

1 Tonini et al.: Niche Modeling of Invasive
2 Termites
3
4 **Environmental Entomology: Population**
5 **Ecology**

11 Francesco Tonini
12 University of Florida, Fort Lauderdale
13 Research and Education Center, 3205
14 College Avenue,
15 Davie, Florida, 33314, U.S.A.
16 Phone: +1 954-577-6392
17 Fax: +1 954-424-6851
18 Email: ftonini@ufl.edu

6
7
8
9
10

19
20

21 **Predicting the Geographical Distribution of Two Invasive Termite Species from**
22 **Occurrence Data**

23

24 FRANCESCO TONINI¹, FABIO DIVINO², GIOVANNA JONA LASINIO³, HARTWIG H.
25 HOCHMAIR¹, RUDOLF H. SCHEFFRAHN¹

26

27 ¹University of Florida, Fort Lauderdale Research and Education Center, 3205 College Avenue,
28 Davie, Florida, 33314, U.S.A.

29 ²Division of Physics, Computer Science, and Mathematics, University of Molise, Contrada Fonte
30 Lappone, 86090, Pesche (IS), Italy.

31 ³DSS, University of Rome "La Sapienza", P.le Aldo Moro 5, 00185 Rome, Italy.

32

33 **Abstract**

34 Predicting the potential habitat of species under both current and future climate change scenarios
35 is crucial for monitoring invasive species and understanding a species' response to different
36 environmental conditions. Frequently, the only data available on a species is the location of its
37 occurrence (presence-only data). Using occurrence records only, two models were used to
38 predict the geographical distribution of two destructive invasive termite species, *Coptotermes*
39 *gestroi* (Wasmann) and *Coptotermes formosanus* Shiraki. The first model uses a Bayesian linear
40 logistic regression approach adjusted for presence-only data while the second one is the widely
41 used maximum entropy approach (Maxent). Results show that the predicted distributions of both
42 *C. gestroi* and *C. formosanus* are strongly linked to urban development. The impact of future
43 scenarios such as climate warming and population growth on the biotic distribution of both
44 termite species was also assessed. Future climate warming seems to affect their projected
45 probability of presence to a lesser extent than population growth. The Bayesian logistic approach
46 outperformed Maxent consistently in all models according to evaluation criteria such as model
47 sensitivity and ecological realism. The importance of further studies for an explicit treatment of
48 residual spatial autocorrelation and a more comprehensive comparison between both statistical
49 approaches is suggested.

50

51 **Keywords:** Bayesian logistic modeling, Maxent, presence-only data, subterranean termite,
52 species distribution models

53

54

55

Introduction

56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78

The use of statistical models to predict a species' potential habitat has seen a growing interest during the past two decades given the importance of monitoring endangered or invasive species and understanding a species' response to different environmental conditions (Guisan and Thuiller 2005). Such models are often referred to as habitat models, ecological niche models, or species distribution models (SDMs) (Elith and Leathwick 2007) and have been applied to a variety of fields such as ecology, conservation, and biogeography. SDMs attempt to model the species-environment relationships by using sites of known occurrence (presence data) and, sometimes, non-occurrence (absence data) together with environmental variables recorded over the whole study area. In most cases, records from atlases, herbaria, or museum databases only contain information on a species' incidental observations (Franklin 2009). A fundamental limitation of presence-only datasets is that the prevalence of a species, i.e. the proportion of occupied sites across the study area, is unknown. In recent years, several statistical methods have been proposed for modeling these types of datasets, such as inhomogeneous Poisson process (IPP), (Warton and Sheperd 2010, Chakraborty et al. 2011), and maximum entropy (Maxent) (Phillips et al. 2004, Phillips et al. 2006). Other approaches use presence-absence models by assuming random samples chosen from the region of interest (background samples) as absences (also called "pseudo-absences") (Elith et al. 2006). However, this assumption has been shown to have substantial problems of model specification, interpretation, and implementation (Warton and Sheperd 2010).

In this work, a recently developed Bayesian logistic regression model adjusted for presence-only data (Divino et al. 2011, Divino et al. 2013) and the widely used maximum entropy

79 approach were used to predict the current and future biotic distributions of two major invasive
80 termite pests within the state of Florida: the Asian subterranean termite (AST), *Coptotermes*
81 *gestroi* (Wasmann), and the Formosan subterranean termite (FST), *Coptotermes formosanus*
82 Shiraki. The Bayesian approach used herein has only been tested on artificial data prior to this
83 study (Divino et al. 2013).

84 The highly invasive AST and FST are, or will become, the most destructive subterranean
85 termites in areas of suitable climate, causing severe damage to wood in service (Evans et al.
86 2013). AST is endemic to southeast Asia and it is currently found mostly in tropical areas (Li et
87 al. 2009). FST is probably endemic to southern China and is found primarily in subtropical and
88 temperate regions (Li et al. 2009). AST and FST are only known to occur sympatrically in
89 Taiwan, Florida, and Hawaii (Li et al. 2010). AST was first found in Florida in 1996 (Dade
90 County) and is a more recent invasive species compared to FST, discovered in Florida in 1980
91 (Broward County) (Scheffrahn 2013). Both species are now well established pests in Florida.

92 Regional predictions of the potential habitat of the two termite species under both current and
93 future climate scenarios are currently lacking in the available literature. A single recent study
94 attempted to predict the ecological niche of AST on a global scale using mostly coarse-precision
95 occurrence data derived from the literature (Li et al. 2013). However, the reliability of such
96 predictions could be affected by the excessive extent of the study area used for both model
97 calibration and estimation, given the small amount of available occurrence data.

98 The format of this paper is as follows. The study area, data, variables, and modeling
99 approaches used are described in the Materials and Methods section. Results and their
100 interpretation are then presented, followed by a final discussion on the advantages and
101 limitations of the models tested herein.

102

103

Materials and Methods

104

Study Area and Species Data

106

107 Florida was selected as a common study area for both AST and FST in order to compare the
108 performance of two different statistical approaches under the same environmental conditions.
109 Termite collection localities, including winged reproductives and/or nonvolent foragers, were
110 taken from the University of Florida Termite Collection at the Fort Lauderdale Research and
111 Education Center. Winged reproductives were taken from within infested structures, and
112 therefore in close proximity to their foraging nest mates and stationary nests. Geographical
113 coordinates of 280 and 411 separate land-based colonies of AST (1996-2012) and FST (1985-
114 2012), respectively, were used in this study (Fig. 1). A few records representing boat infestations
115 (Scheffrahn and Crowe 2011) were excluded. All database samples were collected less than 40m
116 from buildings by R.H.S., pest control professionals, property owners, entomologists, and others
117 interested in species-level identification. About 95% of foraging caste samples were collected
118 within 5-10 m or inside the structures themselves.

