



University of Chester



**This work has been submitted to ChesterRep – the University of Chester’s
online research repository**

<http://chesterrep.openrepository.com>

Author(s): Nicola A Wardrop ; Matthew Geary ; Patrick E Osborne ; Peter M Atkinson

Title: Interpreting predictive maps of disease, highlighting the pitfalls of species distribution models in epidemiology

Date: 2014

Originally published in: Geospatial Health

Example citation: Wardrop, N. A., Geary, M., Osborne, P. E., & Atkinson, P. M. (2014). Interpreting predictive maps of disease, highlighting the pitfalls of species distribution models in epidemiology. *Geospatial Health*, 9(1), 237-246.

Version of item: Authors' post-print

Available at: <http://hdl.handle.net/10034/338525>

1 **Interpreting predictive maps of disease, highlighting the pitfalls of species distribution models**
2 **in epidemiology**

3
4 Nicola A Wardrop¹, Matthew Geary^{1,2}, Patrick E Osborne³, Peter M Atkinson¹.

5
6 ¹Geography and Environment, Faculty of Social and Human Sciences, University of Southampton,
7 Highfield, Southampton, SO17 1BJ, UK

8 ²Current address: Department of Biological Sciences, University of Chester, Chester, CH1 4BJ, UK

9 ³Centre for Environmental Sciences, Faculty of Engineering and the Environment,
10 University of Southampton, Highfield, Southampton SO17 1BJ, UK
11
12
13

14 **Corresponding author**

15 Dr Nicola Wardrop
16 Geography and Environment
17 Faculty of Social and Human Sciences
18 University of Southampton
19 Highfield, Southampton
20 SO17 1BJ, UK
21

22 Tel.: +44 (0)2380 594 612

23 Email: Nicola.Wardrop@soton.ac.uk.

24
25
26 **Abstract**

27 The application of spatial modelling to epidemiology has increased significantly over the past
28 decade, delivering enhanced understanding of the environmental and climatic factors affecting
29 disease distributions and providing spatially continuous representations of disease risk (predictive
30 maps). These outputs provide significant information for disease control programmes, allowing
31 spatial targeting and tailored interventions. However, several factors (e.g. sampling protocols or
32 temporal disease spread) can influence predictive mapping outputs. This paper proposes a
33 conceptual framework which defines several scenarios and their potential impact on resulting
34 predictive outputs, using simulated data to provide an exemplar. It is vital that researchers recognise
35 these scenarios and their influence on predictive models and their outputs, as a failure to do so may
36 lead to inaccurate interpretation of predictive maps. As long as these considerations are kept in
37 mind, predictive mapping will continue to contribute significantly to epidemiological research and
38 disease control planning.

39
40 **Keywords:** spatial epidemiology; predictive modelling; species distribution modelling.

41 Introduction

42 In recent years, there has been a significant increase in the application of spatial modelling tools to
43 disease studies. This has been driven by the increasing availability of epidemiological,
44 environmental and climatic datasets with spatial (and temporal) dimensions, increased
45 computational capacity, the development of geographical information systems (GIS) and a growing
46 number of spatial analytical tools and platforms capable of handling spatial and space-time datasets.
47 Traditional, non-spatial methods of epidemiological analysis can fail to adequately address major
48 determinants of disease risk. The spatial distributions of many diseases are linked explicitly to
49 environmental conditions (e.g. climatic factors or land cover) and these relationships are most
50 effectively explored, quantified and utilised via spatial visualisation and analysis (Bergquist, 2001).
51 The increasing application of spatial analysis is not unique to epidemiology; there is a close parallel
52 in biodiversity studies, where species distribution modelling (SDM) has proliferated (Elith and
53 Leathwick, 2009). Pathogens can be considered in this context: the tools and theories developed in
54 SDM have useful applications in epidemiological research and *vice versa*.

55
56 The cartographic representation of epidemiological data has many benefits over presentation using
57 tables or plots; images are attention-grabbing, of more interest and allow immediate visual
58 interpretation of spatial patterns (Koch, 2005). Detailed information on the spatial distribution of
59 diseases also provides significant benefits for disease control programmes, particularly for spatially
60 heterogeneous disease distributions (Snow et al., 1996; Simarro et al., 2010). However, just as in the
61 mapping of biodiversity, obtaining comprehensive spatial coverage of a disease within a region of
62 interest is not always possible using disease surveillance data (particularly in developing countries
63 here the infrastructure is often poor). Additionally, the large-scale surveys required to provide
64 complete information are commonly impractical due to financial constraints, logistical issues,
65 security needs and time limitations (Snow et al. 1996; Brooker et al. 2000). These limitations may
66 be overcome, at least in part, using predictive modelling, as described below.

67
68 Statistical methods can be used to fit regression models of the relationship between disease and
69 environment; thus, quantifying the effects of covariates (i.e. variables representing environmental,
70 climatic or landscape factors) on epidemiological measures of disease such as occurrence
71 (presence/absence), prevalence or incidence rates. Models based on covariates, which are measured
72 at the same locations for which epidemiological information is available, but where precise
73 geographical coordinates are absent, and their spatial relationships to one another are not accounted
74 for, focus on environmental space (Elith and Leathwick, 2009). Where covariate information is
75 available covering the full area of interest (e.g. as a raster), these models can be interpolated or
76 extrapolated (prediction within or beyond the range of the training data, respectively) over
77 continuous space; hence, predicting disease at locations for which observed data are not available
78 (Elith and Leathwick, 2009). Prediction with respect to new sites is based on the disease's location
79 in environmental space. These types of models provide information regarding factors driving the
80 observed spatial distribution of disease. The resulting output is a predictive map, also known as a
81 "risk map" (Brooker, 2007), which are is widely used (without incorporating the geographical
82 coordinates) in biodiversity studies (Austin, 2002; Elith and Leathwick, 2009). It can be argued that
83 such models are capable of producing predictive (risk) maps because the main processes
84 determining occurrences are aspatial: it is assumed that species do not respond to location *per se*.

