

1 Title: **Gaussian processes with numerical approximations for fast**  
2 **and flexible Bayesian species distribution modelling**

3

4 Running title: **Gaussian process species distribution modelling**

5

6 Word count: **6,139**

7

8 Authors:

9 **Nick Golding**

10 Spatial Ecology and Epidemiology Group,

11 Department of Zoology,

12 University of Oxford

13

14 and

15

16 Centre for Ecology & Hydrology,

17 Crowmarsh Gifford,

18 Wallingford

19

20 **Bethan V. Purse**

21 Centre for Ecology & Hydrology,  
22 Crowmarsh Gifford,  
23 Wallingford  
24  
25 Corresponding author email: `nick.golding@zoo.ox.ac.uk`

## 26 Abstract

27 **1.** Species distribution modelling (SDM) is very widely used in ecology  
28 and predictions from these models often inform both policy and eco-  
29 logical debates. It is therefore important to use methods with high  
30 predictive accuracy that permit sources of bias to be taken into ac-  
31 count and enable biological interpretation.

32 **2.** Gaussian processes (GPs) are a highly flexible approach to statistical  
33 modelling and have recently been proposed for SDM. Typically, these  
34 models require computationally intensive Markov chain Monte carlo  
35 (MCMC) methods for inference, making them unsuitable for many  
36 SDM applications. We propose fitting GP distribution models us-  
37 ing numerical approximations instead of MCMC in order to overcome  
38 these hurdles and make GPs more widely applicable. We provide an  
39 intuitive introduction to GPs for SDM and demonstrate how the ap-  
40 proach can be used to account for common sources of bias in species  
41 distribution data. Using a dataset of 227 terrestrial vascular plant dis-  
42 tributions within the UK we compare the predictive accuracy of GP  
43 models with those of widely used species distribution models.

44 **3.** Predictive accuracy of GP models fitted using numerical approxima-

45 tions was consistently higher than Boosted regression trees and Gen-  
46 eralized additive models when trained on presence-absence plant data  
47 and greater both of these models plus MaxEnt when trained on presence-  
48 only data. This result was consistent under both randomly-stratified  
49 and geographically stratified training and evaluation datasets.

50 **4.** As well as offering greater predictive power than existing methods,  
51 GP models offer users the ability to account for imperfect occurrence  
52 records, incorporate prior knowledge of the species' ecology and auto-  
53 matically estimate prediction uncertainty. We provide an open source  
54 R package GRaF, to allow ecologists to implement these models.

55 **Keywords:** Gaussian processes, Gaussian random fields, Boosted regres-  
56 sion trees, MaxEnt, Generalized additive models, measurement error, Laplace  
57 approximation, R package

## 58 Introduction

59 Species distribution models (SDMs), in their basic form, attempt to model  
60 the distribution of species using environmental conditions as predictors.  
61 Typically these models make use of records of the distribution of the species  
62 in question and gridded datasets of environmental variables to generate maps  
63 of the species' predicted distribution. In recent years SDMs have become  
64 some of the most widely used methods in ecology (Elith & Leathwick, 2009),  
65 providing essential tools for both theoretical and applied research. Among  
66 other applications, SDMs are used to investigate drivers of global biodiver-  
67 sity patterns and to guide conservation policy and public health interven-  
68 tions (Lehmann *et al.*, 2002; Sinclair *et al.*, 2010; Sinka *et al.*, 2010).

69 A wide range of different approaches has been suggested for SDMs, rang-  
70 ing from relatively simple 'envelope' models and commonly used statistical  
71 methods such as logistic regression to more complex methods such as those  
72 developed in the field of machine learning (Elith *et al.*, 2006). These ap-  
73 proaches have a number of features which determine their suitability to  
74 model species distributions including:

75 **Predictive performance.** Predictive accuracy is likely to depend on  
76 a number of factors, amongst which the ability to model complex (non-

linear) effects of and interactions between drivers of species distributions seems to be particularly important (Elith *et al.*, 2006). Preventing the model from overfitting to the training data (modelling random noise, rather than the true ecological response) will also increase predictive performance when the model is applied to new datasets drawn from different time periods or geographical regions (Wenger & Olden, 2012);

**Imperfect data.** SDMs are often applied to distribution records opportunistically collated from a variety of different sources, rather than from planned surveys. Such datasets are prone to various sources of error, such as observation bias, a lack of absence records and uncertainty in the location or reliability of individual records (Newbold, 2010; Elith *et al.*, 2010a). Failure to account for these sources of error can lead to biased model predictions.

Predicted distribution maps are often needed for species where few occurrence data are available (Pearson *et al.*, 2006). In these cases it may be useful to augment these limited data with existing knowledge of the species' ecology (Murray *et al.*, 2009).

**Computational efficiency.** SDMs are often required to be run in

96 batch operations, where the distributions of multiple species are mod-  
97 elled at once. For routine analyses such as these, SDMs which are  
98 relatively quick to run and do not require much user input are advan-  
99 tageous.

100 **Prediction uncertainty.** SDM predictions typically represent the  
101 model's *best guess* at the species' distribution, given the occurrence  
102 and environmental data available. These predictions are usually sub-  
103 ject to multiple sources of uncertainty in the data and model parame-  
104 ters. It is therefore desirable to provide maps quantifying uncertainty  
105 in these predictions (Elith *et al.*, 2002; Guisan & Zimmermann, 2000).

106 Gaussian processes (GPs; also referred to as Gaussian random fields)  
107 provide a flexible approach to fitting complex statistical models (Rasmussen  
108 & Williams, 2006) and they have seen occasional use in ecology for mod-  
109 elling population dynamics (Patil, 2007; Sigourney *et al.*, 2012). GPs have  
110 recently been proposed as an alternative approach for fitting flexible, non-  
111 linear species distribution models (Vanhatalo *et al.*, 2012).