119

120

Figure 1–caption at the end of file

121

122 The study area was divided into roughly 38,000 2-km grid cells and all termite observations
123 falling within a given cell were aggregated to a single point. After aggregation, a total of 65 and
124 160 occurrences were considered for AST and FST, respectively. In this work, grid cells were

125 considered as independent given the explanatory variables and the probability of presence was
126 modeled for each one of them. The chosen spatial resolution attenuates some of the bias caused
127 by spatial dependence between nearby occurrences because termite reproductives from a mature
128 colony fly only a few hundred meters during their annual dispersal flights (Nutting 1969).
129 Moreover, the available environmental data used in this study were obtained at a 2.5-arcmin (~4
130 km) resolution and it is appropriate to consider a sampling unit whose size is equal (or close) to it
131 (Elith and Leathwick 2009). Finally, the number of occurrences available after the
132 aforementioned spatial aggregation ensures robustness of the estimates from the statistical
133 models used herein.

134

135 **Predictor Variables**

136

137 A set of gridded climatic variables was selected (Table 1) based on both its ability to directly
138 influence the ecophysiology of both AST and FST (Gautam and Henderson 2011), and on
139 suggestions taken in consultation with termite experts. Data for historical climatic conditions
140 were extracted from two sources: (i) the PRISM Climate Group database (Daly et al. 2002) and
141 (ii) the WorldClim (1950-2000) database (Hijmans et al. 2005). General annual trends such as
142 annual total precipitation (prec), average daily mean dew point temperature (dew), and average
143 daily maximum (tmax) and minimum (tmin) temperatures were obtained from the PRISM
144 database, representative of average historical conditions for the years of available occurrence
145 records of both AST and FST. Two bioclimatic variables representing extreme or limiting factors
146 such as maximum temperature of the warmest month (bio5) and minimum temperature of the
147 coldest month (bio6) were chosen from the WorldClim database, representative of 1950-2000

148 average historical conditions. Both WorldClim and PRISM data were obtained at a 2.5-arcmin
149 (~4 km) resolution and further resampled down using bilinear interpolation to maintain the
150 higher data resolution of the reference spatial grid over the study area. The available time series
151 of historical climate data from PRISM (1895-Present) allowed us to extract those years that
152 matched historical occurrence records for both AST and FST exactly.

153 In addition to climate variables, the U.S. Geological Survey National Land Cover Database
154 (NLCD) 2006 (Multi-Resolution Land Characteristics Consortium 2012) was also used (see
155 Table 1), which has a native resolution of 30 m. The database comes with 20 land cover classes,
156 which were modified according to the following steps: (1) reduction from 20 to 8 main land
157 cover classes according to the NLCD 2006 product legend; (2) creation of single layers for each
158 land cover class from the previous step; and (3) aggregation of each land cover layer from 30 m
159 to our 2-km reference grid by expressing each cell value as the percentage of land cover
160 contained within.

161 Finally, centroids of grid cells occupied by termite locations were also used in some of the
162 statistical models (see Tables 2 and 3) in order to account for the geographic proximity between
163 collection sites across the geographic space. Locations were expressed by their projected easting
164 and northing values. All layers, including the 2-km reference grid, were mapped using the
165 NAD83 / Florida GDL Albers projection to minimize distance distortions throughout the study
166 area.

167

168 **Table 1–caption at the end of file**

169

170 Most predictor variables in our dataset were highly correlated and their simultaneous
171 presence in statistical models has been proven to cause several problems (e.g. biased parameter
172 estimates or lower efficiency in the estimates) (Farrar and Glauber 1967). Therefore, an *a priori*
173 choice of variables was carried out in order to exclude pairs of highly correlated variables ($r \geq$
174 0.8). A different set of models was also estimated using principal components obtained from the
175 full set of predictor variables considered herein. Principal component loadings are shown in
176 Supp. Table S1 (available in the online version).

177

178 **Future Scenarios**

179

180 The Intergovernmental Panel on Climate Change (IPCC) Fourth Assessment Report (AR4)
181 describes a set of alternative CO₂ emissions scenarios grouped under four main narrative
182 storylines (Intergovernmental Panel on Climate Change 2007). In this study, the A2 emission
183 scenario was used, which forecasts an average increase in global surface temperature of about
184 3.4°C by 2100. This scenario was preferred over others in order to assess the impact of a larger
185 climate change on both termite species' potential distributions and consider it as a benchmark
186 "worst-case" scenario.

187 Given the uncertainty associated with the path of future climate change, average projections
188 of annual precipitation and minimum/maximum temperatures for the years 2040-2069 (referred
189 to as 2050s hereafter) were extracted from the Climate Change, Agriculture and Food Security
190 (CCAFS) web portal (Climate Change Agriculture and Food Security 2013). Projections of
191 annual mean dew point temperatures were not available from any data provider, hence this
192 predictor variable could not be considered for future scenarios.

193 The three following Atmospheric-Oceanic Global Circulation Models (GCMs hereafter)
194 (Diniz-Filho et al. 2009), statistically downscaled using the so-called delta method (Ramirez-
195 Villegas and Jarvis 2010), were selected for the A2 emission scenario and the 2050s time frame:
196 GFDL-CM 2.1, NCAR-CCSM 3.0, UKMO-HadCM3.

197 A projected population growth scenario in 2060 was obtained from the University of Florida
198 GeoPlan Center (University of Florida Geoplan Center 2013). The dataset assumes no further
199 population growth in areas currently urbanized.

200 First, we assumed a change in the climatic variables under the A2 emission scenario given by
201 the three selected GCMs, assuming no change in population. Then, we added a population
202 growth scenario together with climate change, resulting in a total of six future scenarios for each
203 species. To create "consensus" maps of projected probabilities in the 2050s, predictions were
204 averaged over the three GCMs. This method has been shown to significantly increase the
205 accuracy of species distribution forecasts (Marmion et al. 2009).

206

207 **Modeling Approaches**

208

209 In this study, several models were considered using two statistical approaches for presence-
210 only data (see Tables 2 and 3): (i) maximum entropy (Phillips et al. 2006); and (ii) a Bayesian
211 linear logistic regression adjusted for presence-only data, named Bayesian for presence-only data
212 (BPOD) hereafter (Divino et al. 2011, Divino et al. 2013). The former, was presented in Phillips
213 et al. (2004) and it is widely used for modeling distributions of species (Elith et al. 2011). The
214 BPOD, builds upon the work presented in Ward et al. (2009) while using a Bayesian framework.