85
86 One potential problem with the approach discussed above is the inability to account for spatial
87 autocorrelation in the residuals (where values close together in space are more similar than values
88 further apart, which occurs commonly when studying the distributions of infectious diseases). This
89 can (i) violate the underlying assumptions of the statistical methods used; and (ii) result in
90 inaccurate models, biased regression parameters, underestimated standard errors, falsely narrow
91 confidence intervals and an overestimation of the significance of covariates, ultimately leading to
92 misinterpretation of the relationships between observations and covariates (Legendre, 1993;

93 [Thomson et al., 1999](#)). In practice, the effect of spatial autocorrelation on prediction accuracy varies
94 among modelling techniques and represents one source of uncertainty in SDM ([Marmion et al.,](#)
95 [2009](#)). However, extension of traditional modelling methods allows the explicit inclusion of spatial
96 information in the modelling process, e.g., the inclusion of both environmental and geographic
97 space in the model. Such extension deals appropriately with the potential problem above. One
98 potential solution involves inclusion of geostatistical spatial prediction of the residuals in a mixed
99 regression model ([Diggle and Ribeiro Jr, 2007](#)). Geostatistical methods incorporate information on
100 the precise location of each observation in relation to other observations to represent spatial
101 autocorrelation, giving increased accuracy of estimates of covariate effects, measures of uncertainty
102 and predictive outputs ([Diggle et al., 2002](#)).

103
104 Predictive mapping of disease (or species distributions more generally) can help overcome the
105 problems associated with sparse datasets. Data from a sample of locations (surveys or surveillance)
106 can be used to fit a model, and subsequent interpolation or extrapolation can provide a spatially
107 continuous prediction of disease ([Brooker, 2007](#)), alleviating the need for comprehensive and large-
108 scale surveys. These outputs can allow the consideration of spatial heterogeneity in disease
109 distributions during planning, implementation and monitoring of interventions, including targeting
110 interventions to areas with the greatest predicted risk of disease ([Clements et al., 2006](#)),
111 identification of areas with a low risk of disease (which can be considered low priority for
112 intervention) ([Clements et al., 2010](#)) and recognition of areas in which intervention may be
113 detrimental ([Diggle et al., 2007](#)). The consideration of uncertainty in outputs allows the delineation
114 of areas from which additional information is required; thus, allowing targeted data acquisition
115 ([Clements et al., 2006](#)).

116
117 The integration of predictive maps and population distribution data allows the estimation of
118 populations at risk of disease and disease burden, providing information to support the allocation of
119 resources (e.g., delivery of adequate supplies of drugs) as described by [Gething et al. \(2011\)](#). The
120 types of outputs described above can also provide valuable resources for advocacy purposes, aiding
121 communication to Government bodies, international organisations and the general public.
122 Additional benefits from predictive mapping include enhanced understanding of the ecology of
123 disease transmission, identification of landscape risk factors and the implication of environmental
124 factors in the spread or distribution of disease ([Wardrop et al., 2010](#)), each of which can allow the
125 development of tailored interventions for specific epidemiological settings.

126
127 The underlying theoretical basis for SDM and predictive mapping is ecological niche theory,
128 particularly Hutchison's model ([Austin, 2002](#)). [Hutchison \(1959\)](#) envisaged the niche as a hyper-
129 volume in multi-dimensional space (each axis being an environmental characteristic) that defines
130 the conditions, under which a population can maintain a positive net growth rate ([Pearman et al.,](#)
131 [2008](#)). The fundamental niche (constrained by genetics and physiology) is defined as distinct from
132 the realised niche (with limitations on resource-use caused by competing species): the realised niche
133 usually seen as a subset of the fundamental niche (see [Pulliam, 2000](#) for exceptions). Vector-borne
134 diseases are interesting in this context since modelling may focus on the vector, the host(s) and/or
135 the disease itself. Furthermore, the vector and host(s) are essentially part of the niche of the disease
136 and, indeed, may control its survival to such an extent that they act as the full niche in certain parts
137 of the life cycle.

138
139 When predictive models are extrapolated (and to some extent interpolated) to new locations (and
140 time periods), two ecological assumptions are necessary: (a) the species is in equilibrium with the
141 environment in the area used to train the model; and (b) the niche is conserved across space and
142 time, i.e. the species-environment relationship is spatially homogeneous ([Broennimann and Guisan,](#)
143 [2008](#); [Nogues-Bravo, 2009](#)). Assumption (a) is violated when ranges are expanding ([Elith et al.,](#)
144 [2010](#)) or where parts of a range are unoccupied by the species (e.g., due to chance or human

145 intervention), but may otherwise hold. There is considerable uncertainty over the applicability of
146 assumption (b) and, indeed, whether it is the realised niche, the fundamental niche, or both that
147 might vary between areas (Pearman et al., 2008). Additionally, careful consideration should be
148 given to the observed epidemiological data and covariate data used in the modelling process. To
149 illustrate how these theoretical underpinnings affect disease modelling in conjunction with the
150 limitations imposed by incomplete or unrepresentative sampling, we applied predictive modelling
151 methods to a simulated dataset under four scenarios.

152

153 **Materials and methods**

154

155 **Study area and data**

156 A hypothetical disease was simulated across an area of East Africa (between latitude 27° and 5° and
157 longitude 22° and 42°; Figure 1). This choice was arbitrary and the disease simulated is not meant
158 to represent any particular existing disease. Environmental data for the disease distribution
159 simulations were downloaded from Worldclim as raster layers at the spatial resolution of 10' and
160 cropped to the study area (Hijmans et al., 2005). Data for mean monthly temperature and mean
161 monthly precipitation were converted to annual averages. Altitude and mean temperature of the
162 wettest quarter were also used in the modelling.

163

164 The disease was simulated to occur in areas with a mean annual temperature between 18.0 and
165 22.5°C and mean annual precipitation between 60 and 170 mm but was not constrained by altitude.
166 As a result of these choices, approximately one quarter of the study area was classified as suitable
167 for disease transmission (26.4 %; Figure 2).