112 Fitting GP models typically requires the use of computationally expen-  
113 sive Markov chain Monte Carlo (MCMC) methods. Whilst MCMC is a  
114 useful approach to fitting complex models it can be very time consuming  
115 and requires an experienced user to supervise the model fitting process.

116 These limitations make models fitted using MCMC infeasible for the many  
117 SDM users lacking experience with MCMC and for applications which re-  
118 quire running of large batches of models, such as for making predictions of  
119 species richness (Ferrier & Guisan, 2006). We propose that GP models fit-  
120 ted using efficient numerical approximations can overcome these drawbacks  
121 and provide a solution to a number of issues inherent in species distribution  
122 modelling.

123 Below, we illustrate how GP models work, demonstrate how they provide  
124 solutions to some problems commonly encountered in SDM and compare  
125 their predictive ability with other commonly used approaches on a large  
126 dataset of known vascular plant occurrence records. The advantages and  
127 limitations of GPs and potential avenues for future enhancements of the  
128 approach as applied to SDM are discussed. A software package GRaF for  
129 the statistical programming language R (R Development Core Team, 2012)  
130 is also provided to allow ecologists to employ these methods for fitting SDMs.

## 131 **Gaussian process models**

132 While most statistical models attempt to describe the relationship between  
133 covariates and the response variable by parameterising some equation (e.g.  
134 linear regression), GP models instead describe this relationship based on



135 an assumption that observations with similar covariate values will yield a  
136 similar response value. The model then uses the available data to construct  
137 a normally distributed, correlated ‘process’ of variables which give rise to  
138 the observed response variable.

139     The approach of fitting models based on similarity (or dissimilarity) be-  
140 tween sites, rather than fitting directly to the covariates, is shared by a num-  
141 ber of related ‘kernel methods’ such as kernel regression (used in Generalised  
142 dissimilarity matrix models (Ferrier, 2002; Ferrier *et al.*, 2002)), kernel sup-  
143 port vector machines (Evgeniou *et al.*, 2005) and kriging (?). The flexibility  
144 of explicit Bayesian statistical treatment of GP models enables a number  
145 of useful model-fitting procedures and extensions (such as accounting for  
146 uncertainty in covariates, and incorporating prior ecological knowledge, as  
147 detailed below) which would be difficult to implement for many of these  
148 other methods.

149     The GP approach is widely used in the field of model-based geostatis-  
150 tics (Diggle & Ribeiro Jr, 2007), where the covariates are the geographic  
151 location of the observation and the GP therefore models spatial or tempo-  
152 ral correlation. In applying GPs to species distribution modelling, we take  
153 this concept and apply it to *environmental* covariates to model the response  
154 variable - probability of occurrence.

155 A statistical explanation of how such a GP model may be formulated and  
156 fitted is given in Appendix S1. Here we provide a more intuitive illustration  
157 of how GP models differ from other SDMs, using as an example the effect  
158 of temperature on the probability of presence of a hypothetical species.

### 159 **Covariance function**

160 In order to construct the GP we first calculate the environmental distances  
161 between observations. In our example the environmental distances are simply  
162 the difference in temperature between each pair of sites (Fig. 1a). In a  
163 model with more covariates we would calculate the multi-dimensional Euclidean  
164 distance between pairs of sites in environmental space. We convert  
165 these distances into expected correlations in probability of presence between  
166 sites using a *covariance function*. There are a number of different covariance  
167 functions that we could use, but a good choice is the squared-exponential  
168 covariance function since it is easy to parameterise and produces ecologically  
169 plausible smooth curves (Rasmussen & Williams, 2006). As well as the distances  
170 between observations, we must supply the covariance function with a  
171 *lengthscale* parameter for each environmental covariate in the model. These  
172 lengthscales dictate how the correlation between probabilities of occurrence  
173 at pairs of observations decays with environmental distance; and therefore

the complexity of the fitted response curves. Fig. 1b shows how this function converts temperature difference to expected correlation given three different lengthscales. Assuming a lengthscale of one degree celsius, the expected correlation between two observations with a one degree difference in temperature is around 0.6, whereas with a difference of two degrees this drops to around 0.14. With a longer (higher valued) lengthscale, these expected correlations will be higher, resulting in a less complex fitted line (Fig. 2). In practice these lengthscales do not need to be specified in advance as they can be estimated from the data, though we may wish to inform the model of how likely different lengthscales are for the species being modelled.

## Mean function

In addition to these expected correlations, we specify a *mean function*: an initial estimate of how the response variable changes with the covariates which is later updated by the model fitting procedure. If nothing is known about how the probability of presence of a species responds to the temperature gradient, a flat mean function may be used, which assumes an equal probability of presence at all temperatures. If we have some prior knowledge that the species is more likely to be present at low temperatures than at high temperatures, we can incorporate this information into the model. For

193 example, the mean function could be a linear model (with fixed parameters)  
194 relating temperature to probability of presence.

195     Fig. 2 demonstrates the effects of these two different mean functions on  
196 our model, with varying lengthscales. We can see from this illustration that  
197 where there are a sufficient number of observations, the mean function has  
198 little effect on the fitted line, but where there are few data points, such  
199 as toward the limits of the recorded temperature range, the mean function  
200 determines the shape of the fitted response.

## 201 GP species distribution models

### 202 Model structure

203 Machine learning algorithms such as boosted regression trees (BRT; Elith  
204 *et al.* (2008a)) have been shown to perform very well at predicting species  
205 distributions, likely due to their ability to fit complex and highly non-linear  
206 responses to environmental covariates (Elith *et al.*, 2006).

