



## Report Cover Page

<b>ACERA Project</b>		
ARC Project Spatial Methods		
<b>Title</b>		
Reasons for differing performances of species distribution models		
<b>Author(s) / Address (es)</b>		
Jane Elith, Botany, ACERA, University of Melbourne and Catherine Graham, State University of New York at Stony Brook		
<b>Material Type and Status (Internal draft, Final Technical or Project report, Manuscript, Manual, Software)</b>		
Project final report #1		
<b>Summary</b>		
<p>The aim of this project was to develop and evaluate methods to improve the prediction of species in geographic space. For some applications such as climate change and invasive species, these methods must be used to make predictions outside the ranges of the data used to build the models, making validation even more important. New methods have arisen in statistics and machine learning to deal with these issues.</p> <p>This study uses simulated data to compare the performance of modeling methods. It concludes that some of new machine learning methods, in particular, boosted regression trees and maximum entropy methods, outperform statistical and climate matching methods under a number of performance measures, for a wide range of data types, data qualities and prediction circumstances. More importantly, the paper provides advice on how a user might evaluate the performance of predictive models, to make best use of the alternatives available to them for specific purposes.</p> <p>This paper provides details of how to access software and statistical code to implement several of the better behaved and most promising methods.</p>		
ACERA Use only	Received By:	Date:
ACERA Use only	ACERA / AMSI SAC Approval:	Date:
ACERA Use only	ACERA / AMSI SAC Approval:	Date:

1 **Title: Do they / how do they / WHY do they differ? -- on finding reasons for differing**  
2 **performances of species distribution models.**

3

4 Authors: Jane Elith<sup>1</sup> and Catherine Graham<sup>2</sup>

5

6 <sup>1</sup> School of Botany

7 The University of Melbourne

8 Parkville, Victoria. 3010

9 Australia

10 [j.elith@unimelb.edu.au](mailto:j.elith@unimelb.edu.au)

11 (corresponding author)

12

13

14 <sup>2</sup>Department of Ecology and Evolution

15 650 Life Sciences Building

16 Stony Brook University

17 Stony Brook, NY 11794, USA

18

19

20

21

22

23

24

25

26

27

28

## 29 **Introduction**

30 Species distribution models (SDMs) are increasingly being used to address a diverse range of  
31 applied and theoretical questions (Graham *et al.* 2004, Guisan and Thuiller 2005). Their expanding  
32 use means that models are now being fitted to new forms of data, including occurrence records from  
33 museums or herbaria. For some applications, such as climate change or invasive species research,  
34 model predictions are extended beyond the geographic or environmental region from which training  
35 samples were drawn (e.g., Araújo *et al.* 2005). Newer applications also move beyond the traditional  
36 ecological focus of SDM into new fields such as evolutionary biology, where they are used to study  
37 topics such as speciation or hybrid zones (e.g., Kozak *et al.* 2008). As a result of these ongoing  
38 changes in the field, and spurred on by advances in data analysis within statistics and computing  
39 sciences (Breiman 2001a), new modelling methods continue to be implemented. Model complexity  
40 has generally increased over time from simple environmental matching (e.g. BIOCLIM, Busby 1991;  
41 DOMAIN, Carpenter *et al.* 1993) to fitting more complex non-linear relationships between species  
42 presence and the environment (e.g., generalised additive models, GAM, Hastie and Tibshirani 1990,  
43 Yee and Mitchell 1993; and maximum entropy modelling, Maxent, Phillips *et al.* 2006). Recent  
44 emphases on machine-learning and Bayesian methods indicate that new methods will continue to be  
45 developed (Prasad *et al.* 2006, Latimer *et al.* 2006).

46 This wide array of methods, data types and diverse research questions imply different  
47 requirements of modelling methods, creating the need for choice of method to be matched to  
48 application. It has spurred a growing literature aimed at evaluating different modelling methods in  
49 varied contexts. Numerous studies compare how well different methods perform (e.g., Elith *et al.*  
50 2006, Seguardo and Araújo 2004, Moisen and Frescino 2002). Generally, in these types of studies a  
51 set of distributions are predicted for one or more species using a range of methods or different ways of  
52 applying one method, and then predictive performance is tested, usually against withheld or  
53 independently collected data (i.e., independent presence/absence data, additional field surveys, fossils  
54 etc.). This type of approach has also been used to explore the relationship between model  
55 performance and data characteristics such as: the number of species localities for modelling (eg.

56 Hernandez *et al.* 2006, McPherson and Jetz 2007a, Reese *et al.* 2005), species-level attributes (e.g.,  
57 Guisan *et al.* 2007, McPherson and Jetz 2007b), and quality of locality data (e.g., Loiselle *et al.* 2008,  
58 Graham *et al.* 2008). While this extensive effort has yielded new insights and advice for modellers (a  
59 review beyond the scope of this paper), there is still considerable debate as to what methods are useful  
60 for which applications, as evidenced in diverse views and conflicting conclusions in recent papers  
61 comparing methods (see recent opinions in *Ecography* by Peterson *et al.* 2007, Phillips 2008).

62 Additional method comparison studies using the same basic approach of evaluating how well a  
63 series of methods can recover the geographic distribution of a species will likely provide diminishing  
64 returns because what method performs best, or is most appropriate, can vary by study system, data  
65 quality, the way the method is used, or the question addressed. This makes it hard in most cases to  
66 learn general principles about the models and to garner the right sort of knowledge to inform choice of  
67 method. The problem is that, unless there are large amounts of appropriate data available, it is difficult  
68 to separate the effects of the particular data from more general questions about the methods and their  
69 performance. In our opinion, much more progress would be made in understanding the keys links  
70 between model and outcome if model comparison studies focussed more on determining why a given  
71 method performs as it does.

72 This demands a wider array of evaluation criteria. There are several possibilities, and here we  
73 mention a few. First, one could approach the problem by learning how a method works. This is  
74 usually the domain of those who develop the methods, commonly statisticians and computer science  
75 professionals, but practitioners would certainly benefit from such understanding. Given knowledge of  
76 how methods work, one could decide *a priori* which method should be most appropriate given the  
77 question being addressed and the biology of the species. Alternatively, a researcher could use a large  
78 number of settings for a given method and conduct a very thorough evaluation of model performance  
79 using an comprehensive range of performance measures (for useful measures see, for example, Pearce  
80 and Ferrier 2000, and Caruana and Niculescu-Mizil 2006). This can shed light on the model's  
81 behaviour and limitations under various settings, but will not lead to understanding of why the  
82 differences in settings lead to differing performance in the absence of knowledge of how the method

83 works. However it is an important part of using a method well. Effects of varying settings on model  
84 performance are rarely reported in ecological applications of SDM (but see Phillips and Dudik 2008)  
85 yet the information is crucial to identifying the sensitivities and tendencies of methods.

86 Another approach is to use data with known characteristics to test model fit and prediction.  
87 These can be sets of data with known properties (e.g. drawn from a normal distribution with mean  $x$   
88 and standard deviation  $y$ ; see for example Moisen and Frescino 2002) or simulated (artificial) species  
89 (Austin *et al.* 2006). Simulated species are particularly useful if relevant features such as population  
90 processes, competition, or typical landscape properties can be included (e.g. Kearney *et al.* in press,  
91 Tyre *et al.* 2001). In examples of more static correlation-based approaches, researchers define the  
92 true distribution of the species or the relationship between the species and its environment (e.g.  
93 Meynard and Quinn 2007; Reineking and Schröder 2006). The features of a simulated species can be  
94 matched to the aspects of model performance that are important to the researcher.

95 The simulated data approach is particularly important for methods where the structure of the  
96 model itself is unclear. We envisage this will be an increasingly relevant issue as more machine  
97 learning methods are appropriated for ecological prediction tasks. This is because, unlike statistical  
98 models that have underlying and clear assumptions about the data and the processes that generated it,  
99 machine learning and other data mining methods approach the data-generating mechanism (nature) as  
100 unknown, but learn patterns from the observed data to make predictions (Breiman 2001a). These  
101 models are not always easy to interpret. In Breiman's words: "the models that best emulate nature in  
102 terms of predictive accuracy are also the most complex and inscrutable". Ecologists are unlikely to  
103 want inscrutable models; we want and sometimes need to understand what is behind the predictions.  
104 In the machine learning community there is some focus on how to interpret and visualise model  
105 output, and statisticians are reinterpreting some methods from a statistical viewpoint. This is  
106 particularly useful for ecologists (e.g. Cutler *et al.* 2007, Elith *et al.* 2008). Modelling simulated data  
107 provides another opportunity to explore how new methods behave.

108 In summary, we believe that deeper insights into the causes of varying model performance  
109 require an expansion of model evaluation approaches (Araújo and Guisan 2005). Comparative studies

110 on one or two regions and a few species do not tend to give insight into why methods differ, and  
111 without that knowledge it is difficult to make well informed selection of methods. Perhaps we should  
112 de-couple the method comparison from the interesting questions about species and biodiversity.  
113 Method comparisons need to focus on approaches that give insight to why methods differ. A range of  
114 tests is required to determine whether the model does what the application requires.

115 To demonstrate these ideas, here we use a simulated species approach to compare methods. We  
116 pose three applications: (1) understand the relationships between the species and its environment; (2)  
117 predict which parts of the landscape are more or less suitable for the species by creating a map of  
118 relative suitabilities; (3) extrapolate to environmental conditions outside those in the sample space. In  
119 doing so we aim to develop a straightforward example of the type of research that helps to move  
120 beyond standard model comparisons to those where the emphasis is on determining why methods  
121 differ. Our underlying premise is that it is useful to think more about the reasons for variation in  
122 model performance, and to match the evaluation criteria (whether a simulation, particular subsets of  
123 data, certain statistics, or map comparisons) to the questions that the models are intending to answer.

#### 124 ***1. Simulating the species and sampling it.***

125 We created a simulated species and mapped it onto a landscape. The species responds to three  
126 variables: *wetness*, aspect ("*southness*") and *geology*. It represents a southern hemisphere plant  
127 preferring wet and south-facing (shaded) sites and fertile substrates, with an interaction that results in  
128 *wetness* dominating the response (equation 1 and Figure 1, top panel). Details of how each index  
129 (SI.*wetness* etc) was specified is presented in the online supplement, Appendix S1.

130 
$$\text{Suitability of a cell} = \text{SI.wetness} * 0.5 * (\text{SI.southness} + \text{SI.geology}) \quad \text{- eq'n 1}$$

131 where SI = suitability index (see Appendix)

132 The predictor variables (*wetness*, *southness* and *geology*) existed as real mapped data; the  
133 species' response to them was invented. The mapped distribution based on these relationships shows  
134 the suitability (scaled 0 to 1) of each grid cell for the plant (Figure 2, top left). We chose to use  
135 models that use presence-only or presence-absence data, so needed to make a presence-absence

136 realisation of the species. We used the suitability value in each grid cell of the map as the success rate  
137 for one sample of the binomial distribution, i.e., a cell with a suitability of 0.6 has a 60% chance of  
138 being occupied. This uses the "rbinom" function in R (R Development Core Team 2006) and is  
139 preferable to choosing a threshold. A threshold would alter the response of the species to one that has  
140 vertical cut-offs of suitability, not gradual changes as in our "truth". The idea of the binomial  
141 realisation is consistent with an interpretation of species relationships that says: "if the environment is  
142 only partially suitable for the species, in some cases the species will occur there and in some it will  
143 not". Our presence-absence realisation of the simulated species occupied 12% of the cells in the study  
144 region (Figure S2a, online supplement).

145       To provide the data set for modelling we sampled from the simulated presence-absence  
146 distribution by taking 1000 sites at random; these contained 115 presences and 885 known absences  
147 and will from here on be called the PA sample. This number of sites is somewhat arbitrary, but is  
148 enough to fit models well, and is typical of the amount of data often used in modelling. Other  
149 protocols such as stratified sampling across environmental gradients could be used instead of random  
150 sampling, but our aim here is only to demonstrate the insights from this type of model comparison,  
151 not to test sensitivities to data samples. The PA sample is visualised in Figure 3a; consistent with the  
152 full data set (Figure 3b), samples at high *wetness* are relatively rare. For methods only requiring  
153 presence data, the presence records in the PA sample were used as presence-only (PO) data. In the  
154 online supplement (Appendix S5) we describe additional data samples, in which 1000 or 3000  
155 pseudo-absences were selected and combined with the PO data, to test the effect of using presences  
156 with pseudo-absences instead of true absence data.

## 157 **2. Modelling and evaluation methods**

158       Given that our intent is to demonstrate an approach, not to conduct an exhaustive study, we  
159 included only five algorithms. For the PA data we selected a generalised linear model (GLM;  
160 McCullagh and Nelder 1989) as a standard regression model and two more recent methods, boosted  
161 regression trees (BRT; Friedman *et al.* 2000) and random forests (RF; Breiman 2001b). For the PO  
162 data, we used Maxent and GARP as featured in recent comparisons (Peterson *et al.* 2007, Phillips

163 2008). Maxent, GARP and RF are machine learning (ML) methods, BRT is a ML method  
164 reinterpreted into a statistical paradigm, and GLM is a statistical method. More details and references  
165 are given in the online supplement. These methods were selected because they allow a number of  
166 contrasts, including comparison of models developed on PA compared with PO data, comparison of  
167 Maxent and GARP, and comparison of two methods using ensembles of trees, BRT and RF. Each of  
168 the methods are assumed capable for at least some of our three posed applications, given that they  
169 have all been used for related tasks.

170 For details of model building methods see Table 1 and Appendices S2, S3 and S4 in the online  
171 supplement. In those appendices the settings for some methods are tested in detail; this reflects our  
172 need to explore the effect of a range of settings for some methods, either because we were relatively  
173 inexperienced with that method or because the recommended settings did not produce good results.

174 Because one of our applications required output showing modelled relationships and not all of  
175 the methods provide visualization of fitted functions, we created a second set of environmental grids  
176 with an evaluation strip inserted, following Elith *et al.* (2005). The strip is a simple data arrangement  
177 that holds two of the three environmental variables at a constant value (here, a value achieving  
178 maximum suitability for the species) whilst varying the value of the third environmental variable  
179 across its numerical range. To visualise interactions (and to check that results do not depend on the  
180 value at which variables are held constant), pairs of variables can be covaried. For this, we included a  
181 comprehensive set of combinations of *wetness* and *southness* for each of the four classes of *geology*.  
182 Predictions were made to the evaluation strip in its tabled form, for those methods that could predict  
183 to that, or by making a "projection" (Maxent and GARP) to the second set of grids. This process  
184 allowed us to determine the partial response to one or two variables while holding other variables  
185 constant. It is analogous to the partial plots from standard regression methods (see Elith *et al.* 2005 for  
186 details). To determine how the models extrapolate we included values for variables in the evaluation  
187 strip that were outside the variable range in the study area for both *wetness* and *southness*. This is only  
188 one measure of ability to extrapolate, and other possibilities exist - for example, creating landscapes

189 or evaluation data with new combinations of variables. The simple test of extrapolation outside limits  
190 will suffice here.

191 Finally, to evaluate ability to predict habitat, mapped predictions were visually and  
192 quantitatively assessed. Summary statistics were calculated across predictions in all 80000 cells, in  
193 comparison to the presence-absence realisation and the true suitabilities. Analyses of predicted  
194 probabilities in relation to the true, sampled presence-absence data were: (i) the area under the  
195 receiver operating characteristic curve (AUC; Hanley and McNeil 1982); (ii) the remaining per-  
196 observation deviance (i.e. the variation left unexplained, as measured by the mean binomial deviance  
197 across sites; Elith and Leathwick 2007); (iii) the point biserial correlation coefficient (a Pearson  
198 correlation, Elith *et al.* 2006, COR.pa); (iv) elements of the confusion matrix for predictions  
199 converted to binary values, using a threshold that gave prevalence as close to possible to truth but that  
200 gave the same prevalence from each method. Predictions were compared with true suitabilities with a  
201 Pearson correlation coefficient (COR.si).