168

169 **Disease scenario sampling**

170 The four scenarios described in Table 1 were investigated using the hypothetical disease described
171 above. Sampling for each of the disease scenarios was performed using the randomPoints function
172 from the dismo package (Hijmans et al., 2013). In each scenario, 300 presence or absence locations
173 were extracted from a true suitability raster and used for model fitting (see Figure 3). In scenarios a)
174 and c) (full information and missing covariates scenarios, respectively), these points were
175 distributed completely randomly across the study area. For scenario b) (heterogeneous sampling
176 effort) these locations were biased towards Kenya (200 locations) rather than the remaining study
177 area (100 locations). For scenario d) (disease not in equilibrium) the presence or absence values for
178 the locations were manipulated so that the disease was recorded normally in Kenya (present/absent),
179 while all of the locations in the remaining study area were recorded as absent: this could represent a
180 situation where the disease is not occupying its full niche due to chance or human intervention.

181

182 **Model fitting and testing**

183 Generalised linear models were fitted to the observed (presence/absence) data from each of the four
184 scenarios: environmental data were extracted for the sample data locations, and logistic regression
185 analysis was applied to quantify the relations between disease presence and the covariates. In each
186 scenario, mean annual temperature, mean annual precipitation and altitude were strongly correlated
187 with one another (Pearson's $c > 0.5$). To avoid problems associated with collinearity, only mean
188 annual precipitation and altitude were included in candidate models for scenarios a), b) and d), and
189 only altitude and mean temperature of the wettest quarter for scenario c). To make meaningful
190 comparisons across scenarios we chose to fit the same model (or its equivalent in the missing
191 covariates scenario) in each case. Based on prior knowledge of the disease distribution, we included
192 an interaction term between altitude and mean annual precipitation (or mean temperature of the
193 wettest quarter for scenario c). For each scenario, 100 simulated sets of sample data were used in
194 the epidemiological distribution models.

195

196 Models were tested using the area under the curve (AUC) of the receiver operating characteristic
197 (ROC) curve (Fielding and Bell, 1997), where a threshold probability of occurrence of 0.5 was used
198 to classify predicted disease presence (or suitability). AUC scores range between 0 and 1: those
199 greater than 0.5 are considered to have predictive ability better than random (for predicting
200 presence), while scores above 0.7 indicate a good predictive ability. ROC plots were constructed
201 and AUC values were calculated using the ROCR package (Sing et al., 2005). Along with the AUC
202 score we assessed the predicted binary distribution (predicted presence, based on a threshold
203 probability of 0.5) from each modelling scenario against the true suitability and calculated the
204 proportion of the study area predicted correctly. These testing metrics were calculated for each
205 scenario over 100 simulations to obtain a full picture of the variability in predictions for each
206 scenario.

207
208 All modelling was performed in R (R Development Core Team, 2013). Spatial functions from the
209 'raster' (Hijmans, 2013), 'rgdal' (Bivand et al., 2013), 'sp' (Pebesma and Bivand, 2005; Bivand et al.,
210 2008) and 'maptools' (Bivand and Lewin-Koh, 2013) packages were also used during the model
211 simulations.

212 213 **Results**

214 215 **Spatial predictions**

216 The models differed with respect to the spatial predictions across the study area (Figure 4) that can
217 be interpreted as predicted probability of occurrence, or predicted suitability for disease. The full
218 information model (scenario a)) predicted an area which broadly matched the actual spatial
219 distribution. However, the predicted area of suitability was slightly larger, particularly in the South
220 of the study area. The missing covariates model (scenario c)) also predicted an area of similar
221 pattern to the simulated disease. However, in this case the area of predicted suitability was broader
222 still and included a patch in the South-west of the study area, which was unsuitable for the disease.
223 The heterogeneous sampling effort model (scenario b)) predicted inaccurately overall with areas on
224 the edges of the study area, outside of the range of the simulated disease, predicted to be suitable.
225 The disease 'not in equilibrium' model (scenario d)) predicted almost all of the study area to be
226 unsuitable. Some small pockets were predicted to be suitable in the north of the region. However,
227 the majority of these pockets were outside the distribution of the simulated disease.

228 229 **Model testing**

230 Figure 5 shows the ROC curves from each of the scenarios, and Figure 6 shows the proportion of
231 the study area that was correctly predicted. The full information model (scenario a)) produced the
232 highest median scores for both AUC (0.77) and the proportion of the study area predicted correctly
233 (0.71). These scores suggest the model has good predictive power. The missing covariates model
234 (scenario c)) was closest to the full information model in terms of performance (median AUC =
235 0.71; median proportion of the study area predicted correctly = 0.68). The AUC scores for both the
236 disease 'not in equilibrium' model (scenario d)) and the heterogeneous sampling effort model
237 (scenario b)) suggest that they perform no better than random in terms of prediction. The disease
238 'not in equilibrium' model performed more accurately in terms of correct prediction of the study
239 area (median = 0.62) than AUC score (median = 0.5). The heterogeneous sampling effort performed
240 less accurately than the other scenarios for both metrics (median AUC = 0.45; median proportion of
241 study area predicted correctly = 0.53).

242
243 Overall, the full information scenario (scenario a)) performed the most accurately in terms of both
244 the proportion of the study area predicted correctly and the AUC score, followed by the missing
245 covariates scenario (scenario c)). The least accurate model was the scenario representing
246 heterogeneous sampling effort (scenario b)) which, along with the disease spreading/control
247 programme scenario, failed to predict disease suitability in the majority of the study area.

248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299

Discussion

As discussed above, statistical models are now used widely to map disease, supplementing traditional epidemiological methodologies leading to enhanced characterisation and understanding of disease distributions and epidemiology. The four scenarios presented above highlight the dependence of predictive mapping outputs on (i) sampling and data considerations; and (ii) contextual factors, such as temporal disease spread in the study area. It is vital that researchers recognise these factors and their influence on predictive models and their outputs as inadvertent use of incomplete or biased data and the use of inadequate covariates may lead to inaccurate interpretation of predictive maps. In addition, absence of consideration for the on-going dynamics of disease transmission and spread within the study area can easily result in erroneous guidance.