207 A drawback of BRT and similar methods is that they fit ‘jerky’ and  
208 biologically implausible predictive responses which may contribute to their  
209 tendency to overfitting to training data (i.e. they fit to noise in the data  
210 as well as the species’ distribution, Wenger & Olden (2012)). By compar-  
211 ison, more traditional approaches such as univariate generalized additive  
212 models (GAM; Hastie & Tibshirani (1986)) fit more biologically realistic  
213 smooth functions. Whilst some implementations of GAM can fit multi-  
214 variate smoothers (Wood, 2011), they perform poorly in more than a few  
215 dimensions so are unable to account for complex interactions.

216 GPs (using the squared exponential covariance function) offer a attrac-  
217 tive solution to this trade-off between model flexibility and ecological re-  
218 alism by allowing for interactions between and highly non-linear effects of  
219 covariates, whilst fitting biologically plausible smooth predictive surfaces

220 (see Fig. 3).

## 221 **Uncertainty in occurrence data**

222 Modellers often want to make high resolution predictions from high reso-  
223 lution gridded environmental data but are hampered by the low resolution  
224 of the species occurrence data. This mismatch in resolution is problematic  
225 when the modeller needs to extract covariate values corresponding to each  
226 record, since there will be a number of different covariate values from which  
227 to choose. Thus the problem of uncertainty in the record location can be  
228 more usefully considered as uncertainty in the measured value of the envi-  
229 ronmental covariate for each occurrence record. A simple approach to this  
230 problem is to fit the model using the mean of the covariate values. Whilst  
231 straightforward to implement, this approach ignores the uncertainty in the  
232 covariate and can lead to regression dilution, which dampens the apparent  
233 effect of a covariate on the species distribution.

234 The problem of regression dilution has been well studied in statistics  
235 (Frost & Thompson, 2000) and measurement error models have been pro-  
236 posed to deal with this bias in SDMs (McInerney *et al.*, 2011) and other  
237 ecological models (McNamara & Harding, 2004). As well as the mean of the  
238 environmental covariates for each record, measurement error models use an

239 estimate of the error variance (which in the case of spatial uncertainty may  
240 be calculated directly from the multiple covariate values for each record) and  
241 use this information to correct for the regression dilution effect. Whereas  
242 most measurement error models are fitted by MCMC, we can fit these models  
243 in the numerical approximation GP framework at negligible computational  
244 cost by accounting for this error directly within the covariance function (see  
245 Fig. 4). This approach is detailed in full in (Dallaire *et al.*, 2011).

246 GP models also allow users to provide regression weights to individual  
247 records in order to account for the variable reliability of different records or  
248 to account for observation bias (Phillips *et al.*, 2009) in a similar way to the  
249 case weights and bias grids available in other SDMs (Elith *et al.*, 2010a).

## 250 **Incorporating prior ecological knowledge**

251 Often when modelling species distributions with few occurrence data there  
252 are other forms of information about the ecology of the target, or similar,  
253 species which could be used to improve the model. For example experimental  
254 studies may have demonstrated a relationship between the species' ability to  
255 persist and some environmental gradient. In such cases it may be desirable  
256 to incorporate this prior knowledge of the species ecology into the model to  
257 augment the occurrence data. Unfortunately this is not easily accomplished

258 in many current SDMs.

259 Bayesian statistical inference provides a convenient way of incorporating  
260 prior information of this sort into statistical models and has become increas-  
261 ingly popular in ecology (see e.g. McCarthy (2007)). In a Bayesian model  
262 the user specifies a prior probability distribution over each model parameter,  
263 representing their existing knowledge about what values of the parameter  
264 are likely. The model then compares this prior probability with the parame-  
265 ter estimate suggested by the data and produces a form of weighted average  
266 over the two; the posterior distribution. As the amount of data available to  
267 the model increases, the impact of the prior on the posterior diminishes.

268 The GP framework allows the user to incorporate ecological knowledge  
269 into distribution models by manipulating two Bayesian priors: the mean  
270 function and the lengthscale hyperprior. The mean function acts as a prior  
271 over the whole model and can be used to incorporate specific knowledge of  
272 the species' response to environmental gradients. The lengthscale hyperprior  
273 determines how likely different lengthscales are and can be used to inform the  
274 model how rapidly probability of presence is likely to change with different  
275 values of the environmental covariates.

276 In the absence of any prior information, a flat mean function can be used,  
277 as in Fig. 2. Similarly a flat lengthscale hyperprior could be used, indicating



278 that all levels of complexity in the fitted response are *a priori* assumed to be  
279 equally likely. By default the R package GRaF which we provide uses a flat  
280 mean function but an informative lengthscale hyperprior which represents  
281 ecologically plausible response curves (detailed in Appendix S1).

## 282 **Uncertainty in model predictions**

283 As with any model, predictions from SDMs are uncertain estimates of the  
284 probability of presence of the species. Where these predictions are to be  
285 used for some practical purpose it would be beneficial to provide maps rep-  
286 resenting the uncertainty in the predicted distribution map (Elith *et al.*,  
287 2002). Such maps allow users to determine how much confidence they can  
288 place in a given prediction, information which is especially valuable where  
289 the predictions have policy implications.

290 SDM uncertainty estimates can be produced by bootstrapping data (Elith  
291 *et al.*, 2002) though this requires models to be run many hundreds of times  
292 and can therefore be computationally prohibitive. GP models automatically  
293 produce estimates of uncertainty in model predictions, without the need for  
294 bootstrapping procedures, since these are calculated directly from the esti-  
295 mated posterior distribution of the model. Fig. 5 illustrates predictions and  
296 associated uncertainty estimates from a GP distribution model of a plant

<sup>297</sup> species, Bog Myrtle (*Myrica gale*), in the UK.