215 Although different in their theoretical backgrounds, both methods use the Bayes' rule as an
216 important point to calculate the probability of presence of the species conditioned on the
217 environment. An outline of the theory, main assumptions, and modeling settings used in both
218 approaches follows.

219

220 **Maximum Entropy Approach.**

221

222 Maximum entropy (Maxent hereafter) is a machine-learning method that uses species
223 occurrences and a random sample of background environmental data over a region of interest to
224 predict species distributions. Let us define $Pr(X = x | Y = 1)$ to be the probability distribution
225 of covariates, i.e. environmental variables, across locations where the species is observed ($Y = 1$),
226 and $Pr(X = x | Y = 0)$ to be the probability distribution of covariates where the species is
227 absent ($Y = 0$). The quantity of interest is the probability of presence of a species,
228 $Pr(Y = 1 | X = x)$, conditioned on a set of environmental covariates X . Maxent considers the
229 modeling of $Pr(X = x | Y = 1)$ and uses the Bayes' rule to estimate the sought conditional
230 probability distribution:

$$Pr(Y = 1 | X = x) = \frac{Pr(x | Y = 1)Pr(Y = 1)}{Pr(x)}$$

231 The core of the Maxent "raw" model output is the estimate of the ratio $Pr(x | Y = 1)/P(x)$.

232 This is accomplished by seeking an estimate of $Pr(x | Y = 1)$ that is consistent with available

233 occurrence data. Among several possible distributions, one that maximizes the entropy of

234 $Pr(x | Y = 1)$ or, in other words, minimizes the relative entropy of $Pr(x | Y = 1)$ with respect

235 to $Pr(x)$ (measured using the Kullback-Leibler divergence) is chosen. The distribution of

236 maximum entropy, i.e. closest to the uniform probability distribution or most spread out, is

237 estimated while being subject to a set of constraints imposed by the information available from
238 the environmental conditions where the species occurs.

239 Environmental variables or functions thereof are known as "features" and are treated as an
240 expanded set of variables to be added as terms in the model specification. A random sample of
241 background locations informs the model about $Pr(x)$. The set of constraints on $Pr(x | Y = 1)$
242 ensures that empirical averages of each feature approximate their averages at sites where the
243 species is present (or a random sample thereof).

244 The probability distribution of maximum entropy is a Gibbs distribution, which has an
245 exponential form (Della Pietra et al. 1997). Raw exponential values estimated by the model are
246 scale-dependent, e.g. they can be extremely small if the study area is large, and only represent a
247 measure of relative suitability of each site. However, the model can also be transformed from an
248 exponential family model into a logistic model, thus making it more comparable with other
249 machine learning or generalized linear/additive models (Phillips and Dudik 2008).

250 To calculate the final conditional probability of occurrence $Pr(Y = 1|X = x)$, knowledge of
251 the prevalence of the species $Pr(Y = 1) = \pi$, i.e. the proportion of occupied sites across the
252 study area, is required. However, π is unknown with presence-only data (Ward et al. 2009). In
253 this case, the maximum entropy approach sets this quantity arbitrarily to 0.5.

254

255 **Bayesian for Presence-only Data (BPOD) Approach.**

256

257 When dealing with presence-only data, sampling from the reference population of locations
258 cannot be performed under the traditional random sampling design. Specifically, while a random
259 sample of presences is available, a random sample of absences cannot be obtained. Therefore, a

260 random sample of "contaminated controls", i.e. a random sample of locations from the whole
261 reference population (background sample) that can also include some occurrences of the species,
262 is matched with the aforementioned random sample from the available occurrence data
263 (Lancaster and Imbens 1996).

264 In order to estimate the regression parameters, a two-level scheme is used: (1) a first level
265 describing the probability law that generates the population data; and (2) a second level using a
266 stratified case-control design, modified for presence-only data to select samples from the
267 population. In a traditional logistic regression, the response variable $Y = 0$ marks the absence of
268 an attribute of interest in the population, while $Y = 1$ marks the presence of the same attribute.
269 The key point in the BPOD approach is the introduction of a stratum variable Z , considered as
270 the only observable variable. Specifically, $Z = 0$ means that a location is collected from the
271 whole reference population, while $Z = 1$ indicates that a location is collected from the sub-
272 population of presences. $Z = 1$ implies that $Y = 1$, while $Z = 0$ implies that Y is an unknown value
273 $y \in \{0,1\}$. The introduction of the stratum variable Z allows us to define a linear logistic
274 regression, adjusted for presence-only data. Denoting by $Pr(Z = 1|C = 1, X = x)$ the
275 probability that a location is sampled ($C = 1$) from the set of locations where the species of
276 interest is present ($Z = 1$) and with covariates $X = x$, the linear logistic model for presence-only
277 data can be defined as:

$$\text{logit}Pr(Z = 1|C = 1, X = x) = x\beta + q,$$

278 where q is a correction term, depending on the number of presences truly observed and the
279 unknown number of presences hidden in the sample of "contaminated" controls. An
280 approximation of q can be derived iteratively within the estimation algorithm. After

281 prior distributions are defined over the parameters of interest (the linear coefficients β and the
282 unobserved responses in the sample of "contaminated" controls), Bayesian inference can be
283 carried out through Markov Chain Monte Carlo (MCMC) techniques (Robert and Casella 2004).
284 In particular, an algorithm including a data-augmentation step (Tanner and Wong 1987) is used
285 to obtain an estimate of the unknown empirical prevalence π of the species of interest, jointly
286 with linear coefficients of the logistic model.

287

288 **Evaluation of Model Performance.**

289

290 In this study, model performance is evaluated according to three criteria: (i) prediction
291 accuracy of occurrence data, i.e. model sensitivity expressed by the percentage of correctly
292 predicted occurrences in the sample; (ii) goodness of fit, using both the Akaike Information
293 Criterion (AIC, Akaike 1974) and its corrected version (AICc, Burnham and Anderson 2002);
294 and (iii) ecological realism, i.e. assessing predictions against prior biological knowledge of a
295 species. AIC and AICc for all Maxent models were calculated using the ENMTools (Warren and
296 Seifert 2011) which uses Maxent "raw" suitability scores, i.e. exponential values standardized
297 over the study area. Several other traditional statistical evaluation metrics such as Cohen's Kappa
298 (Cohen 1960) or the area under the receiver operating characteristic curve (AUC, Hanley and
299 Mcneil 1982) are commonly used with presence-absence (or pseudo-absence) data. However, in
300 this study we do not make any assumption of pseudo-absence for background data. While model
301 sensitivity was compared across all models and both statistical approaches, AIC and AICc values
302 were only used to compare the relative quality of each model within the same statistical approach
303 in order to provide a mean for model selection. This is crucial because Maxent's model structure

304 is different from BPOD, hence values of both AIC and AICc cannot be compared across models
305 considered in both approaches.