Scenario a), which represents the ideal situation where the disease is in equilibrium, representative samples are available and appropriate covariates are being used, is the ideal situation for predictive modelling of disease, although it is likely that many practical examples do not fulfil these criteria. Scenario b) (heterogeneous sampling effort) should be avoided where possible. Indeed, disease prediction studies do not normally make use of spatially biased data as described in this scenario. In addition to the spatial coverage of sampling, statistical models should not be used to provide predictions in areas which are materially different from the area for which training data are available, as the modelled relationships may not be the same (Fitzpatrick and Hargrove, 2009). The example disease provided in this paper was a simulated disease; hence, full information was available on the covariates driving its spatial distribution. However, in real applications, the precise factors which drive the observed distribution are not necessarily known in advance; thus, the subset of potential covariates is selected based on (1) biological understanding; and (2) statistical modelling. This subset may not always represent the most appropriate subset for the disease under consideration, so a lack of data often results in important covariates being omitted from the modelling altogether. Thus, scenario c) (missing covariates) can be considered a frequent occurrence in practical applications.

The final scenario (disease not occupying its full niche) is likely to be the most common scenario encountered in spatial epidemiology applications. As an example, Rhodesian sleeping sickness (caused by the parasite *Trypanosoma brucei rhodesiense*) has been spreading in Uganda over the past two to three decades, with the movement of infected livestock implicated in the most recent introductions (Fèvre et al., 2001; Wardrop et al., 2010). This indicates that historically, the recorded spatial distribution of Rhodesian sleeping sickness did not cover all areas environmentally suitable for the disease and any predictive modelling based upon this distribution would not necessarily be providing the output intended. Most SDM are blind to the mechanisms that promote dispersal from affected to unaffected areas (e.g., human movements, contact patterns and trade), or factors that may inhibit spatial spread of a disease (e.g., human intervention), so resulting predictions are at best maps of potential risk. Most ecosystems are dynamic, and the spatial dispersal of a disease over time is not uncommon, enabling the disease to occupy a larger proportion of its potential range (Reisen, 2010). The identification and quantification of factors influencing this expansion would be required to ascertain the future risk of disease within currently unaffected areas. As Soberon (2010) argues, the fundamental ecological factors that determine species distributions are environment, biotic interactions and movements; without all three of these, modelled outputs of predicted occurrence and hence risk are compromised.

The four scenarios developed here should be taken into consideration when designing surveys and collecting data, fitting statistical models and during subsequent interpretation of predictive outputs. The goal of mapping should be clear from the outset (e.g., to map the present distribution or to map suitability) due to the impact of data acquisition choices on the final outputs. The consideration of whether an epidemiological situation may incorporate one (or more) of these scenarios should

300 provide greater awareness of the potential impacts on the modelling process and predictive maps.
301 Model coefficients and estimates of uncertainty can only take us so far; the interpretation of these
302 outputs needs to be undertaken with the four scenarios presented in this framework in mind to
303 ensure accurate comprehension of meaning and consequent sound action in relation to decision-
304 making. The premise of SDM is that predictive outputs will represent environmental suitability.
305 However, where input data is not comprehensive, or where dynamic factors have not been taken
306 into account, the predictive outputs may not represent environmental suitability, but may more
307 accurately be described as representing the current distribution of the disease of interest.
308

309 Using a simulated dataset, this paper provides an overview of predictive mapping of disease and the
310 linkages with ecological SDM, and has introduced some important considerations, which are rarely
311 discussed in the predictive mapping literature. Care must be taken when carrying out predictive
312 mapping when the distribution of the disease of interest is changing, and a full understanding of the
313 disease's ecology alongside historical, recent and current spatial distributions of the disease should
314 be used to inform the process of modelling and interpretation. Every mapping scenario will have
315 different complexities which may influence the interpretation of resulting predictions, but time
316 spent considering what the observed data represent and the implications of the possible scenarios
317 detailed above will provide a starting point for more accurate interpretation of predictive maps. As
318 long as the considerations introduced here are kept in mind, predictive mapping will continue to
319 contribute significantly to epidemiological research and disease control planning.
320

321 **Acknowledgements:**

322 This work was supported by the Medical Research Council (PMA, NAW - projects G0902445 and
323 MR/J012343/1). The funders had no role in the decision to publish or in preparation of the
324 manuscript.
325
326
327

328 **References**

- 329
330 Austin MP, 2002. Spatial prediction of species distribution: an interface between ecological theory
331 and statistical modelling. *Ecol Model* 157, 101-118
- 332 Bergquist NR, 2001. Vector-borne parasitic diseases: new trends in data collection and risk
333 assessment. *Acta Trop* 79, 13-20
- 334 Bivand R, Keitt T, Rowlingson B, 2013. Bindings for the geospatial data abstraction library. R
335 package version 0.8-10.
- 336 Bivand, R. & Lewin-Koh, N. Maptools: tools for reading and handling spatial objects. R package
337 version 0.8-25. 2013.
- 338 Bivand, R., Pebesma, E.J., & Gomez-Rubio, V. 2008. *Applied Spatial Data Analysis with R* New
339 York, Springer.
- 340 Broennimann, O. & Guisan, A. 2008. Predicting current and future biological invasions: both native
341 and invaded ranges matter. *Biol Lett*, 4, 585-589
- 342 Brooker, S. 2007. Spatial epidemiology of human schistosomiasis in Africa: risk models,
343 transmission dynamics and control. *Trans R Soc Trop Med Hyg*, 101, 1-8
- 344 Brooker, S., Rowlands, M., Haller, L., Savioli, L., & Bundy, D.A.P. 2000. Towards an atlas of
345 human helminth infection in sub-Saharan Africa: The use of geographical information systems