202 The important point about using more than one measure is that they quantify different aspects  
203 of predictive performance (see Murphy and Winkler 1987 and Pearce and Ferrier 2001 for interesting  
204 discussions of this topic). Amongst the measures tested against PA observations, AUC measures the  
205 ability of predictions to discriminate between observed presence and absence, regardless of the  
206 absolute value of the predictions. COR.pa also measures discrimination, but includes consideration of  
207 the actual value of the prediction, and how it compares to the observation. Deviance puts more  
208 emphasis on the model calibration i.e. to whether predictions reliably predict frequencies of  
209 occurrence. In other words, for deviance the actual value of the prediction is important. These  
210 different metrics measure different attributes of the predictions and are a useful and complementary  
211 set (Murphy and Winkler 1987).

### 212 **3. Results**

213 ***Relationship of species to environment:*** Relationships were tested by predicting to the evaluation  
214 strip. Figure 1 shows the modelled responses to the 3 variables when others are held at their optima,  
215 and in the right column, the response to co-varying *wetness* and *southness* for *geology* class one. Here

216 we analyse the results for the ranges of the variables in the data – i.e., within the blue vertical lines of  
217 Figure 1; in the later section "Extrapolation" we deal with predictions outside the range of the data.

218 We would not expect any of the methods to fully retrieve the species environment relationship  
219 of the simulated species because a relatively small sample of a binary realisation of the data was used  
220 to build the models (Figure 3a). Nonetheless, several of the methods were able to fit reasonably  
221 accurate functions. This was an easier task for the methods using smoothed functions (GLM and  
222 Maxent), given that the true relationships were smooth. All methods except GARP modelled *geology*  
223 correctly (Figure 1, left column). GARP's modelled response to *geology* varied across different data  
224 samples, different runs, and different summaries of the data; see online Supplement Appendix S2).  
225 GARP might be expected to model categorical data properly because it uses logistic rules, but the  
226 implementation is not a true logistic regression and does not properly deal with unordered categorical  
227 variables (Elith 2002 and Peterson pers.comm.). However, the atomic rules provide some opportunity  
228 to model *geology* correctly, and in some cases (online appendix Fig. S3 and S6) the result was better  
229 than the one presented.

230 The low suitability of dry areas (*wetness* < 15) was correctly captured by all methods (Figure 1,  
231 second and fourth columns). At higher levels of *wetness*, where the data are more sparse, BRT  
232 modelled the response most accurately, followed by Maxent (slightly reduced amplitude), RF (smaller  
233 amplitude) and GLM (unnecessary complexity around *wetness* of 50) and last, GARP (small  
234 amplitude and wrongly predicted that high *wetness* was unsuitable).

235 *Southness* (Figure 1, third and fourth columns) was difficult to model well, probably partly  
236 because the response included an interaction between *southness* and *wetness*, and also because it was  
237 less dominant than *wetness* in the suitability equation. RF, BRT and Maxent did best, GLM was good  
238 but gave an increasing response at low values of *southness* and, without an interaction, could not  
239 capture the response to *southness* at high *wetness* (see right side of 3D plot). GARP completely failed  
240 to model the true response and this result was consistent across all tests (Appendix S2).

241 **Mapped predictions; visual assessments:** Differences between the methods were also evident  
242 based on the mapped predictions. We present the results as maps and plots so that any arbitrariness

243 introduced by choice of legend in the map (Figures 2 and 4) can be checked against the plotted results  
244 (Figure 5). Note that the COR.s<sub>i</sub> in Table 2 acts as a summary measure of the plotted data in Figure 5.

245 Because all methods predicted the response to low *wetness* correctly, they all correctly  
246 predicted absence at dry sites (the white areas on the maps, Figure 2). Modelling this part of the  
247 response correctly meant that all methods produced a broadly correct mapped pattern. GARP tended  
248 to predict high values across any areas that had at least some suitability for the species, whereas the  
249 other methods predicted gradations in suitability more accurately (Fig. 2 and, a closeup in Fig. 4).  
250 Given that the true niche of the species is known and consisted of varying suitability depending on the  
251 combination of environmental parameters, failure to predict gradations is an error. The overprediction  
252 can be traced to the errors in retrieving the true underlying species-environment relationships.  
253 Simulated data like these could be used to explore whether another example or different settings for  
254 GARP improve model performance. For these data and other realisations of it, we spent some time  
255 attempting this but could not improve performance substantially beyond that shown (Appendix S2).  
256 The results of the different trials shown in the Appendix are interesting in that they suggest that means  
257 of all runs of GARP are slightly better for these data than the subsets that are generally advised. They  
258 also demonstrate considerable run-to-run variation (where one run is 500 models or summaries  
259 thereof).

260 As expected, none of the methods perfectly retrieved the true mapped suitabilities – sample  
261 size, use of binary data rather than suitabilities, and algorithmic limitations all contribute to this result.  
262 Maxent recreated the general mapped pattern of the simulated species well and only failed, and only  
263 moderately so, with respect to calibration (Figures 2 and 4). Perfect calibration would have resulted in  
264 all records in Figure 5 sitting on the diagonal, but a presence-only method cannot be well calibrated  
265 unless information on the species prevalence in the region is available. In contrast, the three methods  
266 trained on the PA data should be properly calibrated. The results are all reasonable, with BRT  
267 predictions slightly better constrained than RF and GLM (Figures 2, 4 and 5). We tested various  
268 settings for RF (seven combinations are presented in Appendix S4); the ones we present here are

269 those that would have been selected from the out-of-bag error estimates, and are consistent with those  
270 used in other published studies (Appendix S4).

271 ***Mapped predictions; quantitative analyses:*** The summary statistics presented in Table 2 focus on  
272 predictions as continuous values, and these demonstrate some differences between the methods, with  
273 the extent of the difference varying with the statistic. Note that the first row in Table 2 ("Truth")  
274 indicates the best possible performance for AUC, Deviance and COR.pa, because it measures the  
275 relationship between the true presence-absence realisation (sampled for modelling) and the true  
276 suitability. The AUC indicated that all models discriminated the broad patterns of presence and  
277 absence reasonably well, but with some variation between methods evident. The broad success in  
278 discriminating between presence and absence locations probably results from the correct modelling of  
279 dry conditions, because all the pairwise comparisons between predictions in these areas and ones in  
280 the wetter areas would have been correctly ranked. Measures that include consideration of the actual  
281 value of the predictions (all others) emphasise more clearly that Maxent, BRT and GLM all do well,  
282 followed by RF then GARP (Table 2). AUC and elements of the confusion matrix (Table 3) are the  
283 only truly fair basis here for comparing all methods, because the 3 models fitted to the PA data had  
284 more information than the PO methods, and this information allows the models to estimate  
285 probabilities and hence to be better calibrated. Nevertheless, for this exercise Maxent does well  
286 without this information. The fact that the order of model performance is consistent across metrics  
287 shows that errors in prediction are related to both discrimination and calibration. Given that the  
288 general shapes of the fitted functions for RF are reasonable, it is likely that its reduced performance in  
289 these summary metrics results from the slightly noisier fit compared with BRT (Figure 1, right  
290 column) and the slightly poorer calibration.

291 It is possible that using predictions as continuous values unfairly discriminates against GARP,  
292 because GARP works by producing a presence-absence prediction, and non-binary predictions result  
293 from summaries across multiple runs. Because of this, the predictions for all methods were  
294 thresholded and a test of binary predictions applied. We did this without biasing the results towards  
295 methods with good calibration, by using thresholds that gave the same prevalence across the

296 landscape for all methods. This meant that we used a threshold of 1 for GARP (to get prevalence as  
297 close to truth as possible and to restrict the overprediction of GARP), then set others in relation to  
298 that. The results (Table 3) are consistent with other results, though the differences are less extreme. If  
299 we had set separate thresholds to give the best binary realisation for each method (using the known  
300 prevalence), the result for GARP would remain constant and the others would all improve.

301 Information on BRT and GLM models fitted to the presence - pseudo-absence data are  
302 presented in the online supplement, Appendix S5. These show that in some (but not all) cases the  
303 fitted functions are reasonably accurate, and that there are differential effects on the discrimination  
304 and calibration of the models. We note that alternative implementations of BRT that are better suited  
305 to the presence / pseudo-absence data structure will be available soon (Ward *et al.* in press).

306 **Extrapolation:** The extrapolation behaviours of the models is presented in Figure 1, in which the  
307 responses to the highest and lowest values for the *wetness* and *southness*, outside the blue dashed  
308 lines, are where the models are extrapolating to unsampled conditions outside the range of these  
309 variables in the mapped region. The nonsensical negative values of *wetness* and *southness* don't matter  
310 here – the important question is how the models would extrapolate beyond the sampled values of data.  
311 We had no prior knowledge of GARP behaviour, and in this example extrapolation was either a  
312 stepped decline (at extremes of *southness*) or a constant value. The selected rules in the best-subsets  
313 models would explain this behaviour, but they are not accessible in the desktop program so could not  
314 be analysed. Extrapolation patterns varied in GARP according to the dataset, the particular run, and  
315 the selected subsets (online Appendix Figure S3). By contrast, Maxent acts consistently and is  
316 "clamped" so it extrapolates in a horizontal line from the fit at the most extreme environmental value  
317 in the training data, both presence and background. As expected from the way polynomials behave, a  
318 GLM fitted with cubic and quadratic functions extrapolates by continuing the fitted trend beyond the  
319 last observation, sometimes with unwanted results (Austin *et al.* 1990). For example, here the  
320 projected increase in suitability at the lower values of *southness* is not sensible (north-facing (low  
321 *southness*) sites should be least suitable for the species). Classification and regression trees always  
322 extrapolate at a constant value from the last "known" site, as seen for BRT and RF.

#### 323 **4. Reflections on the simulation**

324 Before considering the wider implications of this simulation, we want to first emphasise what it  
325 does not do. Whilst our simulated species had some ecological realism (it is affected by more than one  
326 variable, and reacts to predictors in non-linear and non-additive way), we were not trying to emulate  
327 nature, but to interrogate the behaviour of different algorithms. We also only tested some of the many  
328 ways of fitting these methods - for example, a GLM can include interactions and splines can be used  
329 to control behaviours at the extremes of the sample ranges. We also postulated our application as  
330 mapping suitabilities, whereas mapping binary presence-absence predictions are an alternative output.  
331 Finally, we chose scenarios that give some insight but did not attempt a comprehensive study.

332 What then does this simulation demonstrate? First, knowing what an algorithm is doing can  
333 give insights into various features that are apparent in its predictions; it helps to answer why particular  
334 behaviour is observed. Testing an algorithm's performance in this way reveals problems and may help  
335 developers and modellers to understand and improve model performance. The systematic studies in  
336 Appendices S2, S3 and S5 demonstrate that looking carefully at model performance across settings  
337 helps to expose what is happening.

338 Second, our example clearly illustrates the value of evaluating models from several viewpoints.  
339 The summary statistics indicated which algorithm gave the best mapped predictions, and gave some  
340 hints about why through the different metrics. AUC only measures rank so did not reveal the more  
341 extreme differences between the models that were more related to calibration of the model. We do not  
342 believe that this means that any of these statistics are misleading (see the title of Lobo *et al.* 2008), but  
343 simply that different statistics measure different aspects of performance, and that appropriate statistics  
344 relevant to the application of the model need to be selected. Being able to visualise fitted functions  
345 not only satisfied our application of exploring modelled relationships, but also allowed us to  
346 understand what caused differences among the methods and how different fitted functions influenced  
347 mapped predictions.

348 Third, comparing the results for the partial responses with the quantitative assessments gives  
349 some useful insights. For example, even though the GLM modelled unnecessary complexity in the

350 *wetness* response (giving the wave-like forms in the 3D plots), the evaluation statistics implied it did  
351 nearly as well as BRT. This is in spite of the fact that BRT had a better controlled fit overall. The  
352 distribution of environments (Figure 3b) gives the key: there are relatively few sites in the places  
353 where the GLM has failed to model the true response. Similarly, the environmental distribution of the  
354 training data (Figure 3a) explains why several methods (GLM, BRT and RF) tended to model a  
355 declining response (even if only a small trough) to *wetness* at values early on the plateau of the true  
356 response - the samples are sparse in some parts of this environment, and several have been realised as  
357 "absence" in this particular sample. These demonstrate that understanding the environmental  
358 distribution of the data in the region of interest, in the sample, and in regions that might be used for  
359 projection, is a critical part of understanding the implications for modelling and prediction.

360 Finally, the demonstration prompts a range of questions about what characteristics we want  
361 from models in certain situations. For example, when using species distribution models to predict how  
362 species ranges will shift with climate change or how they will extrapolate to new regions it is critical  
363 that we understand how the algorithm performs when projected into new environmental combinations  
364 combinations not sampled by the training data. In other words, is the way that the algorithm  
365 extrapolates appropriate from an ecological perspective? Different behaviours are apparent in the  
366 methods that we tested, but the choice as to which (if any) is correct is as much an ecological and/or a  
367 physiological question, as a statistical one. For a full investigation of extrapolation or forecasting  
368 behaviour a much larger range of tests is required including prediction to new combinations of  
369 environments, and the species should to be simulated with known responses to these. The important  
370 point is that we need to first recognise what different modelling applications require of SDMs and  
371 then research the best means for achieving what they require. Understanding how the models work  
372 and devising evaluation criteria that are closely matched to the questions being asked can inform  
373 decisions about the best modelling approach.

#### 374 **Conclusion**

375 Given the multitude of applications of SDMs, many of them related to pressing conservation  
376 issues, we hope this forum paper stimulates research on how to effectively continue to develop and

377 test these methods. Additional model comparison studies may not be fruitful unless they start to ask  
378 why certain methods perform better than others. This is particularly important for applications that do  
379 not satisfy the underlying assumptions of species equilibrium with environment (Dormann 2007),  
380 such as range shifts with climate, or species invasions into a new area. By providing an example that  
381 we believe gives useful insight into model performance, we hope that both developers and users will  
382 increasingly question what their models are doing and whether that is most appropriate for the  
383 intended applications and outcomes.

#### 384 **Acknowledgements**

385 Many thanks to Mark Burgman, John Leathwick, Yung En Chee, Stephen Baines, Michael  
386 Kearney, Steven Phillips and Town Peterson for comments on various drafts of the manuscript, and to  
387 Mike Austin, Simon Ferrier and Simon Barry for asking challenging questions that were the starting  
388 points for several of the ideas. The reviewers and editors comments gave important direction and we  
389 appreciate their efforts. Jane Elith was funded by ARC grant DP0772671 and the Australian Centre  
390 of Excellence for Risk Analysis. The colours used in the figures are from a palette developed for  
391 colour-blind people; thanks to its author:[http://jfly.iam.u-tokyo.ac.jp/html/color\\_blind/#stain](http://jfly.iam.u-tokyo.ac.jp/html/color_blind/#stain)

392

#### 393 **References**

- 394 Araújo, M. B. and Guisan, A. 2006. Five (or so) challenges for species distribution modelling. - J.  
395 Biogeogr. 33: 1677-1688.
- 396 Araújo, M. B. and Rahbek, C. 2006. How does climate change affect biodiversity? Science 313:1396  
397 - 1397.
- 398 Araújo, M. B. *et al.* 2005. Validation of species-climate impact models under climate change. - Global  
399 Change Biology 11: 1504-1513.
- 400 Austin, M. 2007. Species distribution models and ecological theory: A critical assessment and some  
401 possible new approaches. - Ecol. Model. 200: 1-19.
- 402 Austin, M. P. *et al.* 2006. Evaluation of statistical models used for predicting plant species  
403 distributions: Role of artificial data and theory. - Ecol. Model. 199: 197-216.
- 404 Breiman, L. 2001a. Statistical modeling: the two cultures. - Statistical Science 16: 199-215.