## 298 **Comparison of GPs with existing SDMs**

299 We compared the predictive ability of GPs (fitted using GRaF) with com-  
300 monly used approaches for modelling species distributions from both pres-  
301 ence/absence and presence-only data under at both interpolation and ex-  
302 trapolation tasks. All model fitting, predictions and statistical analyses  
303 were performed in R version 3.0.0. R code used to carry out these analyses  
304 and to plot all figures in this manuscript are provided in Appendix S5. The  
305 dataset used to perform these comparisons is available to download from  
306 Figshare [DOI/URL tbc].

## 307 **Methods**

### 308 **Data**

309 Gridded presence/absence maps of native terrestrial vascular plant species  
310 of Great Britain at 10 km resolution were obtained from the New Atlas of  
311 the British and Irish Flora (Preston *et al.*, 2002). Of the 1335 distribu-  
312 tions in the original dataset, a subset of 227 species of different genera was  
313 selected for the model comparison. Criteria for selection of these species  
314 are outlined in Appendix S3. The distribution of British plant species is  
315 well characterized at this spatial scale, so records were assumed to represent  
316 known presence or absence which is essential to compare models fairly. The

317 227 plant species selected had a wide range of prevalences (proportion of  
318 grid squares occupied; ranging from 0.075 to 0.925, median 0.394), a fac-  
319 tor known to influence the accuracy of SDMs (McPherson *et al.*, 2004) and  
320 inhabited a wide range of habitats (see Appendix S3).

321       Gridded maps of 10 indices of environmental conditions were used as  
322 covariates for model fitting and prediction. These were derived from a time  
323 series of satellite images of the UK by subsequent Fourier decomposition and  
324 principal components analysis to produce variables representing the major  
325 axes of environmental variation in the UK. The advantage of these abstract  
326 indices is that they compress a large amount of information on conditions  
327 and seasonality (the ten principal components explain 90% of variation in  
328 the Fourier variables) into relatively few orthogonal variables which enable  
329 us to make accurate predictions (Dormann *et al.*, 2008). Whilst they are  
330 difficult to interpret biologically, the first three indices broadly correspond  
331 to gradients from arable land to pasture, lowlands to uplands, and urban  
332 areas to rural areas respectively. Details of this dataset and how it was  
333 produced are given in Appendix S3. All models were fitted using the full  
334 set of 10 covariates.

335       A total of 2774 10 km grid squares contained both distribution data  
336 and environmental data and were used to compare the different modelling

337 approaches.

### 338 **Presence/absence models**

339 For the presence/absence comparison we compared GPs with BRT and  
340 GAM. For each of the 227 species 300 grid squares (10.8% of the 2774  
341 available records) were used to train each of the three models. This number  
342 of observations was selected as a practical compromise between having a  
343 sufficient training set, and the need to have sufficient presence and absence  
344 records in the training and evaluation datasets whilst accommodating the  
345 very rare and very common species.

346 Model predictions for the remaining 2474 grid squares were then used to  
347 compare the predictive ability of the models. GP models were fitted using  
348 GRaF 0.1-12 (Golding, 2013), optimising the lengthscale parameters and  
349 otherwise using default settings. BRT models were fitted using the gbm  
350 package version 2.1 (Ridgeway, 2013) with 5-fold cross validation, a tree  
351 complexity of 5, a learning rate of 0.001 and a minimum of 1000 trees were  
352 fitted (in accordance with Elith *et al.* (2008a)). The optimal number of trees  
353 in the final BRT model was selected from the cross-validation folds using the  
354 gbm.perf function. GAM models were fitted using the gam package version  
355 1.09 (Hastie, 2011) with univariate spline smoothers for each covariate and

356 default settings.

### 357 **Presence-only models**

358 For the presence-only data we compared GPs with MaxEnt (Phillips *et al.*,  
359 2006), one of the most widely used approaches for modelling species dis-  
360 tributions (Yackulic *et al.*, 2012). Since absence data is not available in  
361 presence-only datasets, it is necessary to augment the presence data with  
362 a set of 'background data' against which the models may contrast the en-  
363 vironmental signature of the presence data. For each species 100 presence  
364 records were selected from the dataset along with a further 1000 background  
365 points selected from the remaining data. As for the presence/absence case,  
366 this number of sampling points was a compromise between accommodating  
367 a range of species prevalences and the number of available records; though  
368 real-world applications of presence-only SDM typically use similar numbers  
369 of occurrence points (studies reviewed by Yackulic *et al.* (2012) used a me-  
370 dian of 146 occurrences).

371 GP and MaxEnt models were fitted to these 1100 data points and used  
372 to predict the relative probability of presence at the remaining 1674 grid  
373 squares not used to fit the models. GP models were all fitted as for presence-  
374 absence data and MaxEnt models were fitted using dismo version 0.8-11

375 (Hijmans *et al.*, 2012) with default settings.

## 376 **Defining the evaluation dataset**

377 SDMs are widely used both for interpolation ('filling in' the species distri-  
378 bution in the region from which occurrence data is available) and extrap-  
379 olation (prediction to geographically distinct regions) (Elith & Leathwick,  
380 2009). Previous comparisons of SDM approaches have focussed the inter-  
381 polation capacity of these models by selecting an evaluation set at random  
382 from the available data (Elith *et al.*, 2006). This procedure is unlikely to ac-  
383 curately reflect the extrapolation capacity of the various models and where  
384 the available data are spatially dense, validation statistics calculated using  
385 this procedure are likely to be artificially inflated as a result of spatial au-  
386 tocorrelation between the training and evaluation sets (Wenger & Olden,  
387 2012). We evaluated the interpolation and extrapolation capacities of the  
388 different approaches by running model comparisons using both a randomly  
389 stratified and a geographically stratified dataset for the presence/absence  
390 and presence-only test sets.