- 346 (GIS). *Parasitol Today*, 16, 303-307
- 347 Clements, A.C.A., Kur, L.W., Gatpan, G., Ngondi, J.M., Emerson, P.M., Lado, M., Sabasio, A., &
348 Kolaczinski, J.H. 2010. Targeting Trachoma Control through Risk Mapping: The Example of
349 Southern Sudan. *PLoS Negl Trop Dis*, 4,
- 350 Clements, A.C.A., Lwambo, N.J.S., Blair, L., Nyandindi, U., Kaatano, G., Kinung'hi, S., Webster,
351 J.P., Fenwick, A., & Brooker, S. 2006. Bayesian spatial analysis and disease mapping: tools to
352 enhance planning and implementation of a schistosomiasis control programme in Tanzania. *Tropical*
353 *Medicine & International Health*, 11, 490-503
- 354 Diggle, P., Moyeed, R., Rowlingson, B., & Thomson, M. 2002. Childhood malaria in the Gambia: a
355 case-study in model-based geostatistics. *J Roy Stat Soc C-App*, 51, 493-506
- 356 Diggle, P. J. & Ribeiro Jr, P. J. 2007, "An overview of model-based geostatistics," In *Model-based*
357 *Geostatistics*, New York: Springer, pp. 27-45.
- 358 Diggle, P.J., Thomson, M.C., Christensen, O.F., Rowlingson, B., Obsomer, V., Gardon, J., Wanji, S.,
359 Takougang, I., Enyong, P., Kamgno, J., Remme, J.H., Boussinesq, M., & Molyneux, D.H. 2007.
360 Spatial modelling and the prediction of *Loa loa* risk: Decision making under uncertainty. *Ann Trop*
361 *Med Parasitol*, 101, 499-509
- 362 Elith, J., Kearney, M., & Phillips, S. 2010. The art of modelling range-shifting species. *Methods in*
363 *Ecology and Evolution*, 1, 330-342
- 364 Elith, J. & Leathwick, J.R. 2009. Species Distribution Models: Ecological Explanation and
365 Prediction Across Space and Time. *Annu Rev Ecol Syst*, 40, 677-697
- 366 Fèvre, E.M., Coleman, P.G., Odiit, M., Magona, J.W., Welburn, S.C., & Woolhouse, M.E.J. 2001.
367 The origins of a new *Trypanosoma brucei rhodesiense* sleeping sickness outbreak in eastern
368 Uganda. *Lancet*, 358, 625-628
- 369 Fielding, A.H. & Bell, J.F. 1997. A review of methods for the assessment of prediction errors in
370 conservation presence/absence models. *Environ Conserv*, 24, 38-49
- 371 Fitzpatrick, M.C. & Hargrove, W.W. 2009. The projection of species distribution models and the
372 problem of non-analog climate. *Biodivers Conserv*, 18, 2255-2261
- 373 Gething, P.W., Patil, A.P., Smith, D.L., Guerra, C.A., Elyazar, I.R.F., Johnston, G.L., Tatem, A.J., &
374 Hay, S.I. 2011. A new world malaria map: *Plasmodium falciparum* endemicity in 2010. *Malaria*
375 *Journal*, 10,
- 376 Hijmans, R. J. raster: Geographic data analysis and modeling. R package version 2.1-49. 2013.
- 377 Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., & Jarvis, A. 2005. Very high resolution
378 interpolated climate surfaces for global land areas. *Int J Climatol*, 25, 1965-1978
- 379 Hijmans, R. J., Phillips, S., Leathwick, J. R., & Elith, J. dismo: Species distribution modeling. R
380 package version 0.8-17. 2013.
- 381 Hutchinson, G.E. 1959. Homage to Santa Rosalia, or why are there so many kinds of animals. *Am*
382 *Nat*, 93, 245-249
- 383 Koch, T. 2005. *Cartographies of Disease: Maps, Mapping, and Medicine*, 1 ed. Redlands, CA, ESRI

384 Press.

385 Legendre, P. 1993. Spatial autocorrelation - trouble or new paradigm. *Ecology*, 74, 1659-1673

386 Marmion, M., Luoto, M., Heikkinen, R.K., & Thuiller, W. 2009. The performance of state-of-the-art
387 modelling techniques depends on geographical distribution of species. *Ecol Model*, 220, 3512-3520

388 Nogues-Bravo, D. 2009. Predicting the past distribution of species climatic niches. *Global Ecol*
389 *Biogeogr*, 18, 521-531

390 Pearman, P.B., Guisan, A., Broennimann, O., & Randin, C.F. 2008. Niche dynamics in space and
391 time. *Trends Ecol Evol*, 23, 149-158

392 Pebesma, E. J. & Bivand, R. Classes and methods for spatial data in R. *R News* 5. 2005.

393 Pulliam, H.R. 2000. On the relationship between niche and distribution. *Ecol Lett*, 3, 349-361

394 R Development Core Team. R: A language and environment for statistical computing. 2013.
395 Vienna, Austria, R Foundation for Statistical Computing.

396 Reisen, W.K. 2010. Landscape Epidemiology of Vector-Borne Diseases. *Annu Rev Entomol*, 55,
397 461-483

398 Simarro, P.P., Cecchi, G., Paone, M., Franco, J.R., Diarra, A., Ruiz, J.A., Fevre, E.M., Courtin, F.,
399 Mattioli, R.C., & Jannin, J.G. 2010. The Atlas of human African trypanosomiasis: a contribution to
400 global mapping of neglected tropical diseases. *International Journal of Health Geographics*, 9,

401 Sing, T., Sander, O., Beerenwinkel, N., & Lengauer, T. 2005. ROCR: visualizing classifier
402 performance in R. *Bioinformatics*, 21, 3940-3941

403 Snow, R.W., Marsh, K., & leSueur, D. 1996. The need for maps of transmission intensity to guide
404 malaria control in Africa. *Parasitol Today*, 12, 455-457

405 Soberon, J.M. 2010. Niche and area of distribution modeling: a population ecology perspective.
406 *Ecography*, 33, 159-167

407 Thomson, M.C., Connor, S.J., D'Alessandro, U., Rowlingson, B., Diggle, P., Cresswell, M., &
408 Greenwood, B. 1999. Predicting malaria infection in Gambian children from satellite data and bed
409 net use surveys: The importance of spatial correlation in the interpretation of results. *Am J Trop*
410 *Med Hyg*, 61, 2-8

411 Wardrop, N., Atkinson, P.M., Gething, P.W., Fèvre, E.M., Picozzi, K., Kakembo, A., & Welburn, S.
412 2010. Bayesian Geostatistical Analysis and Prediction of Rhodesian Human African
413 Trypanosomiasis. *PLoS Negl Trop Dis*, 4, e914

414

415

416 **Table titles:**
417 **Table 1.** Four scenarios for disease modelling

418
419 **Figure titles:**

420 **Figure 1.** Map of Africa showing the bounding box of the study area in green.