- 405 Breiman, L. 2001b. Random Forests Technical Report.  
406 <http://oz.berkeley.edu/users/breiman/randomforest2001.pdf>
- 407 Busby, J. R. 1991. BIOCLIM - a bioclimate analysis and prediction system. - In: Margules, C. R. and  
408 Austin, M. P. (eds.), Nature Conservation: Cost Effective Biological Surveys and Data Analysis.  
409 CSIRO, pp. 64-68.
- 410 Carpenter, G., Gillison, A. N. and Winter, J. 1993. DOMAIN: a flexible modelling procedure for  
411 mapping potential distributions of plants and animals. - *Biodivers. Conserv.* 2: 667-680.
- 412 Caruana, R. and Niculescu-Mizil, A. 2006. An Empirical Comparison of Supervised Learning  
413 Algorithms. - In, Proceedings of the 23 rd International Conference on Machine Learning,  
414 Pittsburgh, PA.
- 415 doi:10.1111/j.1365-2699.2007.01779.x. - *J. Biogeogr.* 35: 105-116.
- 416 Dormann, C. F. 2007. Promising the future? Global change projections of species distributions. -  
417 *Basic and Applied Ecology* 8: 387-397.
- 418 Elith, J. 2002. Predicting the Distribution of Plants. PhD thesis. School of Botany. The University of  
419 Melbourne.
- 420 Elith, J. *et al.* 2005. The evaluation strip: a new and robust method for plotting predicted responses  
421 from species distribution models. - *Ecol. Model.* 186: 280-289.
- 422 Elith, J. *et al.* 2006. Novel methods improve prediction of species' distributions from occurrence data.  
423 - *Ecography* 29: 129-151.
- 424 Friedman, J. H., Hastie, T. and Tibshirani, R. 2000. Additive logistic regression: a statistical view of  
425 boosting. - *Ann. Statist.* 28: 337-407.
- 426 Graham, C. H. *et al.* 2004. New developments in museum-based informatics and applications in  
427 biodiversity analysis. - *Trends in Ecology & Evolution* 19: 497-503.
- 428 Graham, C. H. *et al.* 2008. The influence of spatial errors in species occurrence data used in  
429 distribution models. - *J. Appl. Ecol.* 45: 239-247.
- 430 Guisan, A. and Thuiller, W. 2005. Predicting species distribution: offering more than simple habitat  
431 models. - *Ecology Letters* 8: 993-1009.
- 432 Hanley, J. A. and McNeil, B. J. 1982. The meaning and use of the area under a Receiver Operating  
433 Characteristic (ROC) curve. - *Radiology* 143: 29-36.
- 434 Hastie, T. and Tibshirani, R. 1990. *Generalized Additive Models.* - Chapman and Hall.

- 435 Hernandez, P. A. *et al.* 2006. The effect of sample size and species characteristics on performance of  
436 different species distribution modeling methods. - *Ecography* 29: 773-785.
- 437 Kearney, M. *et al.* in press. Modelling species distributions without using species distributions: the  
438 cane toad in Australia under current and future climates. - *Ecography*.
- 439 Kozak, K., Graham, C.H. and Wiens, J.J. In press. Species distribution modeling in evolutionary  
440 biology. *TREE*.
- 441 Latimer, A. M. *et al.* 2006. Building statistical models to analyze species distributions. - *Ecol. Appl.*  
442 16: 33-50.
- 443 Latimer, A. M. *et al.* 2006. Building statistical models to analyze species distributions. - *Ecol. Appl.*  
444 16: 33-50.
- 445 Lobo, J. M., Jiménez-Valverde, A. and Real, R. 2008. AUC: a misleading measure of the performance  
446 of predictive distribution models. - *Global Ecology and Biogeography* 17: 145-151.
- 447 Loiselle, B. A. *et al.* 2008. Predicting species distributions from herbarium collections: does climate  
448 bias in collection sampling influence model outcomes?
- 449 McCullagh, P. and Nelder, J. A. 1989. *Generalized Linear Models*. - Chapman and Hall.
- 450 McPherson, J. M. and Jetz, W. 2007a. Type and spatial structure of distribution data and the perceived  
451 determinants of geographical gradients in ecology: the species richness of African birds. - *Global  
452 Ecology and Biogeography* 16: 657-667.
- 453 McPherson, J. M. and Jetz, W. 2007b. Effects of species' ecology on the accuracy of distribution  
454 models. - *Ecography* 30: 135-151.
- 455 Meynard, C. N. and Quinn, J. F. 2007. Predicting species distributions: a critical comparison of the  
456 most common statistical models using artificial species. - *J. Biogeogr.* 34: 1455-1469.
- 457 Moisen, G. G. and Frescino, T. S. 2002. Comparing five modeling techniques for predicting forest  
458 characteristics. - *Ecol. Model.* 157: 209-225.
- 459 Murphy, A. H. and Winkler, R. L. 1987. A general framework for forecast verification. - *Monthly  
460 Weather Review* 115: 1330-1338.
- 461 Orr, D. W. 2000. Ideasclerosis: Part One. - *Conserv. Biol.* 14: 926-928.
- 462 Pearce, J. and Ferrier, S. 2000. An evaluation of alternative algorithms for fitting species distribution  
463 models using logistic regression. - *Ecol. Model.* 128: 127-147.

- 464 Peterson, A. T., Papes, M. and Eaton, M. 2007. Transferability and model evaluation in ecological  
465 niche modeling: a comparison of GARP and Maxent. - *Ecography* 30: 550-560
- 466 Phillips, S. 2008. Response to "Transferability and model evaluation in ecological niche modelling". -  
467 *Ecography* 31: 272-278.
- 468 Phillips, S. in press. Response to "Transferability and model evaluation in ecological niche  
469 modelling". - *Ecography*.
- 470 Phillips, S. J. and Dudik, M. 2008. Modeling of species distributions with Maxent: new extensions  
471 and a comprehensive evaluation. - *Ecography* 31: 161-175.
- 472 Phillips, S. J., Anderson, R. P. and Schapire, R. E. 2006. Maximum entropy modeling of species  
473 geographic distributions. - *Ecol. Model.* 190: 231-259.
- 474 R Development Core Team 2006. R: A Language and Environment for Statistical Computing. - In, R  
475 Foundation for Statistical Computing.
- 476 Redford, K.H. and Taber, A. 2000. Writing the wrongs: developing a safe-fail culture in conservation.  
477 - *Conserv. Biol.* 14: 1567-1568.
- 478 Reese, G. C. *et al.* 2005. Factors affecting species distribution predictions: A simulation modeling  
479 experiment. - *Ecological Applications* 15: 554-564.
- 480 Reineking, B. and Schröder, B. 2006. Constrain to perform: regularization of habitat models. - *Ecol.*  
481 *Model.* 193: 675-690.
- 482 Segurado, P. and Araujo, M. B. 2004. An evaluation of methods for modelling species distributions. -  
483 *J. Biogeogr.* 31: 1555-1568.
- 484 Tyre, A. J., Possingham, H. P. and Lindenmayer, D. B. 2001. Matching observed pattern with  
485 ecological process: can territory occupancy provide information about life history parameters? -  
486 *Ecol. Appl.* 11: 1722-1738.
- 487 Ward, G. *et al.* in press. Presence-only data and the EM algorithm. - *Biometrics*.
- 488 Wu, J. and Hobbs, R. 2002. Key issues and research priorities in landscape ecology: An idiosyncratic  
489 synthesis. - *Landscape Ecol.* 17: 355-365.
- 490 Yee, T. W. and Mitchell, N. D. 1991. Generalized additive models in plant ecology. - *J. Veg. Sci.* 2:  
491 587-602.
- 492
- 493

494 Table 1 – Details of model fitting procedures and settings

Method	Name for model	Data	Settings and notes on further tests
Genetic Algorithm for Ruleset Prediction	GARP	PO	Used v 1.1.6. Details in online supplement Appendix S2 on tests to compare effects of different settings. Results presented here from: species data = 115 PO samples plus pseudo-absences selected by GARP. 50% data used for training, 50% for extrinsic evaluation. Created 500 models each with a convergence limit of 0.01 and 1000 maximum iterations. Allowed all rule types. From the 500models chose 20% with mid extrinsic omission error and from those 20 with mid commission error. Final prediction is mean of these. For predicting to evaluation strip projected to grids with strip inserted.
Maximum entropy	Maxent	PO	Used version 3.2.1 from the command line. Modelled the 115 PO samples and allowed Maxent to select a random 10000 background samples (the default). All other settings were the defaults except: flagging <i>geology</i> as a categorical variable, providing a separate set of grids to project to that contained the evaluation strip, and using the "-d" flag (see help file for Maxent). The -d flag forces Maxent to calculate the probability distribution over the background samples alone (rather than the default, which calculates it over the joint background and presence data), and providing it with the best chance to be well calibrated. For predicting to evaluation strip projected to grids with strip inserted.
Generalised linear models	GLM.pa	PA	Used R <sup>1</sup> and function <i>glm</i> . Created all possible subsets of models with the options for each variable being: exclude, or (if continuous): linear, quadratic or cubic fits. Used AIC to select the best model.
Boosted regression trees	BRT.pa	PA	Used R <sup>1</sup> and function <i>gbm</i> with custom scripts of Elith <i>et al</i> , 2008 to build an ensemble of regression trees. Selected tree complexity of 3, learning rate of 0.001, using prevalence-stratified cross-validation to determine optimal number of trees (4250). See Appendix S4 for details.
Random forests	RF.pa	PA	Used R <sup>1</sup> and function <i>randomForest</i> to build an ensemble of classification trees. Tests of a range of settings are presented in Appendix S3. Model presented here had 500trees with one variable randomly selected from the 3 candidates at each split. No class weights.

495

496

497

498 Table 2: Comparison of model results with truth, as realised by the presence-absence map  
 499 (columns 2, 3 and 4) and the suitability values (column 5). For all statistics except deviance,  
 500 higher is better.

<b>Model</b>	<b>AUC</b>	<b>Remaining deviance</b>	<b>COR.pa</b>	<b>COR.si</b>
Truth (suitabilities)	0.872	0.514	0.508	1.000
GARP	0.822	3.391	0.401	0.793
Maxent	0.861	0.612	0.467	0.922
GLM.pa	0.863	0.546	0.480	0.941
BRT.pa	0.862	0.537	0.485	0.954
RF.pa	0.834	0.736	0.448	0.875

501

502

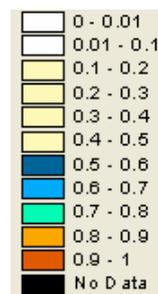
503 Table 2: Method comparisons when predictions are reduced to binary results. Because GARP  
 504 overpredicted, the highest possible GARP threshold (1) was used to convert GARP predictions to  
 505 binary form, then thresholds were selected for all other methods that gave identical prevalence in  
 506 the landscape (17.3%, compared with truth 11.9%). The values in the tables are the proportions  
 507 of predictions that fell into each category (true negative etc) when cross-tabulated with truth.

Prediction:	true negative	true positive	false negative	false positive
GARP	0.771	0.064	0.056	0.109
Maxent	0.778	0.071	0.049	0.102
GLM	0.780	0.073	0.047	0.100
BRT	0.780	0.073	0.047	0.100
RF	0.777	0.070	0.050	0.103

508 **Figure legends:**

509 **Figure 1:** Partial responses to the 3 variables (left) and over co-varying wetness and southness (3D  
 510 plots, right). The true responses (top panel) were generated by using the equations that define the  
 511 simulated species to predict to the evaluation strip, then plotting the results (see Elith *et al.* 2005 for  
 512 details). The blue vertical lines show the extent of the variable values in the mapped region; outside  
 513 these the models are extrapolating. The range on the *wetness* and *southness* axes of the 3D plots is that  
 514 within the blue lines of the 2D ones, and predictions range from 0 to 1.

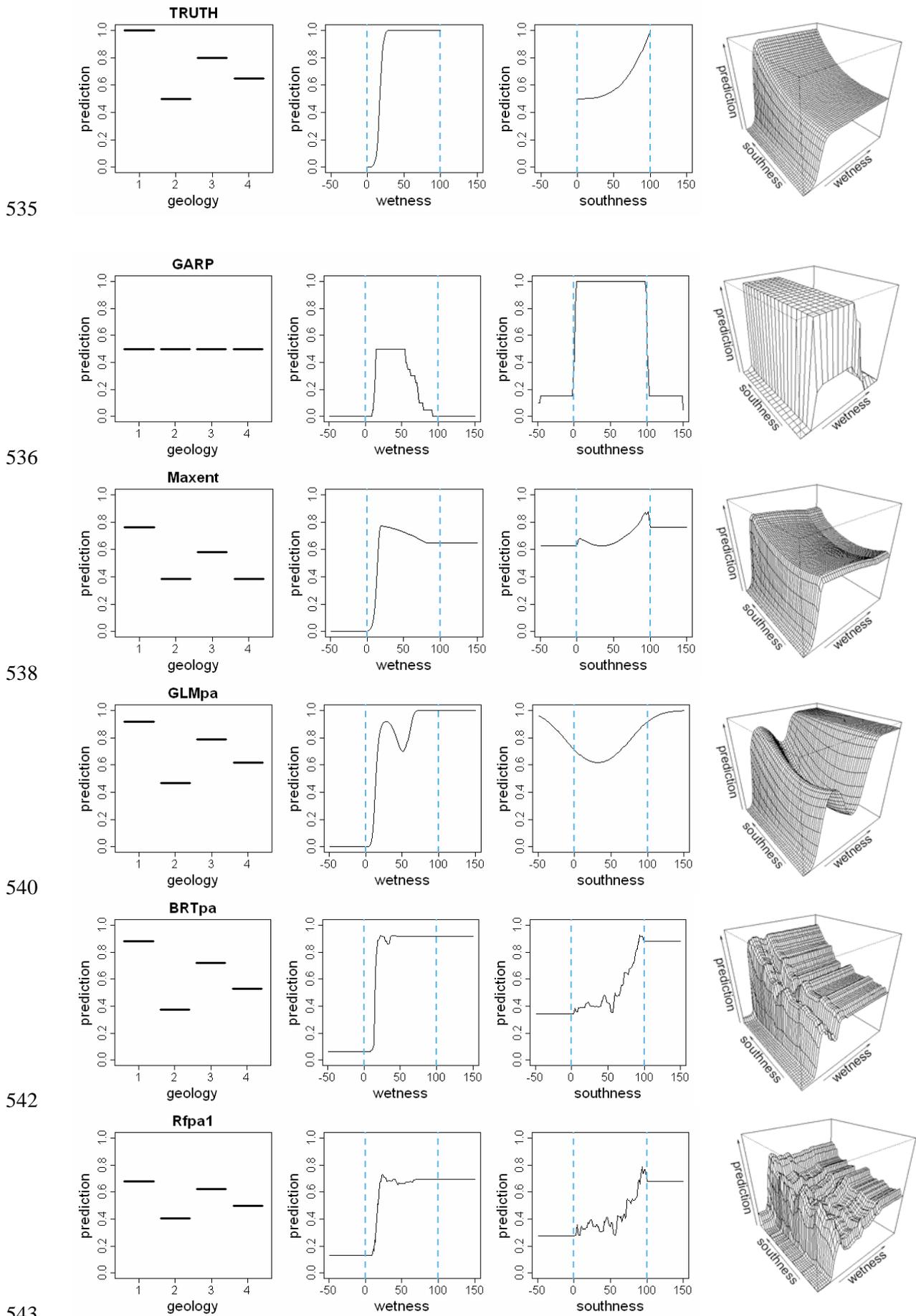
515 **Figure 2:** Mapped distributions of the virtual species (top left) and predictions of relative suitabilities  
 516 from the methods detailed in the text, Legend: white < 0.1, cream 0.1 to 0.5, blue-lightblue-green-  
 517 orange-vermillon at steps of 0.1 from blue (0.5 to 0.6) to vermillon (0.9 to 1) (note to editor - this  
 518 legend could be included):



519  
 520  
 521  
 522 **Figure 3: (a):** The data samples for presence-absence models. Samples are shown at their original  
 523 true suitability value (vertical axis), but were converted to presence (blue) or absence (orange) as  
 524 described in the text. **(b):** The location of all 80000 grid cells in environmental space. The pale  
 525 yellow mesh shows the full suitability surface from the simulated species, for geology class 1. The  
 526 points of varying colours show sites in the four geology classes. Note the few sites with high wetness  
 527 values.

528 **Figure 4:** Close-up of predictions from Figure 2. Choice of location was via random number selection  
 529 for centre grid position. Predictions in greyscale, from white (zero) to black (one); fine grid lines are  
 530 in the same position on each map.