306

307 **Sampling Scheme.**

308

309 The following background sampling schemes were used with respect to Maxent and BPOD
310 modeling approaches.

311 Each Maxent model was run 16 times, with the background sample size set to 10,000
312 randomly selected points. Although there are not set guidelines regarding the ideal number of
313 background points to use in each situation, some recent studies found that predictive accuracy of
314 Maxent was best with about 10,000 points (Barbet-Massin and Jiguet 2012) over areas
315 comparable in size to our study. Moreover, some studies found that predictive accuracy of
316 Maxent was best with about 10,000 points (Barbet-Massin and Jiguet 2012). All other settings in
317 the MaxEnt software have been used with their default values (Phillips et al. 2006).

318 Each BPOD model was run 500 times, with sample size set according to the
319 presence/background ratio of 1:4, as used by Ward et al. (2009). Specifically, in AST a sample of
320 65 observed presences was matched with a background sample of $65 \times 4 = 260$ locations (total
321 sample size $n=325$), while for FST a sample of 160 observed presences was matched with a
322 background sample of $160 \times 4 = 640$ locations (total sample size $n=800$). The MCMC algorithm
323 with data augmentation used 15,000 iterations (10,000 burn-in) to estimate the unknown model
324 parameters.

325 The reason for using different sampling schemes between Maxent and BPOD is due to the
326 fact that the two approaches have different requirements for reaching robust parameter estimates.

327 Specifically, Maxent needs a large background sample, while BPOD needs a large number of
328 model replications. Given these constraints, we chose model settings accordingly and used
329 roughly the same amount of “sampling information” (see Supp. Table S2-S3 available for the
330 online version). In both approaches, parameter estimates were obtained as averages over all
331 model replications.

332

333 **Results**

334

335 Several models were run to predict the current potential distribution of both AST and FST. A
336 list of the best performing models is shown in Table 2 for AST and Table 3 for FST, together
337 with their evaluation metrics.

338

339 **Table 2—caption at the end of file**

340

341 **Table 3—caption at the end of file**

342

343 For AST, the model that reached the highest overall performance in the maximum entropy
344 approach was M1, while in the BPOD approach it was MPC3, which used the first three
345 principal components as covariates. Fig. 2 (a-b) shows the current potential distributions of AST
346 predicted by the best overall models in both approaches, thus BPOD-MPC3 and Maxent-M1,
347 respectively. Southeastern Florida and the Keys Islands show a much higher suitability compared
348 to other areas, matching the general pattern of recorded occurrences. Low probabilities are also

349 predicted along the east coast up to central Florida and on the west coast around urban areas such
350 as Ft. Myers and Tampa.

351

352 **Figure2–caption at the end of file**

353

354 For FST, the model that reached the highest overall performance in the maximum entropy
355 approach was MPC3, while in the BPOD approach it was MPC6, using the first three principal
356 components and all six principal components as covariates, respectively. Fig. 3 (a-b) shows the
357 current potential distributions of FST predicted by the best overall models in both approaches,
358 thus BPOD-MPC6 and Maxent-MPC3, respectively. Highest suitability values are associated
359 with urbanized areas across the entire state. Although no occurrences were recorded in some
360 urban areas, a medium-to-high suitability is predicted for the species in areas such as North-West
361 Florida around Pensacola, along the west coast in Sarasota and Port Charlotte, along the east
362 coast in Melbourne and Palm Coast, and all the Keys islands south-west of Key Largo.

363

364 **Figure3–caption at the end of file**

365

366 Future predicted probabilities of presence were derived using a model from the BPOD
367 approach for both AST and FST. Due to data availability (see Materials and Methods), the
368 BPOD-M1 model was chosen to predict their future distributions. Fig. 4 (a-b) shows the
369 contemporary predictions calculated using model BPOD-M1 for AST and FST, respectively. A
370 visual inspection suggests that predictions are not much different from the best models that used
371 principal components as covariates, with the exception of a few areas for FST such as the Keys

372 Islands or the west coast of Florida where the suitability is slightly lower. Fig. 4 (c-d) shows
373 average "consensus" projected probabilities, i.e. averaged over the three GCMs, for AST and
374 FST, respectively, under climate change conditions for the 2050s time period and given no
375 change in land cover. The climate variables that are projected to the future from M1 are
376 precipitation, bio5, and bio6. Urban areas in southeast Florida seem to have an increased
377 predicted probability of presence for AST, while for FST changes in suitability are less
378 noticeable. The population growth scenario (Fig. 4 e-f) increases the percentage of areal units
379 occupied by developed areas, thus increasing the variable "devel" (refer to Table 1) in our model.
380 The effect of a combined change in climate and developed areas increases the predicted
381 probabilities of presence for both AST and FST. However, the effect is much more noticeable for
382 the latter across the whole study area.

383

384 **Figure4–caption at the end of file**

385

386 **Discussion**

387

388 The performance of the BPOD approach on both species was shown to be consistently better
389 than the widely used maximum entropy method, with a few exceptions, in terms of sampling
390 sensitivity (see Table 3). Whenever the model covariates were highly informative on a species
391 geographical distribution (e.g., for AST), the BPOD approach performed consistently better than
392 maximum entropy. In fact, the highest sensitivity reached by any Maxent model was 61%, hence
393 lower than the worst BPOD model (76%). When the model covariates are less informative for
394 predicting distribution, as for the FST, BPOD performs better than MaxEnt in 78% of the cases.

395 Finally, the best BPOD model gave more realistic predictions from an ecological perspective
396 compared to the best Maxent model for both species. Specifically, for FST maximum entropy
397 tends to over-predict areas across the entire state, far apart from recorded occurrences, and
398 under-predict areas close to them (Fig. 3). Although this phenomenon is less pronounced for
399 AST, areas in the metropolitan southeast Florida are under-predicted nearby recorded
400 occurrences.

