	0.966	1	-	0.206	0.357
NDVI	-	0.388	0	0	0	-	0.210	0.291
Distance to roads	-	0.337	0	0	0.215	-	0.087	0.028
Mean temperature	-	0	0.358	0	0.081	-	0.089	0.026

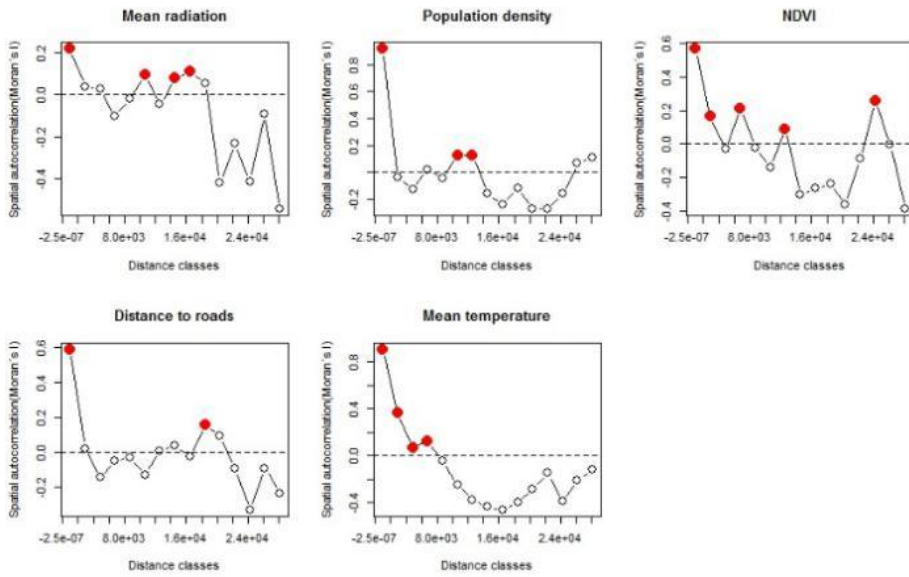
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561

562 Appendix 3. A. SAC of *Culex* predictors. B. SAC of *Stegomyia* predictors. Red circles indicated a  
 563 significant p-value ( $p < 0.05$ ) and indicate distances where the variable was autocorrelated.

564 A.



565 B.



566  
 567  
 568  
 569  
 570  
 571