391 For the randomly stratified tests, training and evaluation data were sam-  
392 pled at random from the dataset, subject to the constraint that at least ten  
393 presence and ten absence points were available in the evaluation set and the

394 same minimum number present in the training set for the presence/absence  
395 comparison.

396 For the geographically stratified tests, a disc with a radius of 250km was  
397 used to divide the training and evaluation datasets. All grid cells falling  
398 within this region were used as evaluation data. Grid cells outside both  
399 this disc, and an additional buffer region of 20km around the disc were used  
400 as the pool from which to draw the training data. Cells falling within the  
401 buffer region were discarded. For each species and test (presence/absence  
402 or presence-only) the disc was centred on a different randomly selected grid  
403 cell, subject to the constraint that at least ten presence and ten absence  
404 points were available in the evaluation set and in the training dataset for  
405 the presence-absence test. This procedure was carried out using a rejection  
406 algorithm (provided as R code in the supplementary material), incorporating  
407 functions from the *sperrorest* R package version 0.2-1 (Brenning, 2012).

#### 408 **Validation statistics**

409 When fitting presence-absence statistical models to presence-background  
410 data, the background data is implicitly assumed to represent absence of  
411 the species. As a result, these models are subject to prediction error caused  
412 by 'contamination' of background points with those in which the species is



413 in fact present (Ward *et al.*, 2009; Elith *et al.*, 2010b). MaxEnt, by contrast,  
414 explicitly considers this data as background data and is unaffected by this  
415 form of bias. However, since the prevalence of each species in the study  
416 area is unidentifiable from such data (Ward *et al.* (2009); Phillips & Elith  
417 (2013) - unlike in the presence-absence case), predictions from both MaxEnt  
418 and presence-absence models applied to presence-only data, are not of the  
419 absolute probability of presence of the species, but only an uncalibrated or  
420 *relative* probability of presence (Elith *et al.*, 2010b).

421 We therefore assessed predictions from presence-absence and presence-  
422 only models using different validation statistics. For presence-absence mod-  
423 els calculated the log-likelihood of the withheld data from the predictions  
424 of each model. The Log-likelihood accurately assesses the quality of prob-  
425 abilistic predictions for predicting the true probability of presence and is  
426 preferable where it can be applied (Lawson *et al.*, 2013).

427 For presence-only models we calculated the Area Under the Receiver  
428 Operating Statistic Curve (AUC; calculated using the pROC R package  
429 - version 1.6.0.1; Robin *et al.* (2011)) which assesses the ability of each  
430 each model to correctly rank sites in order of probability of presence and  
431 is therefore preferable where predictions are of the relative probability of  
432 presence (Lawson *et al.*, 2013).

433 For both statistics higher scores indicate a better fit to the data.

## 434 **Statistical analysis**

435 These validation statistics were analysed by mixed-effects regression, imple-  
436 mented using the nlme R package version 3.1-113 (Pineiro *et al.*, 2012).  
437 In each regression the response variable was the metric of predictive perfor-  
438 mance (log-likelihood or AUC) and the covariates were the SDM model type  
439 (modelled as a fixed effect) and plant species (modelled as a random effect  
440 in order to account for the nested study design). As the residual variances  
441 differed between model types, a separate variance parameter was estimated  
442 for each SDM model type. The statistical significance of differences be-  
443 tween model validation statistics was assessed by t-tests on coefficients for  
444 the SDM model type.

445 In the geographically-stratified presence-absence model comparison, GAM  
446 models for 11 species had very low log-likelihood scores (all less than -3000)  
447 and had to be omitted to enable the mixed effects regression model to con-  
448 verge. These low scores are indicative of a failure of the GAM fitting algo-  
449 rithm to converge to a sensible result, and would likely be rejected by an  
450 SDM user during the model fitting process in a real-world application of  
451 the method. Exclusion of these models resulted in higher average prediction

452 metrics for GAMs; however as the validation statistics for GAM models in  
453 this comparison were markedly lower than other models, the results of the  
454 comparison are still robust.

455 Marginal validation statistic scores were calculated from the residuals  
456 of null models with an intercept term and random effects terms for plant  
457 species, but no fixed effect of model type. These marginal statistics enable us  
458 to visualise the expected predictive capacity from each SDM whilst remov-  
459 ing species-level effects, essentially representing likely model performance  
460 expected for a 'typical' species in the dataset, averaging out the effects of  
461 species-specific ecology or of the species' prevalence.

## 462 **Results**

463 GP models made more accurate predictions to the withheld data than other  
464 models for both presence-absence and presence-only data and both under  
465 randomly and geographically stratified training/evaluation sets (Fig. 6).

### 466 **Random stratification**

467 Predictive log-likelihoods for presence/absence GP models under random  
468 stratification were  $34.95 (\pm 1.68 \text{ SE}, t_{452} = 20.82, p < 0.0001)$  higher than  
469 for BRT and  $157.89 (\pm 6.49 \text{ SE}, t_{452} = 24.35, p < 0.0001)$  higher than for

470 GAM models.

471 Similarly, presence-only GP models had an average AUC score of 0.81;  
472 0.024 ( $\pm$  0.002 SE,  $t_{678} = 14.58$ ,  $p < 0.0001$ ) higher than for BRT models;  
473 0.019 ( $\pm$  0.001 SE,  $t_{678} = 17.81$ ,  $p < 0.0001$ ) higher than for GAM models  
474 and 0.018 ( $\pm$  0.001 SE,  $t_{678} = 14.08$ ,  $p < 0.0001$ ) higher than for MaxEnt  
475 models.

476 GP models explained an average of 28.6% of null deviance for the pres-  
477 ence/absence evaluation sets, compared with 25.8% for BRT and 0.17% for  
478 GAM models.