401 The BPOD approach makes a better use of the information from PCA-derived variables
402 compared to maximum entropy, as its predictive power increases until reaching an optimum in
403 terms of sensitivity and both information criteria (see Tables 2-3). However, such models behave
404 in a slightly different manner between the two species. In particular, BPOD models for AST
405 reach an optimum with a smaller number of PCA-derived variables than FST (3 vs. 6 principal
406 components, respectively). This probably means that the original environmental variables
407 enclosed in the first three principal components are sufficient to explain the ecological niche of
408 AST in Florida. FST, tolerating broader climatic and environmental gradients than AST, has
409 attained generic species status in Florida where it occurs in all major human population centers
410 of the State. This result also suggests that some environmental factors influencing the habitat of
411 FST may be missing from these analyses.

412 Maxent models reported in this paper were estimated by fitting linear responses to
413 relationships between response and predictor variables in order to keep comparability between
414 the two different statistical approaches. Maxent models fitting more complex responses were
415 also tested but had a much lower predictive performance compared to the ones fitting linear
416 features. A major advantage of the BPOD approach over maximum entropy is that the MCMC
417 algorithm does not require the *a priori* knowledge of the population prevalence as it is

418 considered as a further parameter in the model. This overcomes the issue of prevalence
419 specification pointed out by Ward et al. (2009). A Bayesian modeling framework allows
420 flexibility in the treatment of uncertainty while making full inference on the probability of
421 presence possible. However, a more formal comparison between the two statistical approaches
422 based on artificial data is suggested for future studies.

423 In this paper, statistically downscaled climate projections for the 2050s were preferred over
424 dynamically downscaled projections, such as the CLARReS10 dataset for the Southeast United
425 States (Stefanova et al. 2012). Although the latter are able to incorporate regional-scale
426 processes, their spatial resolution (~10 km) was too coarse to assess the effect of variation in
427 climate and urbanization on the same scale used for contemporary predictions for both termite
428 species. Climate change under the A2 scenario for the 2050s has a moderate effect on both
429 species' geographical distribution. Conversely, a combined effect of climate change with a
430 population growth scenario has a larger impact on their projected probabilities, especially for
431 FST. This suggests that both termite species are influenced by urban development much more
432 than by climate alone.

433 Two issues not fully addressed in this work are the residual autocorrelation that may still
434 persist among neighboring occurrences and the problem of observer bias (Syfert et al. 2013). In
435 order to reduce spatial autocorrelation, we chose a spatial resolution at which termite occurrences
436 can be assumed independent of each other given the explanatory variables (see Materials and
437 Methods). Spatially explicit models, i.e. models with spatial autoregressive component (Cressie
438 1993) or latent spatially structured component (Zuur et al. 2009), might be available to refine our
439 final predictions. However, a reasonable way of generating pseudo-absences must be found and
440 these models are computationally intensive to estimate. The issue of observer bias would be hard

441 to address in the models developed herein because the data comes from different sources and
442 involves multiple data collectors (see Materials and Methods). All samples were not collected
443 using road accessibility criteria, hence standard solutions, e.g. adding information on road
444 distance within the models (Phillips and Dudik 2008), could not be implemented in this study.
445 The treatment of such a complex issue is deferred to future work.

446

447 **Acknowledgements**

448

449 The authors would like to thank both anonymous reviewers for their valuable comments and
450 suggestions to improve the quality of the paper. We would like to thank David Bucklin and
451 James Watling from the University of Florida for their valuable help and suggestions throughout
452 the development of this work. The authors would also like to thank Lydia Stefanova from the
453 Center for Oceanic-Atmospheric Prediction Studies at Florida State University for her valuable
454 input on downscaled climate projections. The MCMC algorithm used to estimate parameters in
455 the Bayesian logistic regression approach was developed in FORTRAN by Dr. Fabio Divino (co-
456 author, email: fabio.divino@unimol.it).

457

458

459

460

461

462

463

464 **References**

- 465
- 466
- 467 Cohen, J. 1960. A coefficient of agreement for nominal scales. *Educ. Psychol. Meas.* 20: 37–46.
- 468 Farrar, D. E., and R. R. Glauber. 1967. Multicollinearity in regression analysis: the problem
469 revisited. *Rev. Econ. Stat.* 49: 92-107.
- 470 Nutting, W. L. 1969. Flight and colony foundation. *In* K. Krishna, F. M. Weesner (eds.), *Biology*
471 *of termites*. Academic Press, New York, 233-282.
- 472 Akaike, H. 1974. A new look at the statistical model identification. *IEEE Transactions on*
473 *Automatic Control* 19: 716–723.
- 474 Hanley, J. A., and B. J. McNeil. 1982. The meaning and use of the area under a receiver
475 operating characteristic (ROC) curve. *Radiol.* 143: 29-36.
- 476 Tanner, M., and W. Wong. 1987. The calculation of posterior distribution by data augmentation.
477 *J. Am. Stat. Assoc.* 82: 528-550.
- 478 Cressie, N. 1993. *Statistics for spatial data*. John Wiley & Sons, Inc., New York.
- 479 Lancaster, T., and G. Imbens. 1996. Case-control studies with contaminated controls. *J. Econ.*
480 *71: 145-160.*
- 481 Della Pietra, S., V. Della Pietra, and J. Lafferty. 1997. Inducing features of random fields. *IEEE*
482 *Trans. Pattern Anal. Mach. Intell.* 19: 380-393.
- 483 Burnham, K. P., and D. R. Anderson. 2002. *Model selection and multimodel inference: a*
484 *practical information-theoretic approach*, 2nd ed. Springer-Verlag.
- 485 Daly, C., W. P. Gibson, G. H. Taylor, G. L. Johnson, and P. Pasteris. 2002. A knowledge-based
486 approach to the statistical mapping of climate. *Clim. Res.* 22: 99-113.
- 487 Phillips, S. J., M. Dudík, and R. E. Schapire. 2004. A maximum entropy approach to species
488 distribution modeling. *Proc. 21st Int. Conf. Mach. Learn.* 655-662.
- 489 Robert, C. P., and G. Casella. 2004. *Monte Carlo statistical methods*, 2nd ed. Springer-Verlag
490 New York, Inc., Secaucus, NJ.
- 491 Guisan, A., and W. Thuiller. 2005. Predicting species distribution: offering more than simple
492 habitat models. *Ecol. Lett.* 8: 993-1009.
- 493 Hijmans, R. J., S. E. Cameron, J. L. Parra, P. G. Jones, and A. Jarvis. 2005. Very high resolution
494 interpolated climate surfaces for global land areas. *Int. J. Clim.* 25: 1965-1978.
- 495 Elith, J., C. H. Graham, R. P. Anderson, M. Dudík, S. Ferrier, A. Guisan, R. Hijmans, F.
496 Huettmann, J. R. Leathwick, A. Lehmann, J. Li, L. G. Lohmann, B. A. Loiselle, G.
497 Manion, C. Moritz, M. Nakamura, Y. Nakazawa, J. M. Overton, A. T. Peterson, S. J.
498 Phillips, K. Richardson, R. Schachetti-Pereira, R. E. Schapire, J. Soberon, S. Williams,
499 M. S. Wisz, and N. E. Zimmermann. 2006. Novel methods improve prediction of species'
500 distributions from occurrence data. *Ecogr.* 29: 129-151.
- 501 Phillips, S. J., R. P. Anderson, and R. E. Schapire. 2006. Maximum entropy modeling of species
502 geographic distributions. *Ecol. Model.* 190: 231-259.
- 503 Elith, J., and J. Leathwick. 2007. Predicting species distributions from museum and herbarium
504 records using multiresponse models fitted with multivariate adaptive regression splines.
505 *Divers. Distrib.* 13: 265-275.
- 506 Intergovernmental Panel on Climate Change. 2007. *IPCC fourth assessment report (AR4)*. Url:
507 <http://goo.gl/3C2amF>. Accessed: 03/24/2014
- 508 Phillips, S. J., and M. Dudik. 2008. Modeling of species distributions with maxent: new
509 extensions and a comprehensive evaluation. *Ecogr.* 31: 161-175.