421 **Figure 2.** Environmental suitability for hypothetical disease: suitable areas are shown in green and
422 unsuitable areas in grey.

423 **Figure 3.** Example sample data used for the disease modelling scenarios showing (a) the full
424 information scenario, (b) the heterogeneous sampling effort scenario, (c) the missing covariates
425 scenario and (d) the disease not in equilibrium scenario: 100 simulated datasets were created for
426 each scenario. Presence records are shown in red and absence records in black. Actual
427 environmental suitability for disease transmission is shown in green.

428 **Figure 4.** Actual suitability for disease occurrence (top) and predicted probability of disease
429 presence across the study area for each of the four scenarios: scenario a) full information (centre
430 left), scenario b) heterogeneous sampling effort (centre right), scenario c) missing covariates
431 (bottom left) and scenario d) disease not in equilibrium scenario (bottom right).

432 **Figure 5.** Mean ROC curves for 100 simulations of the four scenarios with 95% confidence
433 intervals (dotted lines) showing (a) the full information scenario, (b) the heterogeneous sampling
434 effort scenario, (c) the missing covariates scenario and (d) the disease not in equilibrium scenario.

435 **Figure 6.** Results of 100 simulations for the four scenarios showing (a) proportion of the study area
436 for which predictions were correct (based on a cut-off probability of 0.5) and (b) AUC scores.

437
438

Table 1.

Scenario	Situation
<i>a)</i> Full information	The disease is in equilibrium with its environment and data are available for a spatially representative sample of its range.
<i>b)</i> Heterogeneous sampling effort	The disease is in equilibrium with its environment, but there is spatial bias in the detection of the disease (i.e. a heterogeneous sampling effort).
<i>c)</i> Missing covariates	The disease is in equilibrium with its environment and there is a spatially representative sample available, but the covariates used for prediction do not fully reflect the species environmental constraints.
<i>d)</i> Disease not in equilibrium with the environment	The disease is not in equilibrium with its environment due to either successful disease control (disease no longer occupying its full niche) or on-going spatial spread (the disease does not yet occupy its full niche).