## 479 **Geographic stratification**

480 Predictive log-likelihoods for presence/absence GP models under geographic  
481 stratification were 16.31 ( $\pm$  5.08 SE,  $t_{441} = 3.21$ ,  $p < 0.0014$ ) higher than for  
482 BRT and 231.46 ( $\pm$  22.50 SE,  $t_{441} = 10.29$ ,  $p < 0.0001$ ) higher than for GAM  
483 models.

484 Under geographic stratification, presence-only GP models had an average  
485 AUC score of 0.68; 0.018 ( $\pm$  0.004 SE,  $t_{678} = 4.38$ ,  $p < 0.0001$ ) higher than  
486 for BRT models; 0.008 ( $\pm$  0.003 SE,  $t_{678} = 2.31$ ,  $p < 0.0212$ ) higher than for  
487 GAM models and 0.021 ( $\pm$  0.004 SE,  $t_{678} = 5.83$ ,  $p < 0.0001$ ) higher than for  
488 MaxEnt models.

489        Under the more stringent geographic stratification test, all models ex-  
490    plained less of the null deviance in the presence/absence evaluation sets  
491    than under the random stratification, with GP models explaining an aver-  
492    age of 11.6%, BRT models 9.3% and GAM models 3.6%.

## 493 Discussion

### 494 SDM comparison

495 In our comparison GP SDMs clearly outperformed a number of popular  
496 SDM approaches, including BRT which has been shown to be one of the  
497 best performing of existing SDM approaches (Elith *et al.*, 2008b).

498 In this comparison, we fitted each model following best-practice guide-  
499 lines where available and default settings otherwise. We also compared mod-  
500 els across a very large dataset of species distributions with varying ecologies  
501 and prevalences. We therefore consider this to be a fair comparison of the  
502 SDMs considered. However further comparisons of these methods using dif-  
503 ferent datasets and by different modellers, would be useful in order to eval-  
504 uate the performance of GP species distribution models across a broader  
505 swathe of SDM applications.

### 506 Advantages of a Bayesian approach

507 The ability of GP models to incorporate prior ecological knowledge and  
508 to account for uncertainty in occurrence locations stems from the use of  
509 a Bayesian statistical approach. These features are likely to prove useful  
510 where there are few occurrence records (Murray *et al.*, 2009) and where

511 there is a well-understood environmental driver of a species' distribution,  
512 such as a temperature limit on the distribution of a pathogen (Gething  
513 *et al.*, 2011). This method of incorporating prior knowledge could also be  
514 used to integrate process-based ecological models (Dormann *et al.*, 2012)  
515 with the more commonly used correlative SDMs.

516 By taking a Bayesian approach, GPs can produce estimates of uncer-  
517 tainty in model predictions, by considering a probability distribution over  
518 the predicted values (accounting for uncertainty in the shape of the GP)  
519 rather than making a single 'best guess' prediction, as is the case with most  
520 SDMs.

521 Whilst predictions from GPs fitted using the numerical approximation  
522 procedures described here (and implemented in the GRaF R package) ac-  
523 count for uncertainty in the shape of the GP (allowing us to produce credible  
524 intervals around species response curves), they do not account for uncer-  
525 tainty in the lengthscale hyperparameters (which control how dynamic the  
526 effects of environmental covariates are on probability of presence), but are  
527 conditional on an optimum estimate calculated from the dataset. Condi-  
528 tional posterior predictions are likely to satisfy most SDM users' require-  
529 ments since most widely used SDMs do not provide any measure of predic-  
530 tion uncertainty, let alone accounting for uncertainty in hyperparameters. If

531 required, predictions accounting for uncertainty in hyperparameters can be  
532 approximated by numerical integration (e.g. a deterministic algorithm as in  
533 (Rue *et al.*, 2009) or by Monte Carlo). Such a procedure will inevitably be  
534 more computationally intensive, but is likely to be far more efficient than  
535 alternative approaches such as MCMC.

## 536 **Computational efficiency**

537 GP models fitted using numerical approximations are reasonably compu-  
538 tationally efficient, with model fitting times to the 300 presence-absence  
539 datasets in our comparison similar to those of BRT models (GP: mean 39.95  
540 seconds  $\pm$  9.32 SD cf. BRT: 10.81 seconds  $\pm$  0.89 SD,  $n = 227$ ). GPs are  
541 fairly efficient for datasets of up to a few thousand data points and running  
542 our implementation on a desktop computer is likely to be sufficiently fast  
543 for the majority of users.

544 This computational efficiency is achieved by carrying out an efficient  
545 numerical approximation routine (by default Laplace approximation) to fit  
546 the model to the data. Predictions resulting from these models are there-  
547 fore subject to a degree of approximation error which may limit their per-  
548 formance. In Appendix S2 we discuss this issue further and compare the  
549 accuracy and of these numerical approximations with those of an MCMC



550 approach which would be computationally and technically prohibitive for  
551 the vast majority SDM applications, by users without MCMC experience  
552 and limited computing resources.

553     A downside to GP models is that in the naive case they scale cubically with  
554 the size of the dataset (due to multiple matrix decompositions of  $\mathcal{O}(n^3)$  com-  
555 plexity), so for very large data sets, GP models can be disproportionately  
556 slow. For users who wish to fit GP models to very large datasets efficiently,  
557 substantial speed-ups can be achieved by exploiting parallel computing. Our  
558 R implementation GRaF uses R’s base functions for linear algebra (which in  
559 turn call third-party linear algebra libraries) and consequently these compu-  
560 tations can be efficiently parallelised, without any additional coding, simply  
561 by linking R to a parallel linear algebra library (see e.g. Schmidberger *et al.*  
562 (2009)). Planned updates to GRaF include optional sparse approximation  
563 methods (see e.g. Vanhatalo *et al.* (2010)) which should enable major im-  
564 provements in computational efficiency for large models.