531 **Figure 5:** Predictions (y axis) versus the true suitability for all 80000 grid cells in the maps in Figure  
 532 2, covering five modelling methods described in Table 1. The blue diagonal line shows the 1:1  
 533 relationship.



535

536

538

540

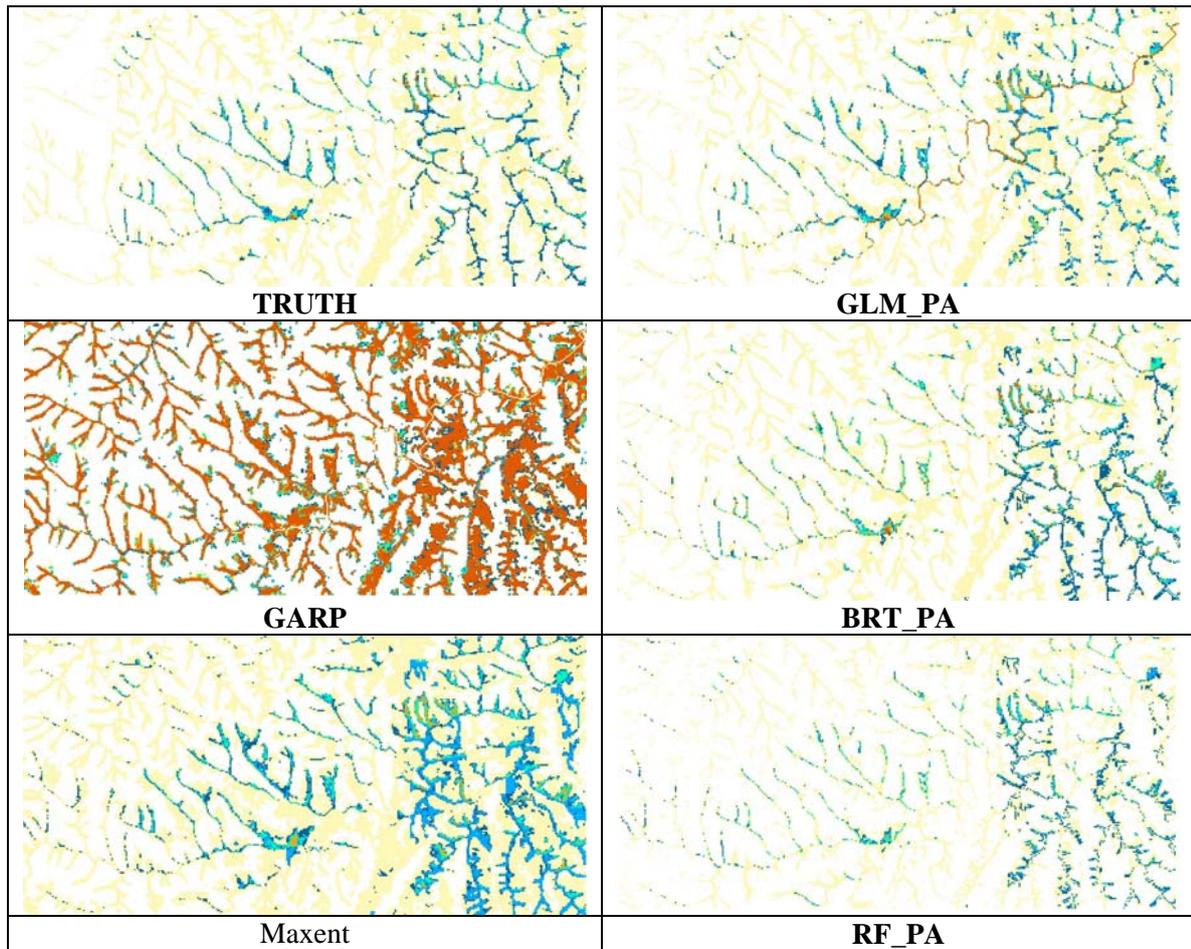
542

543

544

**Figure 1**

545



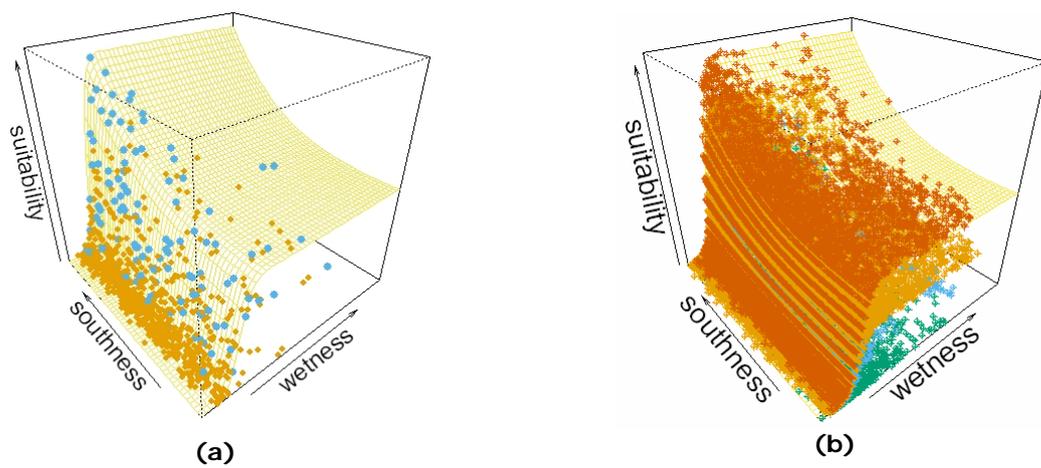
546

547

548

549

Figure 3



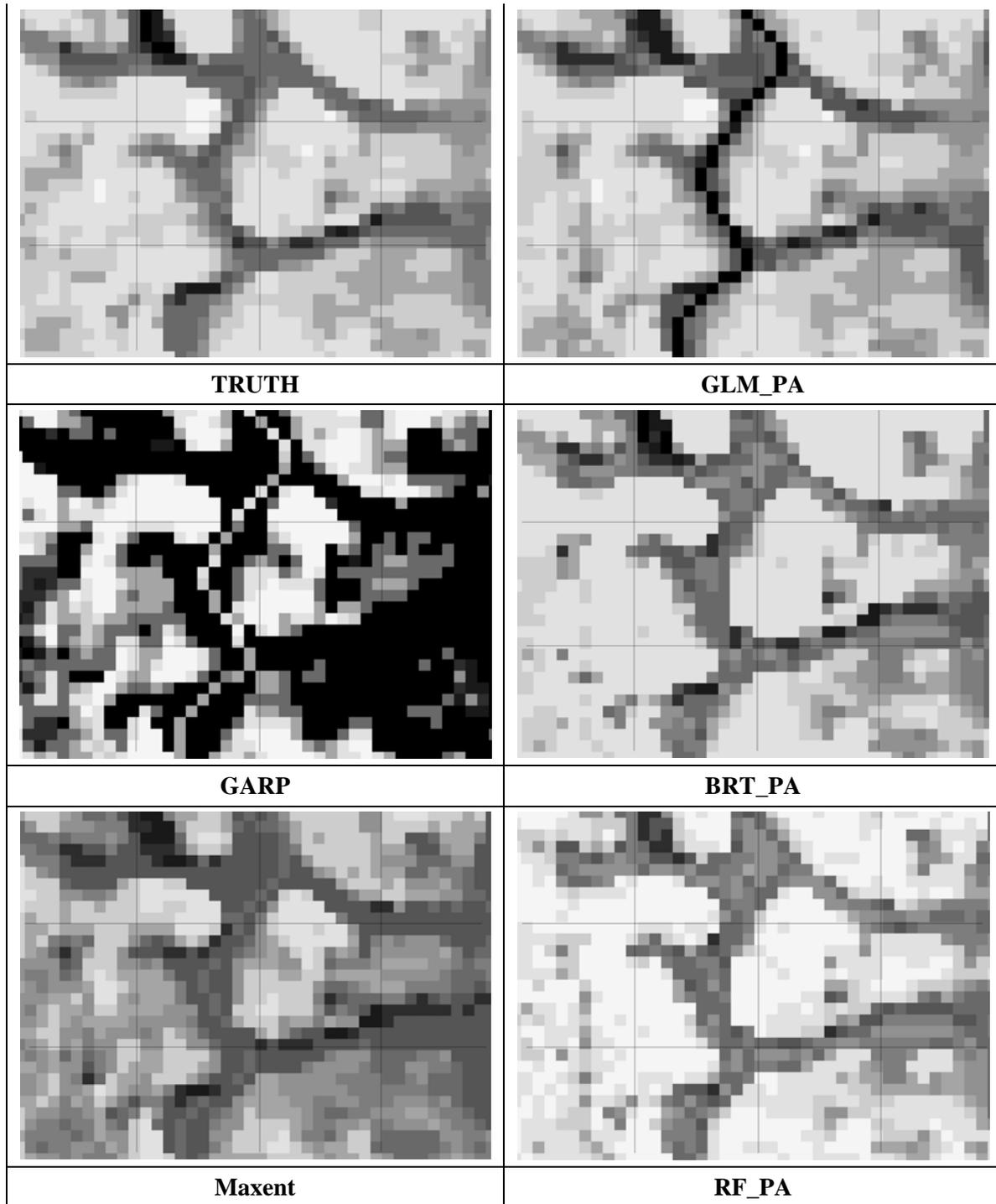
550

551

552

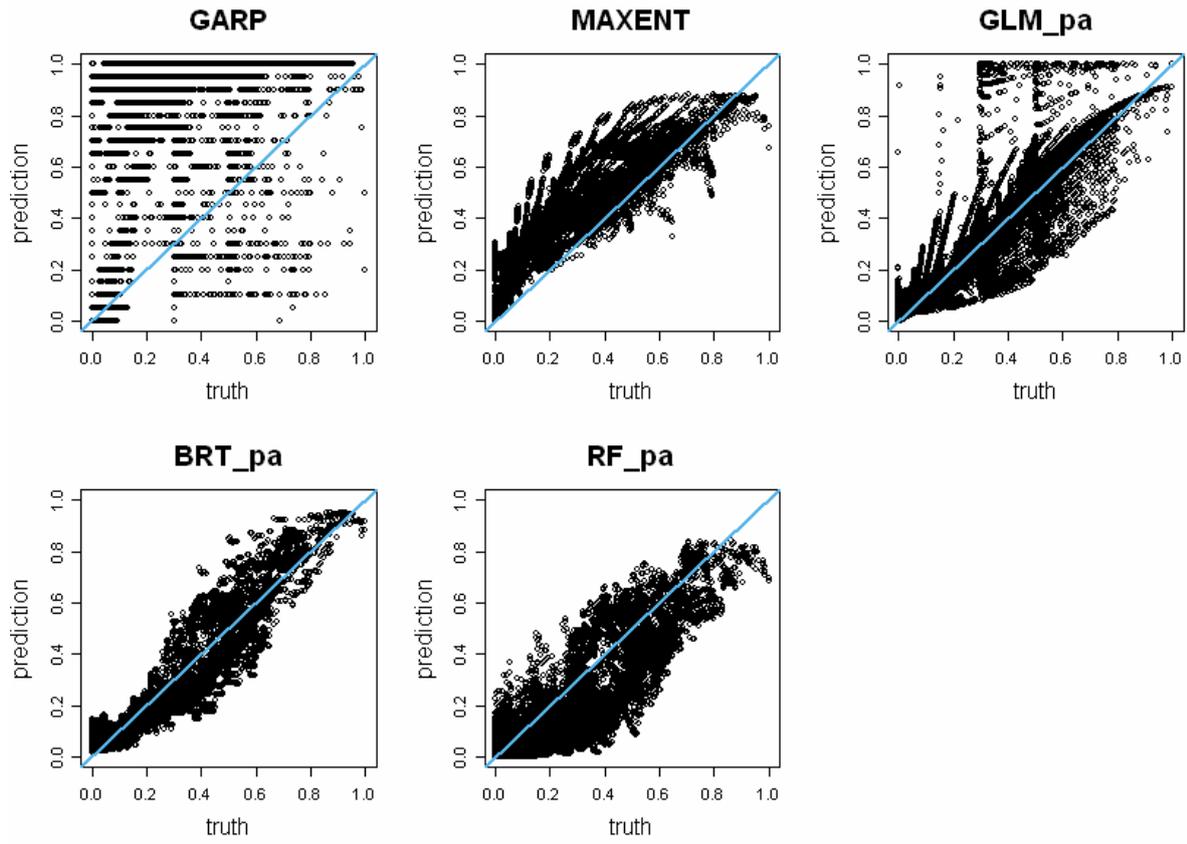
Figure 3

553



554  
555  
556

Figure 4



557  
558

559

560

Figure 5

**Online supplementary material: Appendix S1: Details of the data generation**

The species, described in Elith and Graham (2008: equation 1 and Figure 1), responds to three environmental gradients:

$$\text{Suitability (SI)} = \text{SI.wetness} * 0.5 * (\text{SI.southness} + \text{SI.geology}) \quad \text{- equation S1}$$

where SI = suitability index

SI.wetness is the individual response to wetness, varying non-parametrically between 0 and 1 (Figure S1a), and SI.southness is the response to how south-facing a site is, varying parametrically between 0 and 1 (Figure S1b):  $\text{SI.southness} = 0.000001 * (\text{southness}^3)$ . The response to geology (SI.geology) is simply set at four levels: response to class 1 = 1, class2 = 0, class3 = 0.6, class4 = 0.3.

The overall suitability is not a simple addition of these terms but involves an interaction between wetness and the sum of the responses to southness and geology. This implies that southness and geology substitute for one another but wetness overrides both.

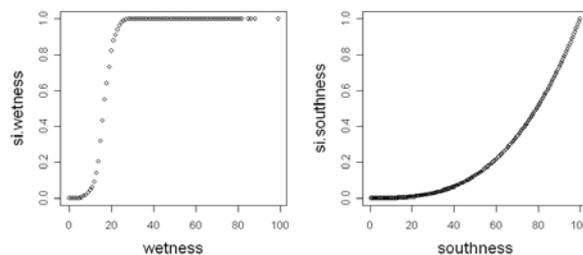


Figure S1. Individual suitability indices for (a) wetness and (b) southness

Using mapped grids of wetness, southness and geology (400 columns by 200 rows) from a real region in south-east Australia, we created the suitability indices for each, and a composite SI for the simulated species, from equation S1. The 80000 SI values were mapped, and also converted to binary values (using the function *rbinom* in R; R Development Core Team 2006; Figure S2a), which were then sampled (Elith and Graham 2008). In addition to the presence-absence (PA) and presence-only (PO) samples described in Elith and Graham (2008), we also created pseudo-absence samples. For these we randomly sampled from the 80000 grid cells in the region, with sample sizes of 1000 and 3000 sites. These will be called P0.1000 and P0.3000. In each of these pseudo-absence samples, some sites will, in reality, be inhabited by the species (i.e. as expected for pseudo-absences like this, we will create contaminated absences). In our samples there were 116 true presence sites used as pseudo-absences in P0.1000 and 355 in P0.3000. However in the model we treat all as absences to be consistent with the case where true absence is not known (Figure S3b).