Figure 1

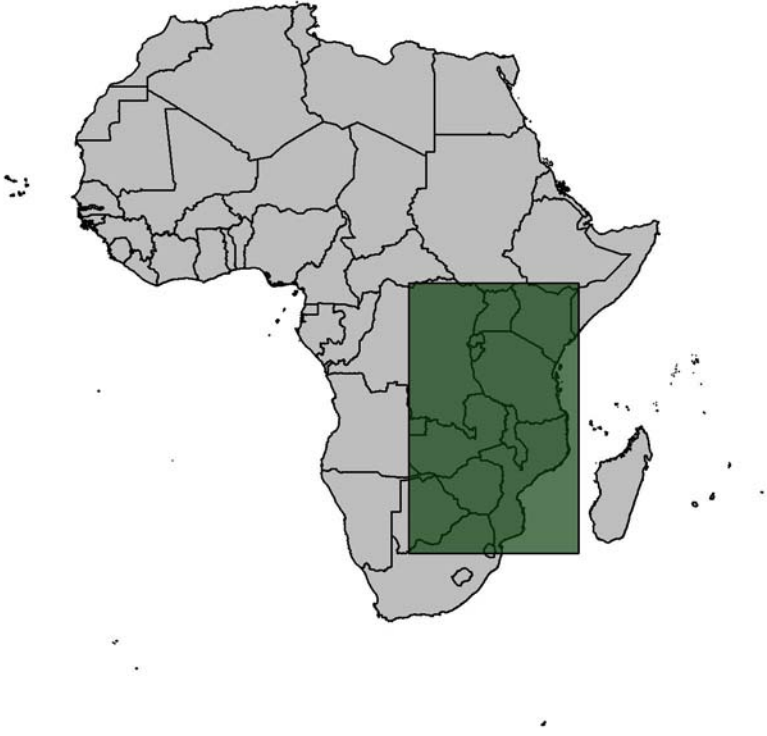


Figure 2

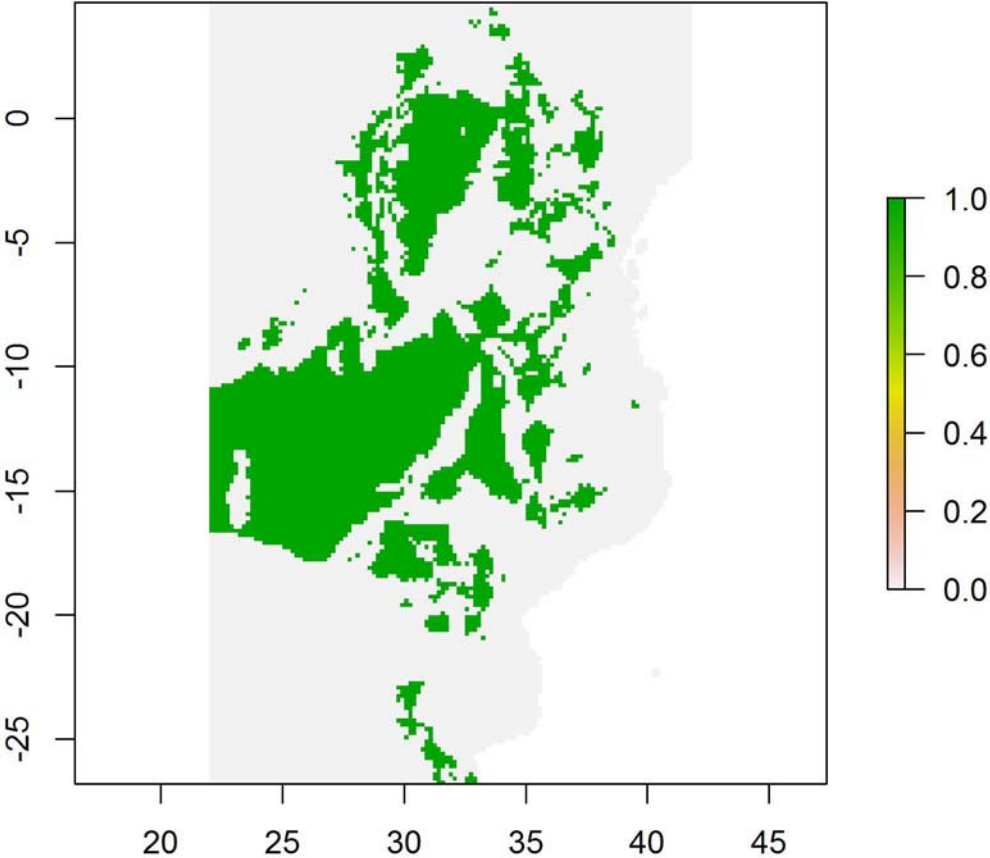


Figure 3

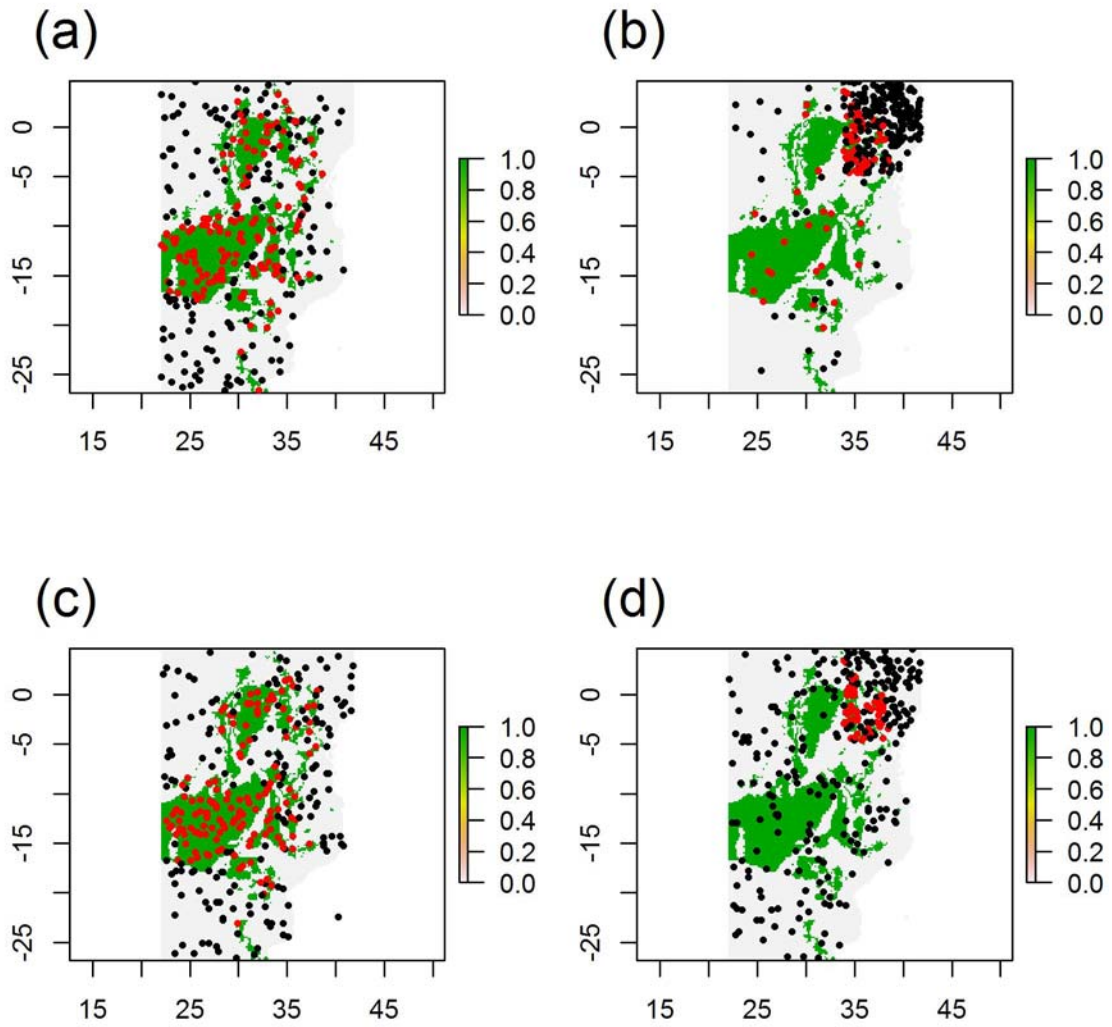


Figure 4

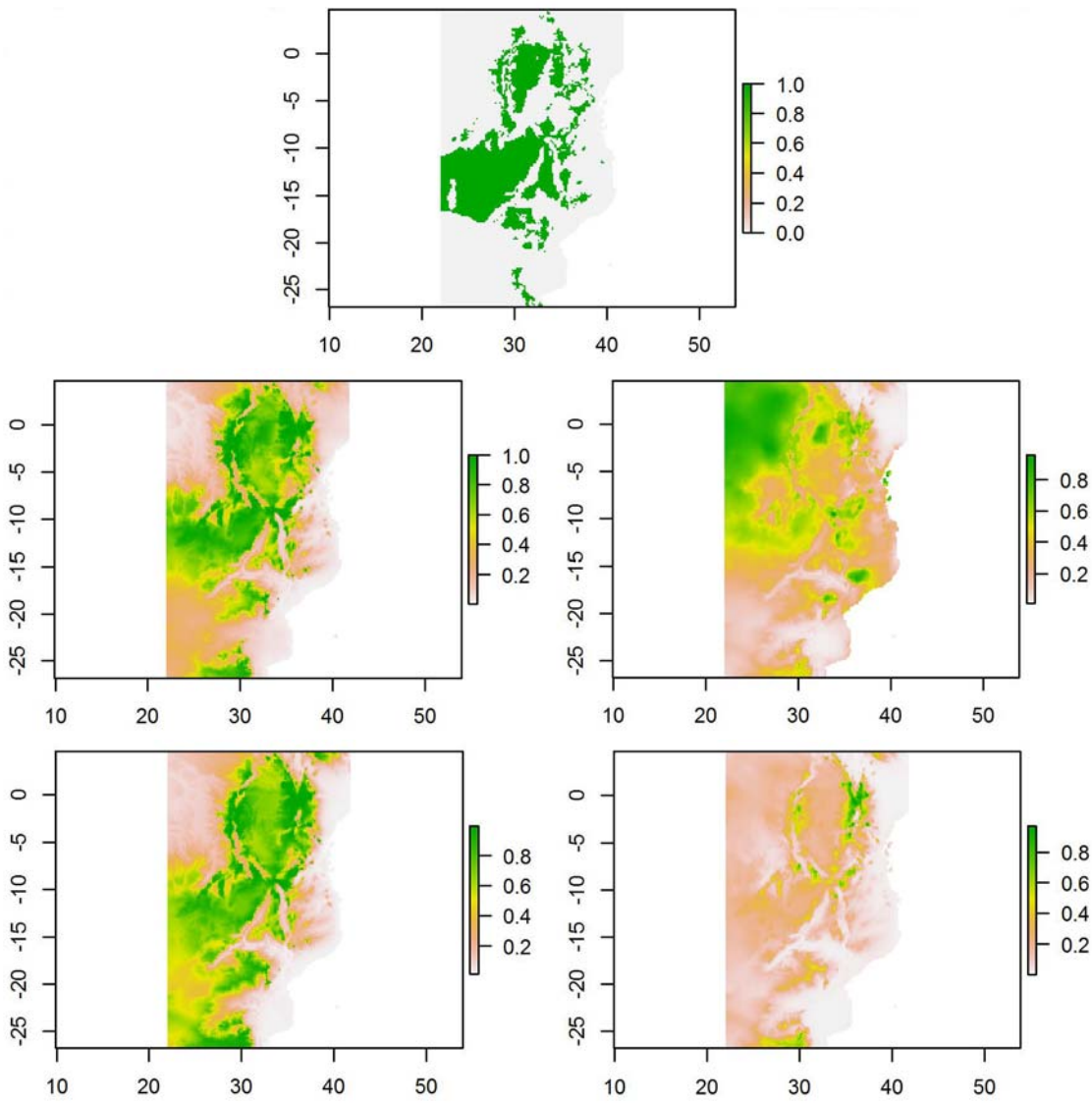