510 Diniz-Filho, J. A., L. M. Bini, T. F. Rangel, R. D. Loyola, C. Hof, D. Nogués-Bravo, and M. B.
511 Araújo. 2009. Partitioning and mapping uncertainties in ensembles of forecasts of species
512 turnover under climate change. *Ecogr.* 32: 897-906.

513 Elith, J., and J. R. Leathwick. 2009. Species distribution models: ecological explanation and
514 prediction across space and time. *Annu. Rev. Ecol. Evol. Syst.* 40: 677-697.

515 Franklin, J. 2009. Mapping species distributions: spatial inference and prediction. Cambridge
516 University Press, Cambridge, UK.

517 Li, H.-F., W. Ye, N.-Y. Su, and N. Kanzaki. 2009. Phylogeography of *Coptotermes gestroi* and
518 *Coptotermes formosanus* (Isoptera: Rhinotermitidae) in Taiwan. *Ann. Entomol. Soc. Am.*
519 102: 684-693.

520 Marmion, M., M. Parviainen, M. Luoto, R. K. Heikkinen, and W. Thuiller. 2009. Evaluation of
521 consensus methods in predictive species distribution modelling. *Divers. Distrib.* 15: 59-
522 69.

523 Ward, G., T. Hastie, S. C. Barry, J. Elith, and J. R. Leathwick. 2009. Presence-only data and the
524 EM algorithm. *Biom.* 65: 554-563.

525 Zuur, A., E. N. Ieno, N. Walker, A. A. Saveliev, and G. M. Smith. 2009. Mixed effect models
526 and extensions in ecology with R. Springer.

527 Li, H.-F., R. L. Yang, and N. Y. Su. 2010. Interspecific competition and territory defense
528 mechanisms of *Coptotermes formosanus* and *Coptotermes gestroi* (Isoptera:
529 Rhinotermitidae). *Env. Entomol.* 39: 1601-1607.

530 Ramirez-Villegas, J., and A. Jarvis. 2010. Downscaling global circulation model outputs: the
531 delta method. *Decis. Policy Anal. Working Pap.* No. 1.

532 Warton, D. I., and L. C. Sheperd. 2010. Poisson point process models solve the “pseudo-absence
533 problem” for presence-only data in ecology. *Ann. Appl. Stat.* 4: 1383-1402.

534 Chakraborty, A., A. E. Gelfand, A. M. Wilson, A. M. Latimer, and J. A. Silander. 2011. Point
535 pattern modelling for degraded presence-only data over large regions. *J. R. Stat. Soc.:*
536 *Series C (App. Stat.)* 5: 757-776.

537 Divino, F., N. Golini, G. J. Lasinio, and A. Penttinen. 2011. Data augmentation approach in
538 Bayesian modelling of presence-only data. *Procedia Environ. Sci.* 7: 38-43.

539 Elith, J., S. J. Phillips, T. Hastie, M. Dudík, Y. E. Chee, and C. J. Yates. 2011. A statistical
540 explanation of maxent for ecologists. *Divers. Distrib.* 17: 43-57.

541 Gautam, B. K., and G. Henderson. 2011. Relative humidity preference and survival of starved
542 Formosan subterranean termites (Isoptera: Rhinotermitidae) at various temperature and
543 relative humidity conditions. *Env. Entomol.* 40: 1232-1238.

544 Scheffrahn, R. H., and W. Crowe. 2011. Ship-borne termite (Isoptera) border interceptions in
545 Australia and onboard infestations in Florida, 1986–2009. *Fla. Entomol.* 94: 57-63.

546 Warren, D. L., and S. Seifert. 2011. Environmental niche modeling in maxent: the importance of
547 model complexity and the performance of model selection criteria. *Ecol. Appl.* 21: 335-
548 342.

549 Barbet-Massin, M., and F. Jiguet. 2012. Selecting pseudo-absences for species distribution
550 models: How, where and how many? *Methods Ecol. Evol.* 3: 327-338.

551 Multi-Resolution Land Characteristics Consortium, 2012. National land cover database 2006.
552 <http://www.mrlc.gov/nlcd2006.php>.

553 Stefanova, L., V. Misra, S. Chan, M. Griffin, J. J. O'Brien, and T. J. I. Smith. 2012. A proxy for
554 high-resolution regional reanalysis for the southeast united states: Assessment of

555 precipitation variability in dynamically downscaled reanalyses. *Climate Dynamics* 38:
556 2449-2466.

557 Climate Change Agriculture and Food Security. 2013. GCM downscaled data portal. Url:
558 <http://www.ccafs-climate.org/data/>. Accessed: 08/21/2013

559 Divino, F., N. Golini, G. Jona Lasinio, and A. Penttinen. 2013. Bayesian modeling and MCMC
560 computation in linear logistic regression for presence-only data. (submitted). Cornell
561 University Library. Url: <http://arxiv.org/abs/1305.1232>. Accessed: 08/21/2013

562 Evans, T. A., B. T. Forschler, and J. K. Grace. 2013. Biology of invasive termites: a worldwide
563 review. *Ann. Rev. Entomol.* 58: 455-474.

564 Li, H.-F., I. Fujisaki, and N.-Y. Su. 2013. Predicting habitat suitability of *Coptotermes gestroi*
565 (Isoptera: Rhinotermitidae) with species distribution models. *J. Econ. Entomol.* 106: 311-
566 321.