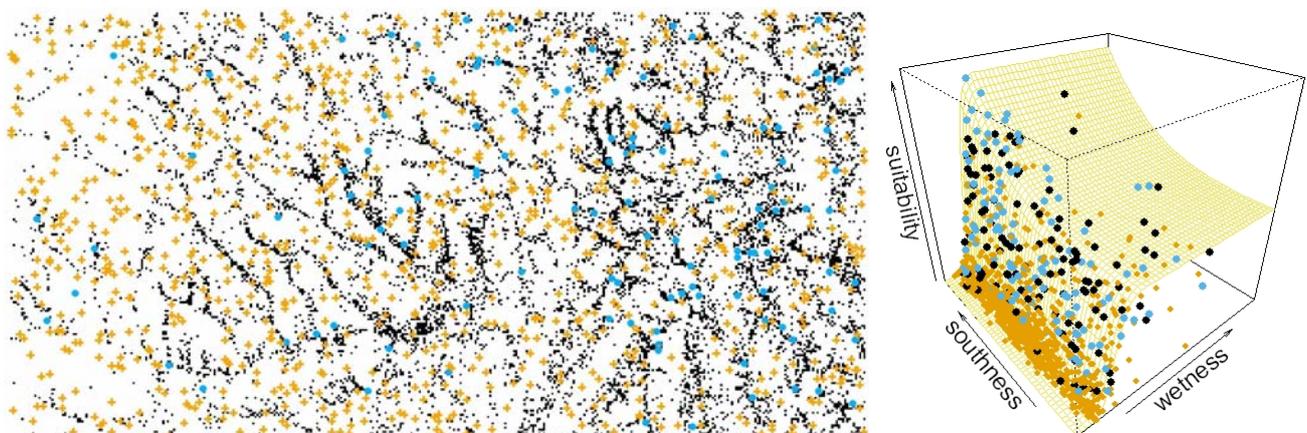


Figure S2: (a) Left: Map showing realised presence-absence data (black presence white absence, with the PA sample: blue (presence) and tan (absence)); (b) Right: The data sample for presence plus pseudo-absence 1000. Samples are shown at their original true suitability value (vertical axis), but were converted to presence (blue and black) or absence (orange) as described in the text. The black points are those in the pseudo-absence sample that were subsequently changed to "absence" for modelling.

**Online supplementary material:****Appendix S2: Running and testing GARP**

We used the current version of Desktop GARP (Genetic Algorithm for Rule-set Production; Stockwell and Noble 1992; version 1.1.6 from late 2007) and followed Peterson *et al.* (2007) for parameter settings. This meant that the sampled presence data (115 records) were supplied for model fitting, and we allowed GARP to select the pseudo-absences. We used 50% data for training, and 50% for extrinsic evaluation, and created 500 models each with a convergence limit of 0.01 and 1000 maximum iterations (500 models took about 8 hours on desktop PC). We offered all three environmental predictors and allowed all rule types, then for predictions evaluated the mapped predictions for performance measures and also projected to grids with the evaluation strip included. We chose two subsets from the 500 models: first, following Peterson *et al.* (2007), selected the 20% of models<sup>1</sup> with the lowest extrinsic omission error, and then selected from that subset of approximately 100 models the twenty models with commission errors in the middle of the range of commission indices. Second, selecting models based on low omission is considered best for predicting potential distributions might not be optimal for this application where we are attempting to accurately model records of the true niche of the species. Therefore, we took a second subset like the first, but the first 20% were selected as those in the middle of the range of extrinsic omission errors. We processed both these subsets of 20 (using the mean prediction as the prediction per grid cell), and found that the results were similar (Figure S3, rows 2 and 3). In the paper we present the second variant – i.e. mid external omission and mid commission error - because conceptually it seemed more consistent with the task required. This meant that in the paper geology was not modelled well, but the test results (Table S1) were better - compare rows 4 and 5, Table S1. We also tested various other combinations of the 500 models, in an attempt to reduce commission error (Figure S3). None of our attempts improved GARP performance so that it was comparable to the other methods, though we note that the mean of 500 (fourth row, Figure S3) provided the best results. We did not include it in Elith and Graham (2008) because it is not the method recommended by those most experienced at running GARP.

The fitted responses across a range of pairwise values of wetness and southness (Figure S3, right column) revealed that, except for the minimum set, the responses wetness and southness are consistent across a range of values of the other and so are not dependent on the precise value at which wetness or southness were kept constant. The relevant predictive maps are shown in Figure S5.

To explore the consistency of these results we repeated the analysis (about 500 runs, same settings) on 2 repeats of the same data and one new sample of the simulated species, using 125 new presence records from the sample of 1000 (see earlier). Results were summarised for mid-omission and mid-commission errors as above (Figure S6). There is some variation amongst runs but again no results do as well as the other methods tested in the paper.

Table S1: Comparison of model results with truth, as realised by the presence-absence map (columns 2, 3 and 4) and the suitability values (column 5). For all statistics except deviance, higher is better. Models 2 to 5 are those described in paragraph 1, above, and 6 to 8, in paragraph 2. Models 5 to 7 (highlighted) are run with the same presence records with the same settings, so any differences are due to stochasticity in the model.

Model	AUC	Remaining deviance	COR.pa	COR.si
1.Truth (suitabilities)	0.872	0.514	0.508	1.000
2.GARPmin	0.564	3.233	0.224	0.441
3.GARPmean	0.842	1.856	0.407	0.805
4.GARPlow omission	0.812	4.812	0.385	0.752
5.GARPmid omission (paper)	0.822	3.391	0.401	0.793
6.GARPrepeat1	0.807	4.470	0.388	0.767
7.GARPrepeat2	0.814	3.944	0.391	0.772
8.GARPnew125pres	0.819	4.034	0.386	0.757

<sup>1</sup> If at the bounds of the subset there were multiple sites with the same extrinsic omission error, we expanded the sets to include all with that error rate. However when selecting the second set based on commission error, strictly selected 20 sites.

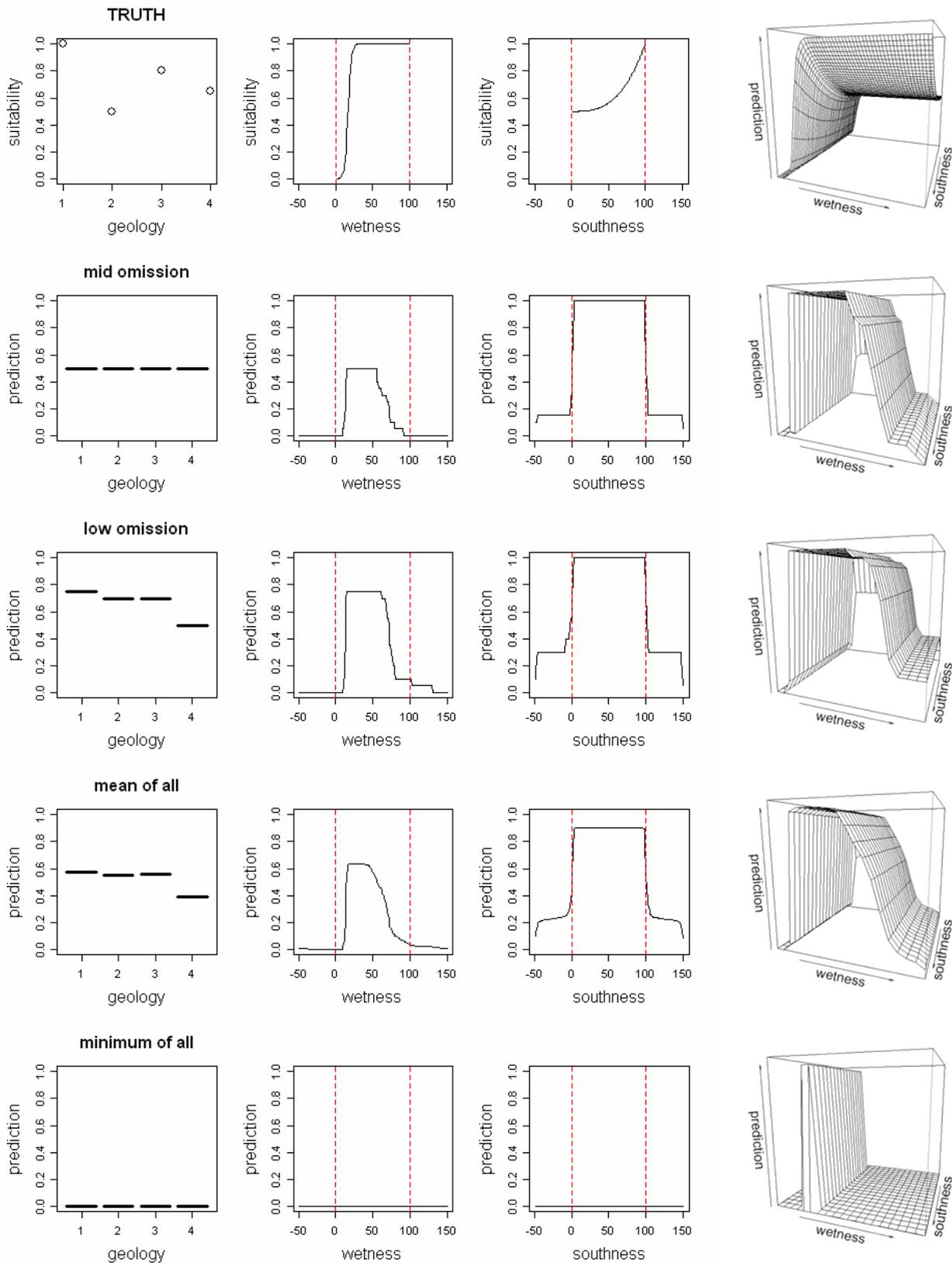


Figure S3: Fitted functions for truth plus four different summaries of the GARP run of 500 models. The second top one ("mid omission") is the one presented in Elith and Graham (2008). Related pairwise plots and predicted maps are in Figures S4 and S5 (not that these are rotated to a different perspective compared with those in Elith and Graham 2008).

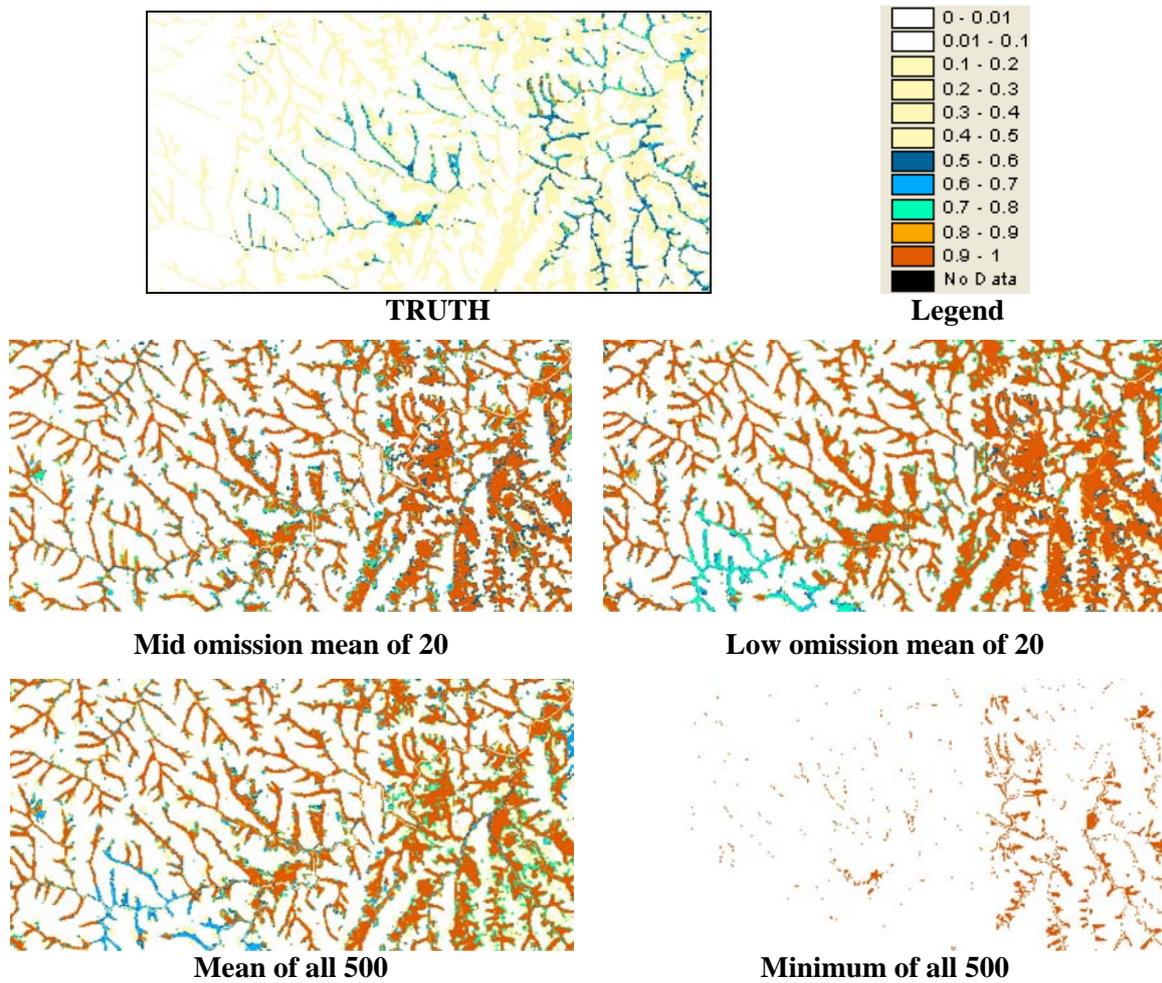


Figure S4 – Maps of predicted distributions, from the models illustrated in Fig. S3 and summarised in Table S1. Legend same as for main paper, and shown in top row

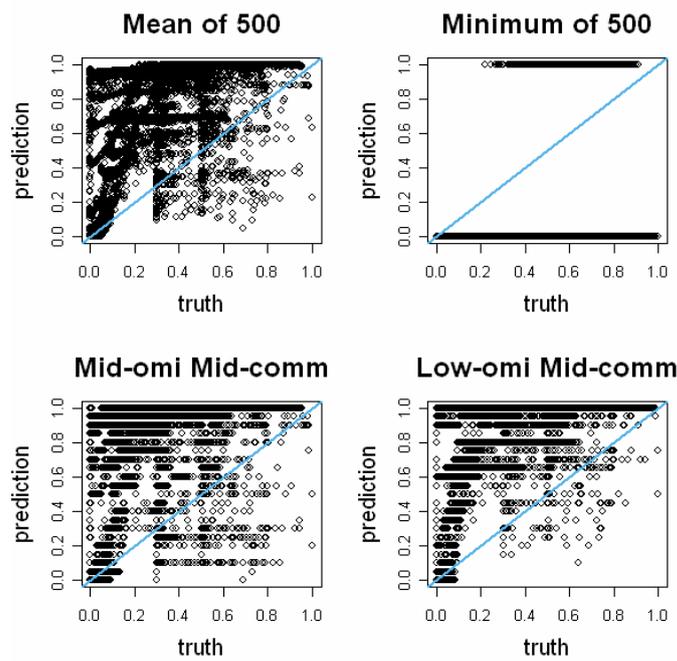
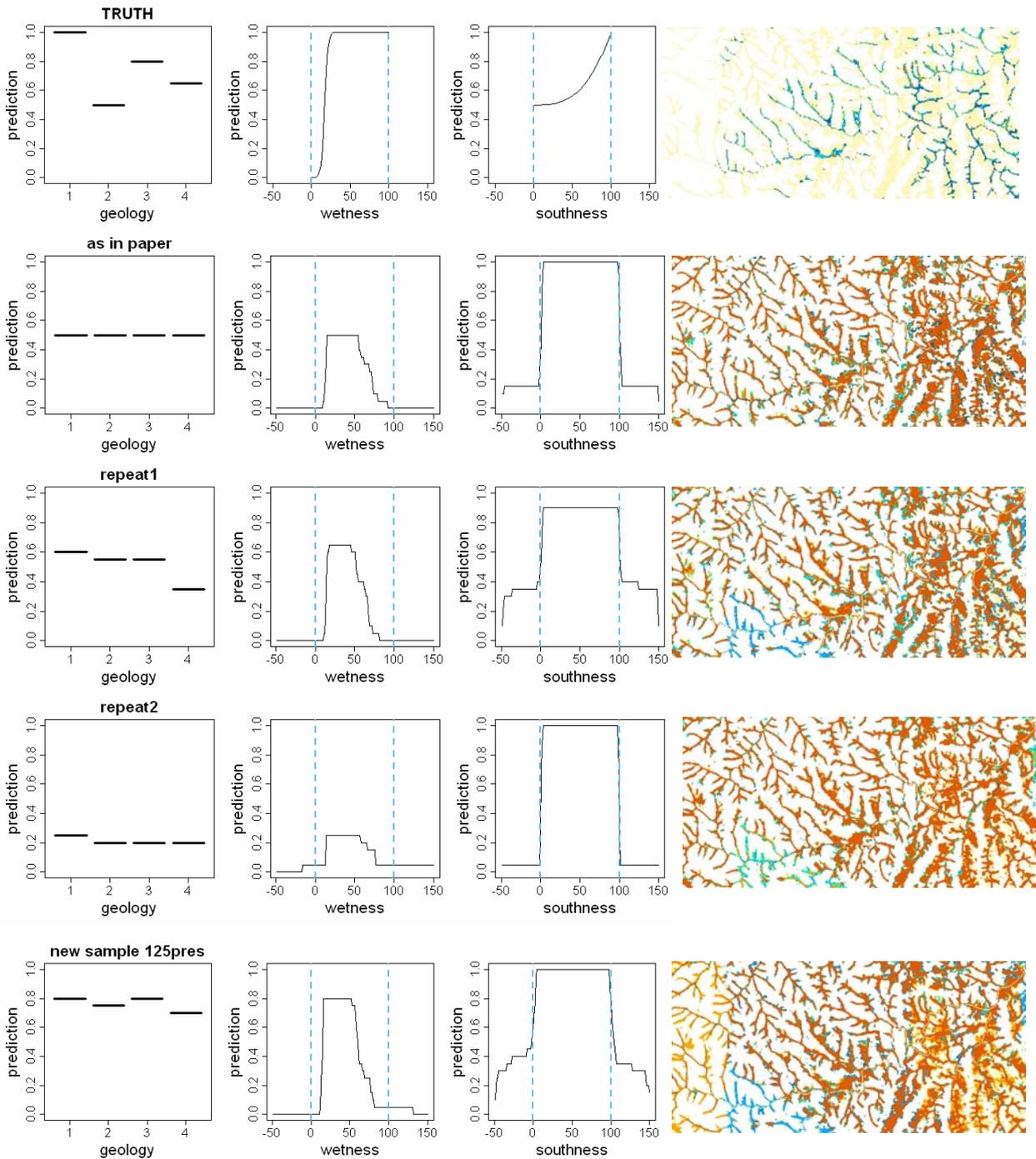


Figure S5: Predictions (y axis) versus the true suitability for all 80000 grid cells in the maps in Figure S4, for the models from Table S1 and Figures S3 and S4. The blue diagonal line shows the 1:1 relationship.