## 565 **Model complexity**

566 GP models can account for highly complex interactions between covariates,  
567 a feature which probably contributes greatly to their strong predictive per-  
568 formance. When using complex models such as this it is important to avoid

569 overfitting to the training data - finding patterns in random noise and there-  
570 fore producing biased predictions to new datasets. Common approaches to  
571 dealing with this problem include covariate selection procedures, which seek  
572 to find a parsimonious subset of the available covariates which explain the  
573 data well; and regularization (e.g. the ‘lasso’ which is used in MaxEnt,  
574 among others (Tibshirani, 1996)) which is used during model fitting to in-  
575 clude only the most important covariates.

576 It is common to tune hyperparameters of GP models using the model  
577 marginal likelihood (Rasmussen & Williams (2006); this is the approach  
578 taken in GRaF implementation) an approach which automatically reaches  
579 an optimal trade off between model fit and complexity. GPs therefore reduce  
580 the influence of environmental covariates with little explanatory power and  
581 prevent model overfitting, without requiring the user to carry out a covariate  
582 selection procedure.

583 A downside of fitting models with high-dimensional and non-linear in-  
584 teractions is that it becomes harder for the user to interpret the relationship  
585 between the species and its environment. This problem is not unique to  
586 GPs, but is an inevitable trade-off when modelling complex data.

587 The GRaF R package provides functions to help users to interrogate  
588 GP models. The effects of individual covariates on probability of presence

589 can be visualised, as can two-way interactions between covariates. Mod-  
590 els can also be compared using deviance information criteria (Spiegelhalter  
591 *et al.* (2002); an equivalent of commonly used information criteria, such as  
592 Akaike’s (Akaike, 1973), for hierarchical models) to quantify the relative  
593 importance of covariates in driving the species’ distribution. A worked ex-  
594 ample modelling the distribution of Bog Myrtle and demonstrating these  
595 features is provided in Appendix S4.

## 596 **Future development of the GRaF R package**

597 Though GRaF is currently designed to fit the kinds of SDMs that are most  
598 commonly used at present, the approach could be extended in a number  
599 of ways. Since latent GP models include, as a special case, a large sub-  
600 set of generalised linear mixed-effects models (Rue *et al.*, 2009), GRaF  
601 could easily be extended to model abundance data, nested study designs  
602 and spatio-temporal autocorrelation. Multiple-response GP models could  
603 also be implemented to allow users to fit models for whole communities of  
604 species, parameterising and accounting for correlations between their dis-  
605 tributions (Wisz *et al.*, 2013; Kissling *et al.*, 2011). GraFs open source R  
606 implementation will facilitate these extensions.

## 607 **Acknowledgements**

608 The authors acknowledge funding from the NERC Centre for Ecology & Hy-  
609 drology (CEH) Environmental Change Integrating Fund Programme. We  
610 thank Chris Preston and David Roy from the Biological Records Centre,  
611 CEH, Wallingford and the Botanical Society of the British Isles for provid-  
612 ing access to the plant atlas data and kindly agreeing to its dissemination  
613 with this paper. We thank the maintainers of the CEH NEMESIS com-  
614 puting cluster and David Rogers, Miles Nunn, Luigi Sedda, David Harris,  
615 Bob O'Hara and Marianne Sinka who all provided helpful comments on the  
616 manuscript.

## 617 References

- 618 Akaike, H. (1973) Information Theory and an Extension of the Maximum  
619 Likelihood Principle. B.N. Petrov & F. Caski, eds., *Proceedings of the*  
620 *Second International Symposium on Information Theory*, pp. 267 – 281.  
621 Budapest.
- 622 Brenning, A. (2012) Spatial cross-validation and bootstrap for the as-  
623 sessment of prediction rules in remote sensing: the r package 'sperror-  
624 est'. *IEEE International Symposium on Geoscience and Remote Sensing*  
625 *IGARSS*. In press.
- 626 Dallaire, P., Besse, C. & Chaib-draa, B. (2011) An approximate inference  
627 with Gaussian process to latent functions from uncertain data. *Neuro-*  
628 *computing*, **74**, 1945–1955.
- 629 Diggle, P.J. & Ribeiro Jr, P.J. (2007) *Model-based geostatistics*. Springer,  
630 New York.
- 631 Dormann, C.F., Purschke, O., García Márquez, J.R., Lautenbach, S. &  
632 Schröder, B. (2008) Components of uncertainty in species distribution  
633 analysis: a case study of the Great Grey Shrike. *Ecology*, **89**, 3371–86.
- 634 Dormann, C.F., Schymanski, S.J., Cabral, J., Chuine, I., Graham, C.H.,

- 635 Hartig, F., Kearney, M., Morin, X., Römermann, C., Schröder, B. &  
636 Singer, A. (2012) Correlation and process in species distribution models:  
637 bridging a dichotomy. *Journal of Biogeography*, **39**, 2119–2131.
- 638 Elith, J., Burgman, M.a. & Regan, H.M. (2002) Mapping epistemic uncer-  
639 tainties and vague concepts in predictions of species distribution. *Ecolog-  
640 ical Modelling*, **157**, 313–329.
- 641 Elith, J., Graham, C.H., Anderson, R.P., Dudík, M., Ferrier, S., Guisan,  
642 A., Hijmans, R.J., Huettmann, F., Leathwick, J.R., Lehmann, A., Li, J.,  
643 Lohmann, L.G., Loiselle, B.A., Manion, G., Moritz, C., Nakamura, M.,  
644 Nakazawa, Y., Overton, J.M.M., Townsend Peterson, A., Phillips, S.J.,  
645 Richardson, K., Scachetti-Pereira, R., Schapire, Robert, E., Soberón, J.,  
646 Williams, S., Wisz, M.S. & Zimmermann, N.E. (2006) Novel methods im-  
647 prove prediction of species distributions from occurrence data. *Ecography*,  
648 **29**, 129–151.
- 649 Elith, J., Kearney, M. & Phillips, S.J. (2010a) The art of modelling range-  
650 shifting species. *Methods in Ecology and Evolution*, **1**, 330–342.
- 651 Elith, J. & Leathwick, J.R. (2009) Species Distribution Models: Ecological  
652 Explanation and Prediction Across Space and Time. *Annual Review of  
653 Ecology, Evolution, and Systematics*, **40**, 677–697.