Figure 5

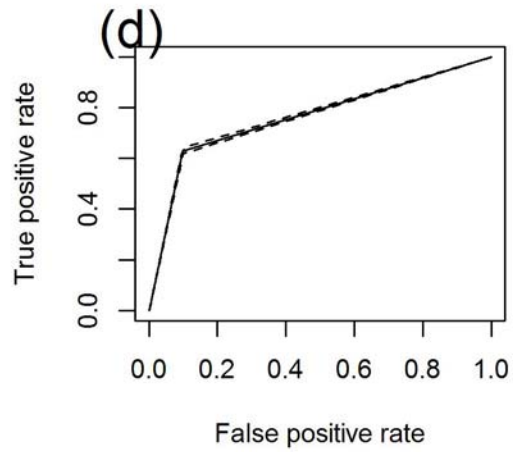
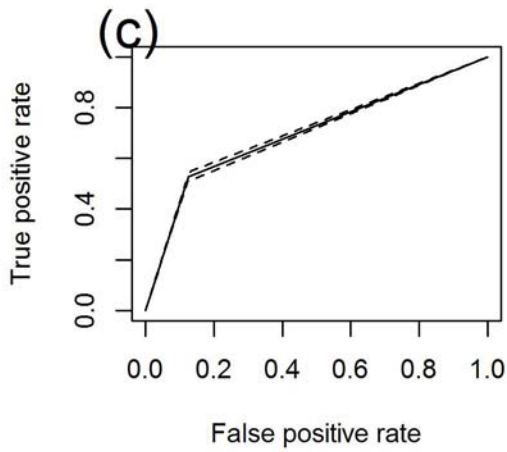
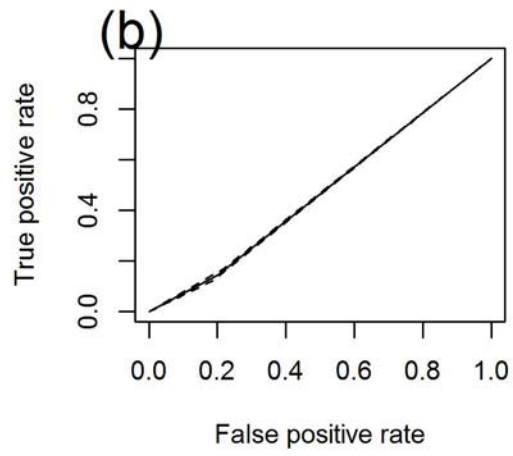
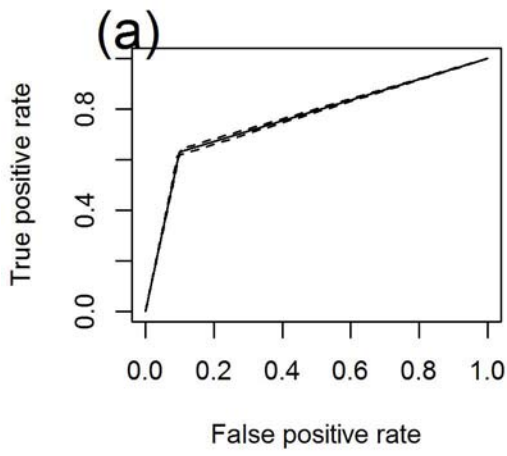


Figure 6