**Figure S6:** Results from 3 independent runs of GARP on the sample of 115 presence records (rows 2 to 4), and a new sample of 125 records. In each 500 models were produced then subsetting as described in the text. The models are based on the mid external omission / mid commission scenarios. Fitted functions (left panels) and mapped distributions (right). The legend for the distribution is the same as that in Figure S4.

**Online supplementary material: Appendix S3: Running and testing Random Forests**

Random Forests (RF) is a machine learning method that builds an ensemble of classification or regression trees. It uses *bagging* (bootstrap aggregation) to form the ensemble, taking a new bootstrap sample of the training data for each new tree. The reason for making many trees (a "forest") is that the variance of single trees, a known problem, is reduced by bagging. A useful side-effect of using bagging is that, at each step, there is an "out-of-bag" sample (i.e. those records not selected) that can be used for testing the model. RF are called "random" forests because at each split only a random subset of the candidate predictors are considered. This de-correlates the trees and improves the variance reduction. Trees are fully grown and not pruned. For regression trees the results are averaged, for classification, each tree casts a vote for the predicted class. For binary data such as ours, classification trees are used but the final votes can give a probability rather than a binary output. Further details on the theory of RF can be found in the publications mentioned below.

RF were run using the R library *randomForest*. JE ran the models and understood the theory of RF but had little experience. Recent publications were read (Prasad *et al.* 2006, Benito Garzon *et al.* 2006, Breiman 2001, Cutler *et al.* 2007), and experts consulted in person or via web pages.

Random forests are generally considered easy to tune. The most important choices are the "*mtry*" and "*ntree*" settings in the R version, representing how many variables are randomly selected at each split of the tree as it is grown (*mtry*), and how many trees are allowed in the ensemble (*ntree*). Rules of thumb are used to estimate a good value for *mtry*, for classification often either  $\sqrt{n}$  (Cutler *et al.* 2007), where  $n$  = number of candidate variables, or  $\log(n)$  (D. Margineau, pers.comm.); *mtry* can be as low as 1. Cutler *et al.* (2007) suggest that *ntree* as low as 50 can be suitable; in R the default is 500, and Prasad *et al.* (2006) used 1000 because it stabilised their results. There is also an R function called *tuneRF* that can be used to set *mtry* in relation to error rates. Cutler *et al.* (2007) comment on it in their appendix: "We have not used this function, in part because the performance of RF is insensitive to the chosen value of *mtry*, and in part because there is no research as yet to assess the effects of choosing RF parameters such as *mtry* to optimize out-of-bag error rates on the generalization error rates for RF".

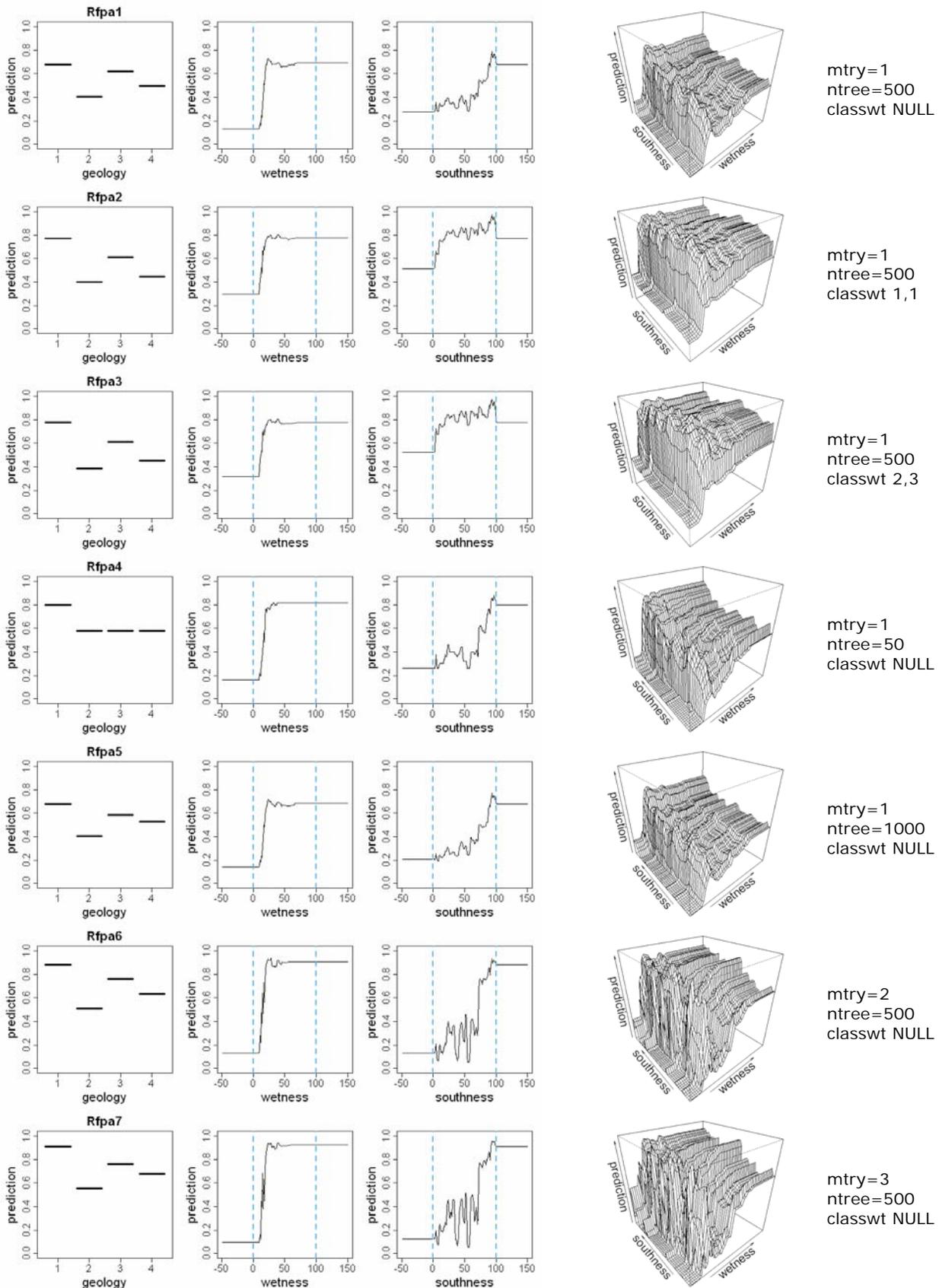
For classification trees, the error rates on the classes in the out-of-bag estimates can be balanced, if this is appropriate for the application, by putting priors on the class weights. In other words, for a binary outcome the model can try to predict "0" as well as it predicts "1".

In our modelling we explored the effect of changing *mtry* and *ntree*, in various combinations. We also looked at the effect of balancing error rates (by use of class weights in R) and tested *tuneRF*. Our results were sensitive to *mtry*, with the best results using *mtry* = 1. This is consistent with  $mtry = \log(3)$  and with the results from *tuneRF*, but not so clearly with  $\sqrt{3} = 1.7$ , which perhaps might have been rounded to 2. Class weights and *ntrees* also affected the outcome. The best results were obtained with the settings used in the paper (*ntrees* = 500, *mtry* = 1 and no class weights) or the comparable model with 1000 trees. These are compared below with examples of some of the other settings tested. The results show that it is important to test settings. Models with *mtry* > 1 were more chaotic than those with *mtry*=1. With 50 trees the response to geology was not modelled properly. Out of bag estimates or cross-validations could be used to systematically test a range of settings to get best predictive performance for the given application.

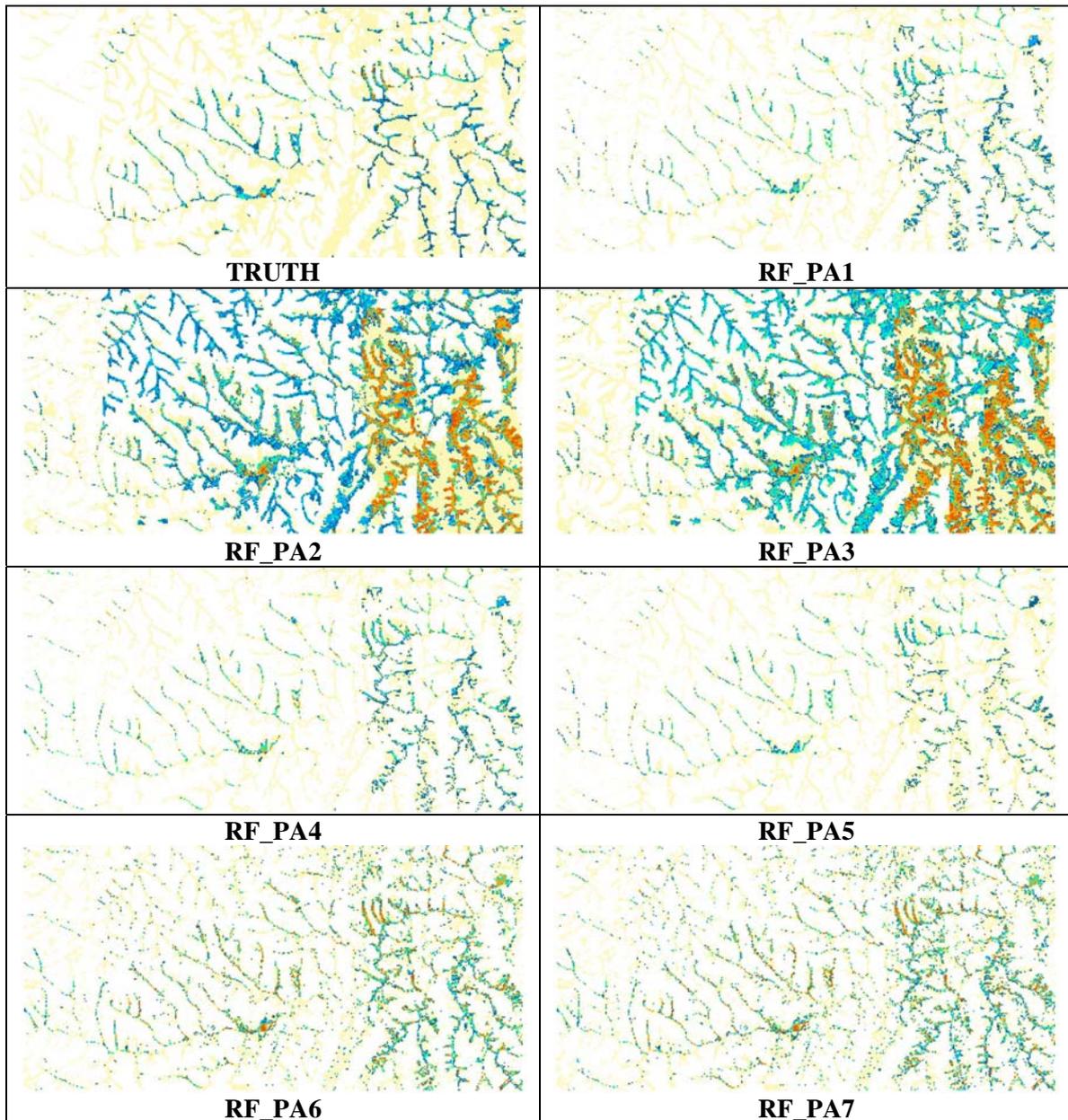
**Table S2:** Details of the models described in the text above, showing the effect of varying parameter settings (first 3 rows) on the error estimates (rows 4 to 6) and evaluation statistics (rows 7 to 10). The meaning of the statistics is described in Elith and Graham (2008).

Model:	rfpa1 (paper)	rfpa2	rfpa3	rfpa4	rfpa5	rfpa6	rfpa7
1. <i>mtry</i>	1	1	1	1	1	2	3
2. <i>ntrees</i>	500	500	500	50	1000	500	500
3. classwt (0,1)	null	1,1	2,3	null	null	null	null
4. oob <sup>1</sup> error overall	0.098	0.184	0.221	0.103	0.105	0.107	0.121
5. oob <sup>1</sup> error class1	0.011	0.168	0.220	0.018	0.008	0.038	0.054
6. oob <sup>1</sup> error class2	0.765	0.304	0.226	0.757	0.852	0.635	0.635
7. auc	0.834	0.843	0.838	0.814	0.835	0.828	0.823
8. remaining deviance	0.736	0.785	0.902	0.948	0.712	0.748	0.769
9. cor.pa	0.448	0.434	0.421	0.438	0.447	0.425	0.412
10. cor.si	0.875	0.853	0.827	0.855	0.874	0.822	0.793

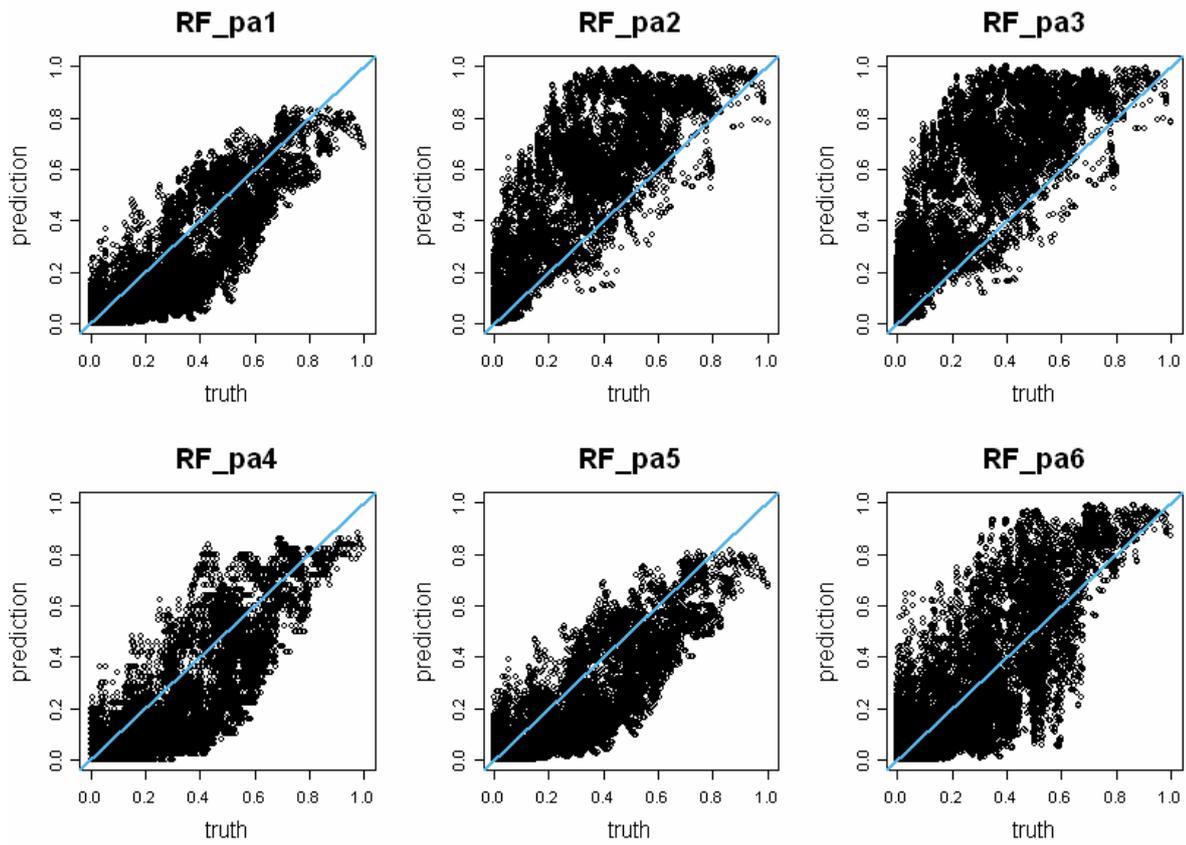
<sup>1</sup> oob = out-of-bag estimate from R



**Figure S7:** Fitted functions for the seven random forest models described in the text and Table 3. Note the effects of changing mtry (rows 1,6 & 7), ntree (rows 1, 4 & 5) and class weights. The first model was presented in Elith and Graham 2008.