567 Scheffrahn, R. H. 2013. Overview and current status of non-native termites (Isoptera) in Florida.
568 *Fla. Entomol.* (in press).

569 Syfert, M. M., M. J. Smith, and D. A. Coomes. 2013. The effect of sampling bias and model
570 complexity on the predictive performance of maxent species distribution models. *PLoS*
571 *ONE* 8.

572 University of Florida GeoPlan Center. 2013. Florida projected population growth - 2060 Url:
573 <http://www.fgdl.org>. Accessed: 08/21/2013

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589 **Table 1.** Climatic and environmental variables used and their data source for both AST and FST.

Database	Variable	Description
PRISM	prec	Annual total precipitation
	dew	Average daily mean dew point temperature
	tmax	Average daily maximum temperature
	tmin	Average daily minimum temperature
WORLDCLIM	bio5	Maximum temperature of the warmest month
	bio6	Minimum temperature of the coldest month
NLCD 2006	water	Open water or permanent ice/snow cover
	devel	High percentage ($\geq 30\%$) of constructed materials (e.g. asphalt, concrete, buildings, etc.).
	barren	Bare rock, gravel, sand, silt, clay, or other earthen material, with little or no "green" vegetation
	forest	Tree cover (> 6 m tall). Tree canopy accounts for 25% to 100% of the cover
	shrub	Natural/semi-natural woody vegetation with aerial stems (< 6 m tall)
	herb	Natural/semi-natural herbaceous vegetation (75% - 100% of the cover)
	cultiv	Herbaceous vegetation that has been planted or is intensively managed for the production of food, feed, or fiber (75% - 100% of the cover)
	wetlands	Soil or substrate is periodically saturated with or covered with water

590

591

592 **Table 2.** List of main models used for AST, their sampling sensitivity, and information criteria. M1: X
 593 (easting), prec, bio5, bio6, all land cover variables. M2: prec, bio5, bio6, all land cover variables. M3: X
 594 (easting), Y (northing), prec, bio5, all land cover variables. MPCx: *x* stands for the number of principal
 595 components used as covariates. The best models are highlighted in bold.

Approach	Model	Sampling Sensitivity	AIC	AICc
	M1	0.96	62.2	38.2
	M2	0.96	63.2	38.4
	M3	0.96	67.7	41.1
	MPC1	0.76	118.3	112.4
BPOD	MPC2	0.90	78	70.2
	MPC3	0.97	48.8	39
	MPC4	0.97	50.1	38.3
	MPC5	0.97	51.4	37.8
	MPC6	0.97	52.5	36.9
	M1	0.61	823.4	825.2
	M2	0.61	828.7	830.6
	M3	0.60	857.6	860
	MPC1	0.40	1066	1066
Maxent	MPC2	0.61	942	942
	MPC3	0.60	828.8	829.2
	MPC4	0.60	829.1	829.8
	MPC5	0.58	831.7	832.5
	MPC6	0.57	830.2	831.1

596

597 **Table 3.** List of main models used for FST, their sampling sensitivity, and information criteria. M1: X
 598 (easting), prec, bio5, bio6, all land cover variables. M2: prec, bio5, bio6, all land cover variables. M3: X
 599 (easting), Y (northing), prec, bio5, all land cover variables. MPCx: *x* stands for the number of principal
 600 components used as covariates. The best models are highlighted in bold.

Approach	Model	Sampling Sensitivity	AIC	AICc
	M1	0.73	391.1	366
	M2	0.73	391.4	367
	M3	0.73	391.5	367.2
	MPC1	0.05	689.2	683.2
BPOD	MPC2	0.53	529.7	521.8
	MPC3	0.71	396.3	386.4
	MPC4	0.71	395.1	383.2
	MPC5	0.72	388.6	374.8
	MPC6	0.74	382.8	367
	M1	0.55	2712.9	2714.2
	M2	0.55	2713.1	2714
	M3	0.57	2712	2713.2
	MPC1	0.66	1177	1177.1
Maxent	MPC2	0.59	1047.4	1047.6
	MPC3	0.66	956.8	957.2
	MPC4	0.57	968.5	969.1
	MPC5	0.61	963.6	964.6
	MPC6	0.58	961.1	962.5

601

602

603 **Figure Captions:**

604

605 **Fig. 1.** Florida occurrences of AST (green) and FST (purple). Available in color online.

606

607 **Fig. 2.** Current predicted probabilities of presence for AST. (a) BPOD-MPC3 model. (b)

608 Maxent-M1 model. Darker red areas correspond to areas with higher probabilities. Available in

609 color online.

610

611 **Fig. 3.** Current predicted probabilities of presence for FST. (a) BPOD-MPC6 model. (b) Maxent-

612 MPC3 model. Darker red areas correspond to areas with higher probabilities. Available in color

613 online.

614

615 **Fig. 4.** Current and average projected probabilities of presence for the 2050s time period. (a)

616 BPOD-M1 contemporary predictions for AST. (b) BPOD-M1 contemporary predictions for FST.

617 (c) BPOD-M1 projected predictions for AST averaged over the GFDL-CM 2.1, NCAR-CCSM

618 3.0, and UKMO-HadCM3 global circulation models under the A2 emission scenario. (d) BPOD-

619 M1 projected predictions for FST averaged over the GFDL-CM 2.1, NCAR-CCSM 3.0, and

620 UKMO-HadCM3 global circulation models under the A2 emission scenario. (e) BPOD-M1

621 projected predictions for AST averaged over the GFDL-CM 2.1, NCAR-CCSM 3.0, and

622 UKMO-HadCM3 global circulation models under the A2 emission scenario and population

623 growth. (f) BPOD-M1 projected predictions for FST averaged over the GFDL-CM 2.1, NCAR-

624 CCSM 3.0, and UKMO-HadCM3 global circulation models under the A2 emission scenario and

625 population growth. Darker red areas correspond to areas with higher probabilities. Available in
626 color online.