**Figure S8:** Mapped distributions of the virtual species (top left) and predictions of relative suitabilities from the methods detailed in the text, Legend: white < 0.1, cream 0.1 to 0.5, blue-lightblue-green-orange-vermillon at steps of 0.1 from blue (0.5 to 0.6) to vermillon (0.9 to 1), as in Figure S3



**Figure S9:** Predictions versus truth for all 80000 grid cells in the maps in Figure S8, for the models from Table S2 and Figures S7. The blue diagonal line shows the 1:1 relationship.

**Online supplementary material: Appendix S4: Running Maxent, BRT and GLM**

Maxent, boosted regression trees (BRT) and generalised linear models (GLMs) were straightforward to run and did not need extended testing, because the modeller was familiar with the methods. Whilst it is possible that performance might have improved with some changes in parameterisation, we were confident that the approaches used were well representative of the capacity of the methods.

Maxent (Phillips *et al.* 2006, Phillips and Dudik 2008) is a machine learning method, using the principles of maximum entropy models to model the species distribution. The idea is to set some constraints that enable the prediction to reflect patterns in the sample, and then select a model that maximises entropy (a uniform or spread out distribution) given that those constraints are met (either exactly or approximately). It considers the region (in this case, the full grid) then models the distribution of the species across that with a density estimation approach. The approach can be thought of as modelling the probability of the covariates (the predictor variables) conditional on species presence. Further information can be found in the papers cited above. Maxent version 3.2.1 was used, and run from the command line (specifically: `-e env -s po.csv -t ge -o res -j proj -r -a -d -P`). Settings for Maxent are presented in the paper, and are mostly the recommended defaults. Because Maxent is set up to model species distributions and the settings have been tested on large data sets, the defaults tend to perform well. The program sets feature selection and regularisation parameters (Phillips *et al.* 2006) in relation to the number of presence records supplied. In this case, with 115 records, it allowed all feature types with fairly strong regularisation (control over) the threshold features. This allows it to model flexible relationships without overfitting. The only exception to the usual defaults is that we used the "-d" flag (see help file provided with the program), which does not add the samples to the background data. This gives Maxent the best chance to have a well calibrated output.

The boosted regression trees were run in R (v. 2.6.1) with the *gbm* library and custom code written by John Leathwick and JE (Elith *et al.* 2008). That paper and others by the authors (e.g. Leathwick *et al.* in press) give detailed descriptions of the method. Briefly, an ensemble of regression trees are formed in a forward stagewise procedure ("stochastic gradient boosting"), where at each step the tree that is added is the one that best explains the residuals from the previous tree(s). The method models binary data accurately by using a logit link function, just as in a GLM. BRT's need careful choice of settings, but once the principles are understood, this is not difficult. We chose settings that would give at least 1000 trees, and that would grow trees deep enough to model interactions. The algorithm uses cross-validation to choose how many trees to add, stopping before it is too overfit (Elith *et al.* 2008). The final model comprised 4250 trees with a learning rate (shrinkage) of 0.001 and a tree complexity of 3.

The generalised linear model (GLM) was run in R version 2.6.1 using function *glm*. For all models we created all possible subsets of models where the allowable fit for each variable was: (1) geology: in as a 4-level factor, or out; (2) wetness and southness: out, in as linear, quadratic or cubic function. From all these models, we selected the model with the lowest Akaike Information Criterion (AIC).

**Online supplementary material: Appendix S5: Using pseudo-absences**

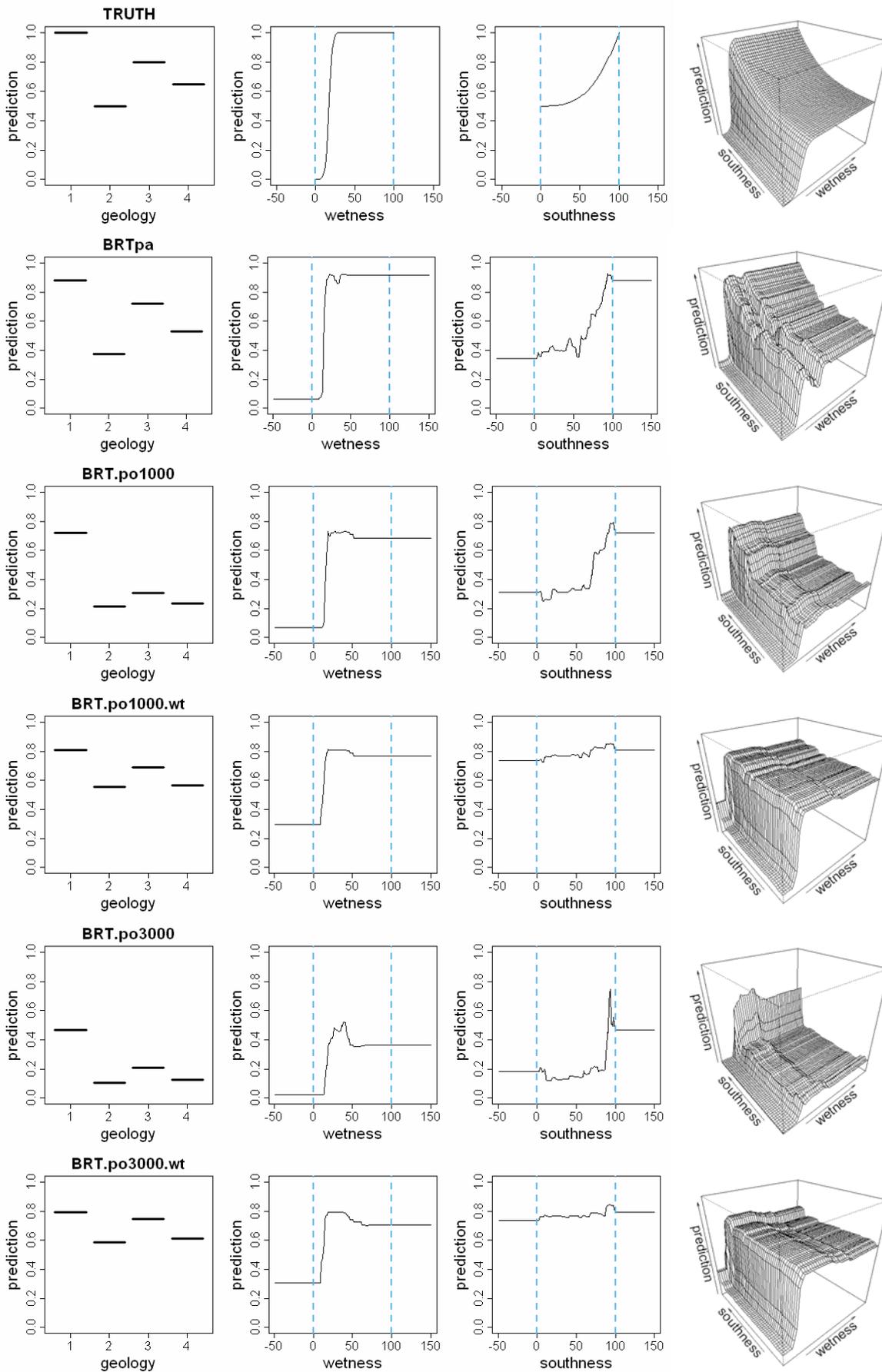
The test with pseudo-absences is a demonstration of the effect on model structure and performance of using presence records and pseudo-absences (appendix S1) in what has been described as a "naïve" model (Ward *et al.* in press, Phillips *et al.* in press). In this, a logistic model is fit to the presence (PO) and background data. If a species is rare, the background data will resemble true absences and the naive model will be close to the true model. But with higher levels of "contamination" (presences in the background sample) the naive model can be biased. Both GLM and BRT are used here as logistic models. The models were fit with the same settings as described in Appendix S4, with the pseudo-absences replacing true absences. For each set of pseudo-absence data (PO.1000 and PO.3000, see Appendix S1), we made two models from each method, in one applying weights on the data so that the sum of the weights on the presence records is the same as the sum of the weights on the absence records (PO.1000.wt etc). This has often been done (e.g Ferrier *et al.* 2002) and produces fitted values and predictions that are distributed across the possible range of the response (here, 0 to 1), rather than predicting many very low values, as occurs if using many more pseudo-absences than presences.

The results are briefly summarised here and we suggest that they are worth pondering in some detail (Table S3; Figures S19 to S15). The results for the unweighted models with the sample of 1000 pseudo-absences were almost as good as those with pa data, mainly because the number of pseudo-absence samples compared with the 115 presences is relatively close to the true prevalence of the species (i.e. there were 885 true absences in the pa data), and because the species is not common in the landscape. With a more common species, contamination of the pseudo-absence sample would have a larger effect. Weighting the data in the models had no effect on the discrimination of the models as long as variable selection wasn't affected. So, for BRT the AUC for the unweighted / weighted pairs are very close, whereas the GLM tended to identify the best model as one without geology when the data were weighted. Dropping geology as a predictor negatively affected all evaluation statistics for the GLMs (Table S3). BRT models the data reasonably in all cases but models the response to southness as much too muted in the weighted models. The GLM never models the response to southness properly at high wetness values (see right side of 3-dimensional plots) but this doesn't affect the evaluation statistics much because there are few data in that part of the environmental space (see Figure 2, Elith and Graham 2008)

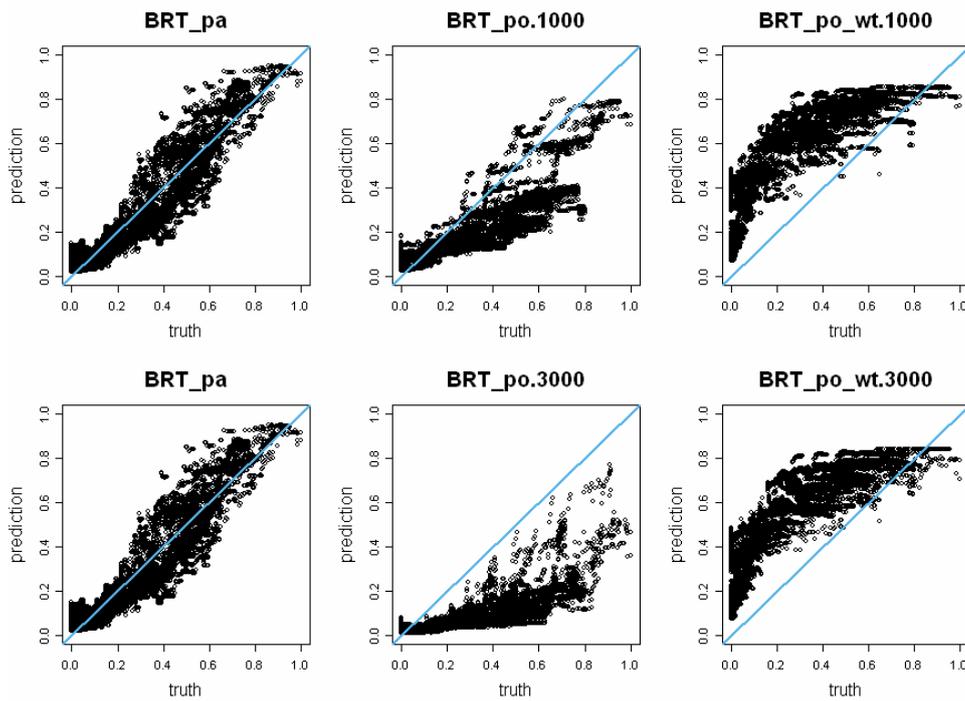
We note that there are statistical solutions to modelling data such as these with more statistical rigour. They are described in Ward *et al.* (in press), and software for running presence-only BRTs with these will be released soon (Gill Ward, pers.comm.).

Table S3: Comparison of model results with truth, as realised by the presence-absence map (columns 2, 3 and 4) and the suitability values (column 5). For all statistics except deviance, higher is better.

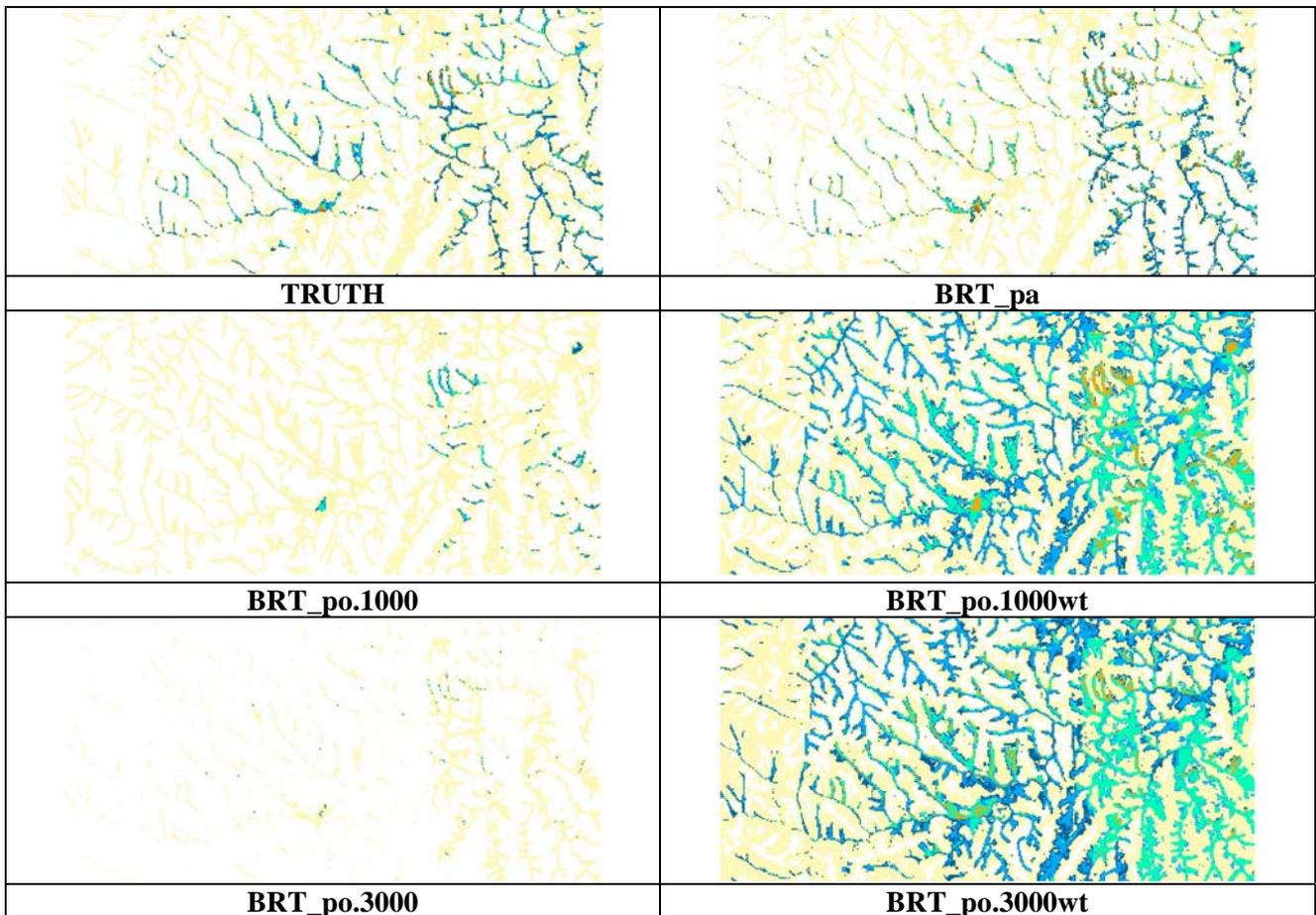
Model	AUC	Remaining deviance	COR.pa	COR.si
Truth (suitabilities)	0.872	0.514	0.508	1.000
Maxent	0.861	0.612	0.467	0.922
BRT.pa	0.862	0.537	0.485	0.954
BRT.po.1000	0.856	0.568	0.464	0.915
BRT.po.1000.wt	0.858	0.842	0.442	0.872
BRT.po.3000	0.855	0.703	0.435	0.851
BRT.po.3000.wt	0.857	0.848	0.441	0.871
GLM.pa	0.863	0.546	0.480	0.941
GLM.po.1000	0.853	0.560	0.468	0.922
GLM.po.1000.wt	0.843	0.924	0.419	0.825
GLM.po.3000	0.855	0.691	0.455	0.896
GLM.po.3000.wt	0.841	0.932	0.417	0.821



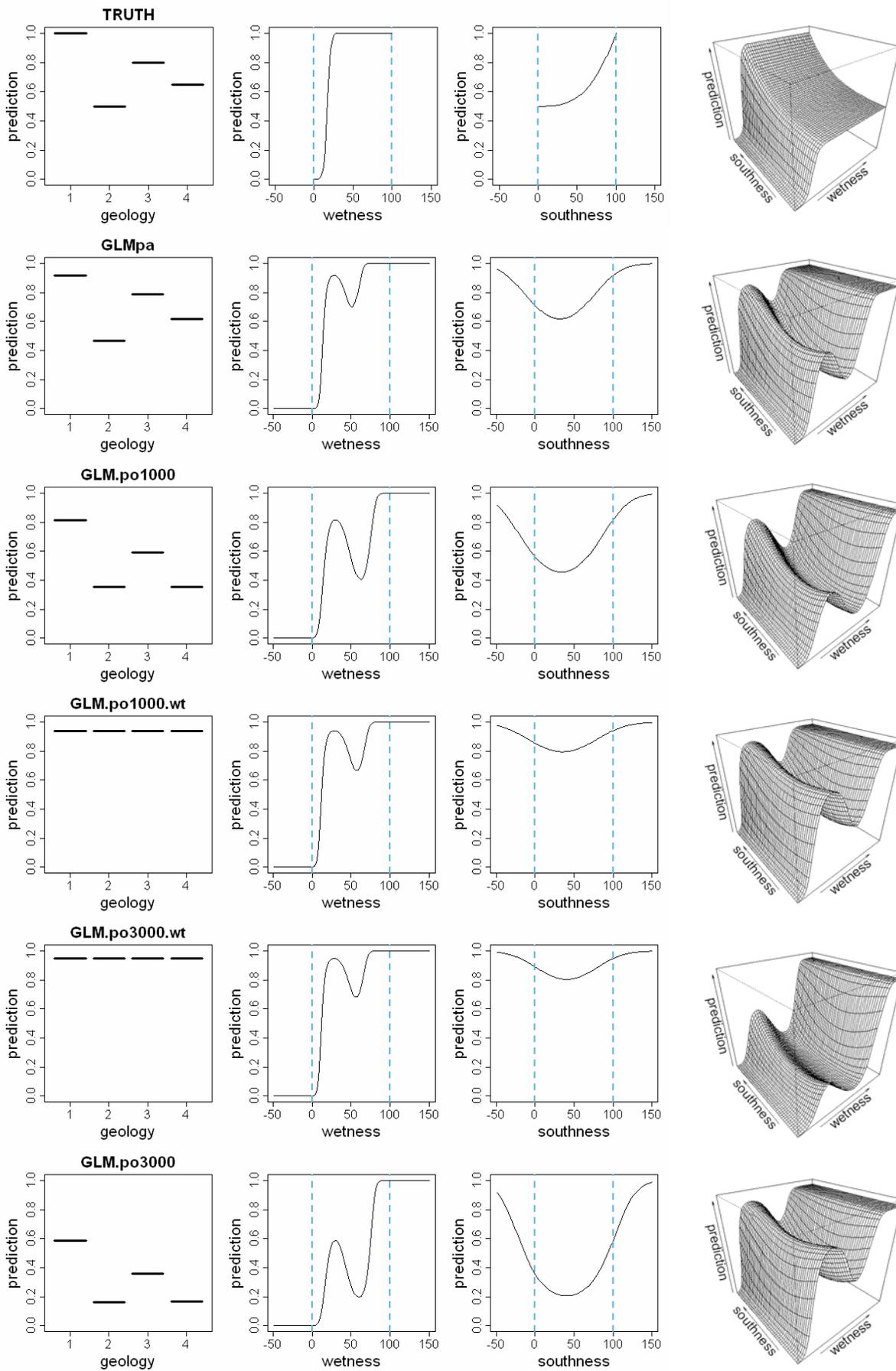
**Figure S10:** Fitted functions for the boosted regression tree models described in the text and presented in Table S3



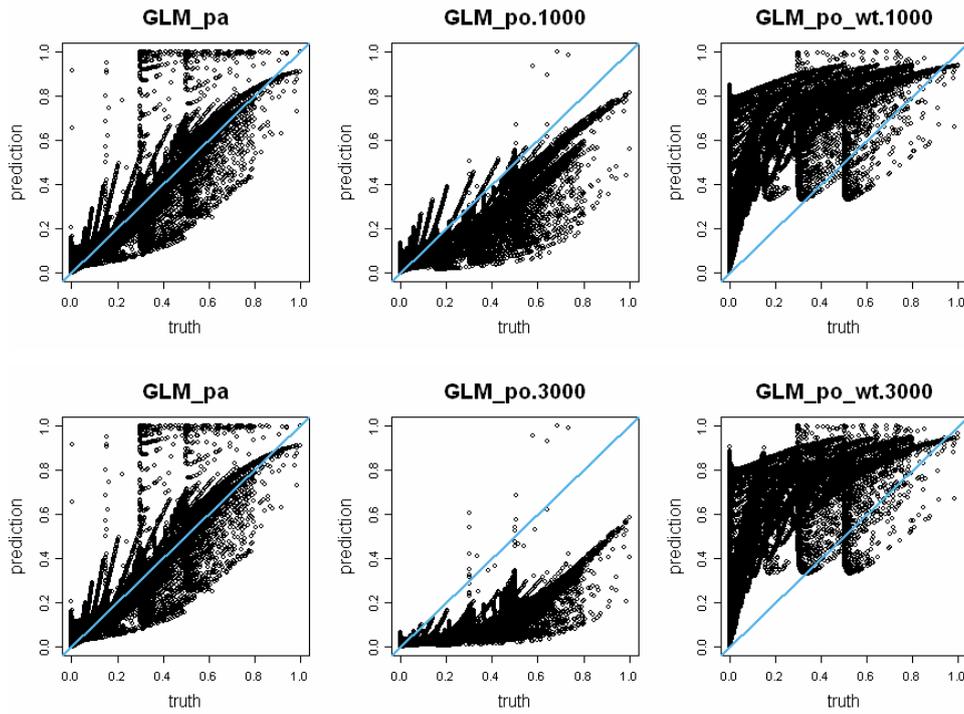
**Figure S11:** Predictions versus truth for all 80000 grid cells in the maps in Figure S12, for the models from Table S3 and Figure S10. The blue diagonal line shows the 1:1 relationship.



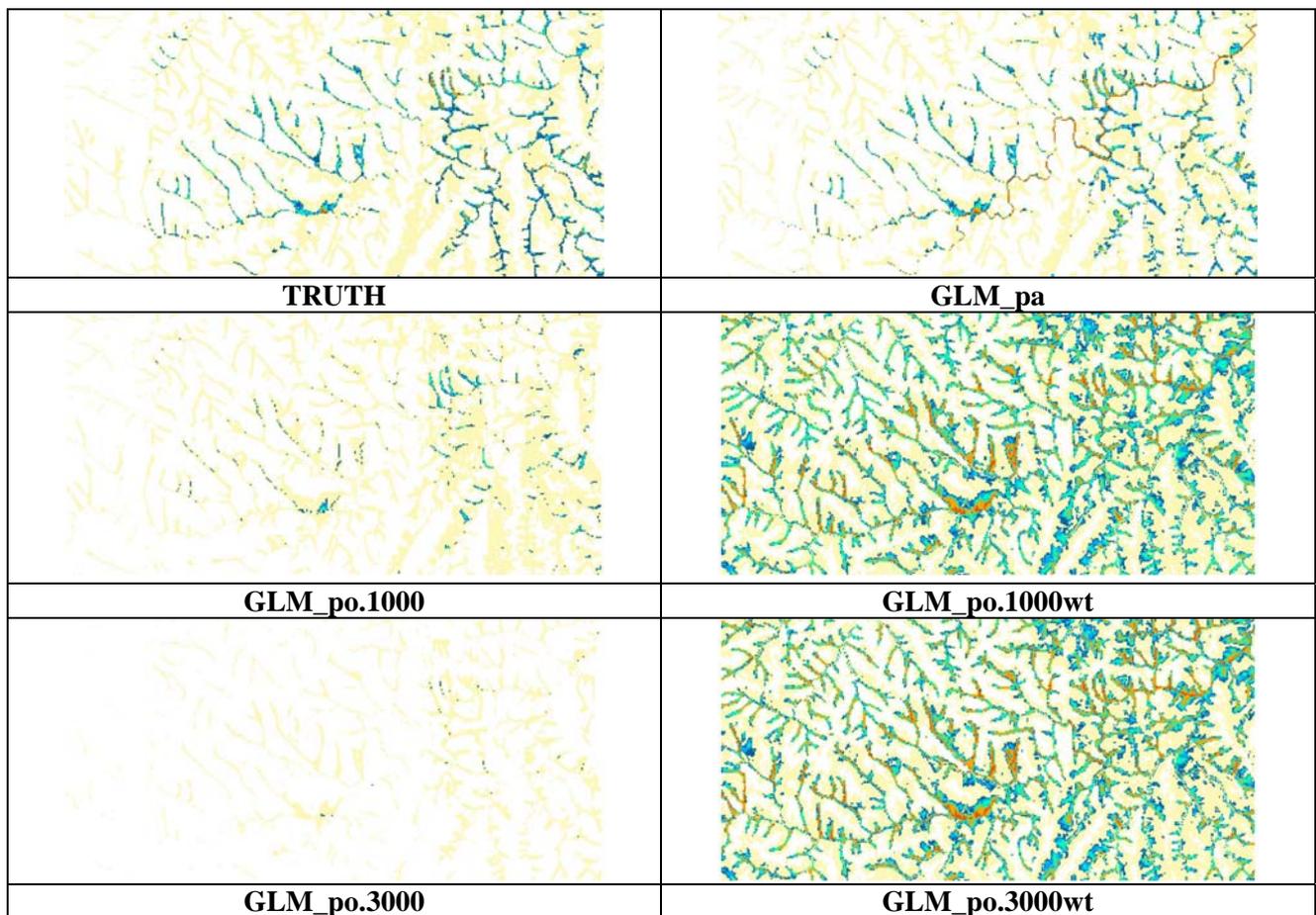
**Figure S12:** Mapped distributions of the virtual species (top left) and predictions of relative suitabilities from the methods detailed in the text, Legend: white < 0.1, cream 0.1 to 0.5, blue-lightblue-green-orange-vermillon at steps of 0.1 from blue (0.5 to 0.6) to vermillon (0.9 to 1), as in Figure S3



**Figure S13:** Fitted functions for the generalised linear models described in the text and presented in Table S3



**Figure S14:** Predictions versus truth for all 80000 grid cells in the maps in Figure S14, for the models from Table S3 and Figure S13. The blue diagonal line shows the 1:1 relationship.



**Figure S15:** Mapped distributions of the virtual species (top left) and predictions of relative suitabilities from the methods detailed in the text, Legend: white < 0.1, cream 0.1 to 0.5, blue-lightblue-green-orange-vermillion at steps of 0.1 from blue (0.5 to 0.6) to vermillion (0.9 to 1), as in Figure S3. Note that the effect of dropping geology as a predictor is to lose the definition of poorer habitat towards the west (left)

## References

- Benito Garzon, M. *et al.* 2006 in press. Predicting habitat suitability with machine learning models: The potential area of *Pinus sylvestris* in the Iberian Peninsula. - *Ecol. Model.* 197: 383-393
- Breiman, L. 2001. Random Forests Technical Report.  
<http://oz.berkeley.edu/users/breiman/randomforest2001.pdf>.
- Cutler, D. R. *et al.* 2007. Random forests for classification in ecology. - *Ecology* 88: 2783-2792.
- Elith, J. and Graham, C. 2008. On the need for a buyers guide to species distribution modelling. - *Ecography*.
- Elith, J., Leathwick, J. R. and Hastie, T. in press. A working guide to boosted regression trees. - *J. Anim. Ecol.*
- Ferrier, S. *et al.* 2002. Extended statistical approaches to modelling spatial pattern in biodiversity: the north-east New South Wales experience. I. Species-level modelling. - *Biodivers. Conserv.* 11: 2275-2307.
- McCullagh, P. and Nelder, J. A. 1989. *Generalized Linear Models*. - Chapman and Hall.
- Peterson, A. T., Papes, M. and Eaton, M. 2007. Transferability and model evaluation in ecological niche modeling: a comparison of GARP and Maxent. - *Ecography* 30: 550-560.
- Phillips, S. J. *et al.* in press. Sample Selection Bias and Presence-Only Models Of Species Distributions. - *Ecol. Appl.*
- Phillips, S. J., Anderson, R. P. and Schapire, R. E. 2006. Maximum entropy modeling of species geographic distributions. - *Ecol. Model.* 190: 231-259.
- Prasad, A. M., Iverson, L. R. and Liaw, A. 2006. Newer Classification and Regression Tree Techniques: Bagging and Random Forests for Ecological Prediction. - *Ecosystems* 9: 181-199.
- R Development Core Team 2006. *R: A Language and Environment for Statistical Computing*. - In, R Foundation for Statistical Computing
- Ridgeway, G. 2006. Generalized boosted regression models. Documentation on the R package "gbm", version 1.5-7. <http://www.i-pensieri.com/gregr/gbm.shtml>.
- Stockwell, D. R. B. and Noble, I. R. 1992. Induction of sets of rules from animal distribution data: a robust and informative method of data analysis. - *Math. Comput. Simulat.* 33: 385-390.
- Ward, G. *et al.* in press. Presence-only data and the EM algorithm. - *Biometrics*.