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

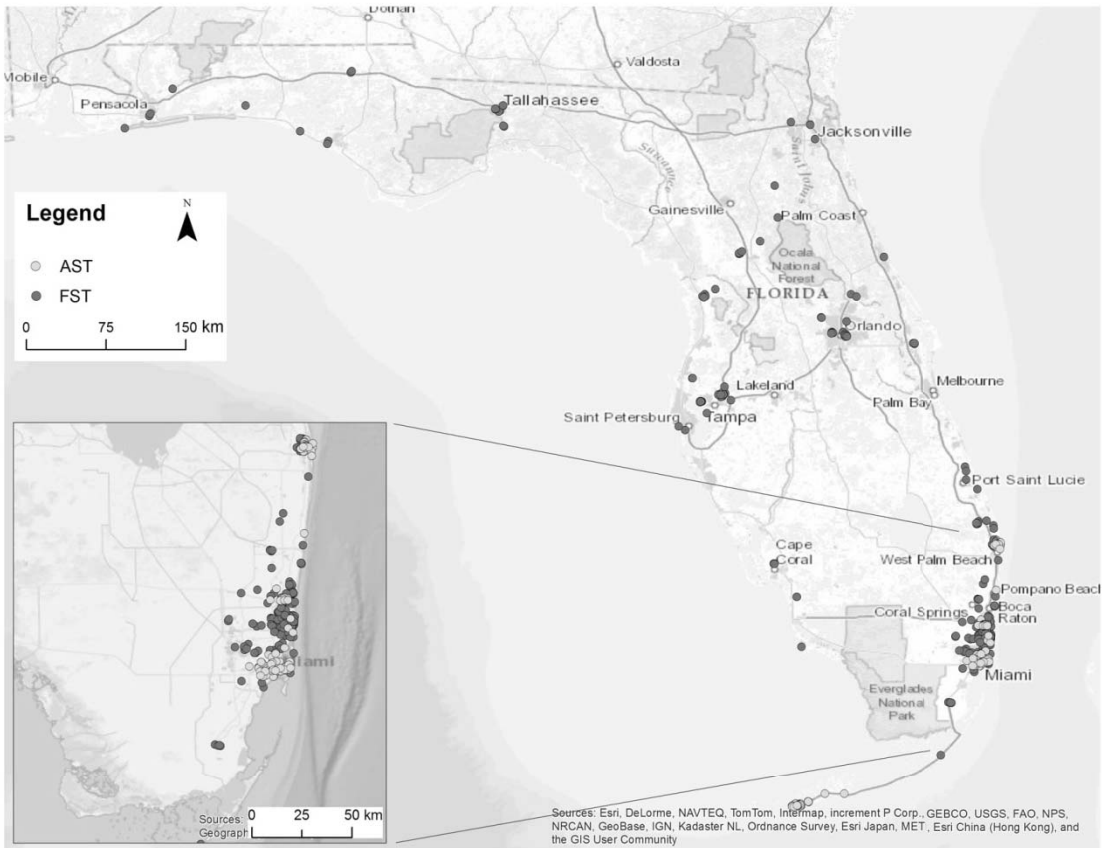
645

646

647

648

FIGURE 1



649

650

651

652

653

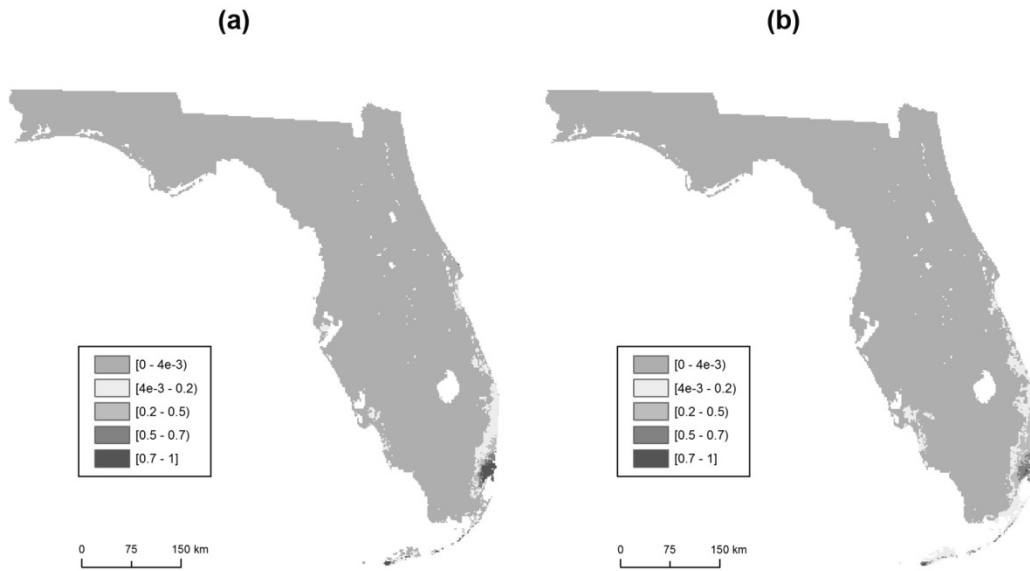
654

655

656

657

FIGURE 2



659

660

661

662

663

664

665

666

667

668

669

670

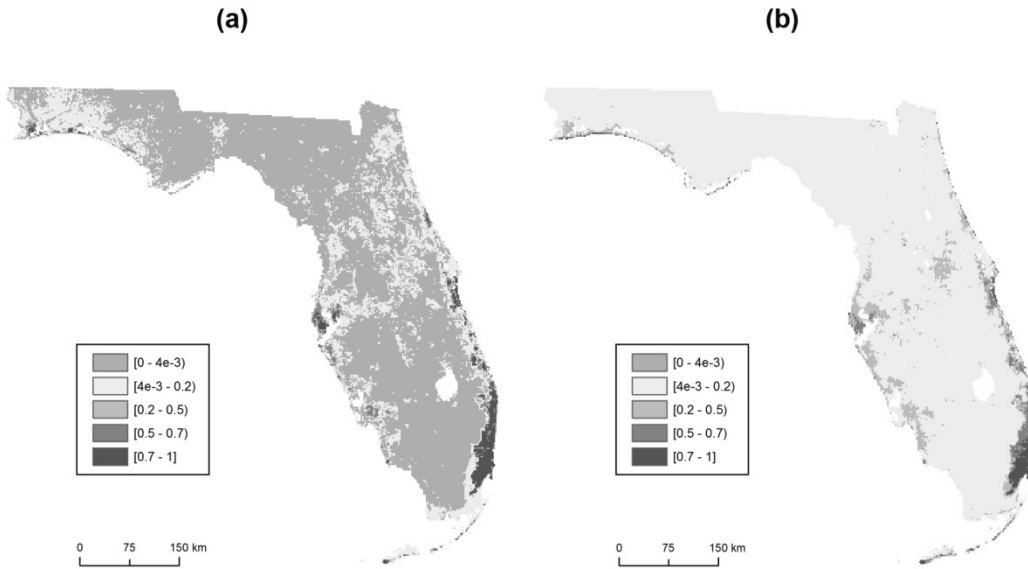
671

672

673

674

FIGURE 3



675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

FIGURE 4