- 654 Elith, J., Leathwick, J.R., Hastie, T. & R. Leathwick, J. (2008a) A working  
655 guide to boosted regression trees. *Journal of Animal Ecology*, **77**, 802–813.
- 656 Elith, J., Leathwick, J.R., Hastie, T. & R. Leathwick, J. (2008b) Elith,  
657 Leathwick & Hastie A working guide to boosted regression trees - Online  
658 Appendices Page 1. *Journal of Animal Ecology*, **77**, 802–13.
- 659 Elith, J., Phillips, S.J., Hastie, T., Dudík, M., Chee, Y.E. & Yates, C.J.  
660 (2010b) A statistical explanation of MaxEnt for ecologists. *Diversity and*  
661 *Distributions*, **17**, 43–57.
- 662 Evgeniou, T., Micchelli, C. & Pontil, M. (2005) Learning multiple tasks with  
663 kernel methods. *Journal of Machine Learning Research*, **6**, 615–637.
- 664 Ferrier, S. (2002) Mapping spatial pattern in biodiversity for regional con-  
665 servation planning: where to from here? *Systematic Biology*, **51**, 331–363.
- 666 Ferrier, S. & Guisan, A. (2006) Spatial modelling of biodiversity at the  
667 community level. *Journal of Applied Ecology*, **43**, 393–404.
- 668 Ferrier, S., Watson, G., Pearce, J.L., Drielsma, M. & Manion, G. (2002) Ex-  
669 tended statistical approaches to modelling spatial pattern in biodiversity  
670 in northeast New South Wales. I. Species-level modelling. *Biodiversity &*  
671 *Conservation*, **11**, 2275–2307.

- 672 Frost, C. & Thompson, S.G. (2000) Correcting for regression dilution bias:  
 673 comparison of methods for a single predictor variable. *Journal of the*  
 674 *Royal Statistical Society: Series A (Statistics in Society)*, **163**, 173–189.
- 675 Gething, P.W., Van Boeckel, T.P., Smith, D.L., Guerra, C.a., Patil, A.P.,  
 676 Snow, R.W. & Hay, S.I. (2011) Modelling the global constraints of temper-  
 677 ature on transmission of *Plasmodium falciparum* and *P. vivax*. *Parasites*  
 678 *& Vectors*, **4**, 92.
- 679 Golding, N. (2013) *GRaF: Species distribution modelling using latent Gaus-*  
 680 *sian random fields*. R package version 0.1-0.
- 681 Guisan, A. & Zimmermann, N.E. (2000) Predictive habitat distribution  
 682 models in ecology. *Ecological modelling*, **135**, 147–186.
- 683 Hastie, T. (2011) *gam: Generalized Additive Models*. R package version  
 684 1.06.2.
- 685 Hastie, T. & Tibshirani, R. (1986) Generalized additive models. *Statistical*  
 686 *Science*, **1**, 297–318.
- 687 Hijmans, R.J., Phillips, S., Leathwick, J. & Elith, J. (2012) *dismo: Species*  
 688 *distribution modeling*. R package version 0.7-17.
- 689 Kissling, W.D., Dormann, C.F., Groeneveld, J., Hickler, T., Kühn, I., McIn-



690 erny, G.J., Montoya, J.M., Römermann, C., Schiffrs, K., Schurr, F.M.,  
 691 Singer, A., Svenning, J.C., Zimmermann, N.E. & OHara, R.B. (2011) To-  
 692 wards novel approaches to modelling biotic interactions in multispecies  
 693 assemblages at large spatial extents. *Journal of Biogeography*, **39**, 2163–  
 694 2178.

695 Lawson, C.R., Hodgson, J.A., Wilson, R.J. & Richards, S.A. (2013) Preva-  
 696 lence, thresholds and the performance of presence-absence models. *Meth-  
 697 ods in Ecology and Evolution*, **4**.

698 Lehmann, A., Leathwick, J. & Overton, J. (2002) Assessing New Zealand  
 699 fern diversity from spatial predictions of species assemblages. *Biodiversity  
 700 & Conservation*, pp. 2217–2238.

701 McCarthy, M. (2007) *Bayesian Methods for Ecology*. Cambridge University  
 702 Press, Cambridge.

703 McNerny, G.J., Purves, D.W. & McIntyre, K.M. (2011) Fine-scale environ-  
 704 mental variation in species distribution modelling : regression dilution ,  
 705 latent variables and neighbourly advice. *Methods in Ecology and Evolu-  
 706 tion*, **2**, 248–257.

707 McNamara, J.M. & Harding, K.C. (2004) Measurement error and estimates  
 708 of population extinction risk. *Ecology Letters*, **7**, 16–20.

- 709 McPherson, J.M., Jetz, W. & Rogers, D.J. (2004) The effects of species'  
710 range sizes on the accuracy of distribution models: ecological phenomenon  
711 or statistical artefact? *Journal of Applied Ecology*, pp. 811–823.
- 712 Murray, J.V., Goldizen, A.W., OLeary, R.A., McAlpine, C.A., Possingham,  
713 H.P. & Choy, S.L. (2009) How useful is expert opinion for predicting the  
714 distribution of a species within and beyond the region of expertise? A case  
715 study using brush-tailed rock-wallabies *Petrogale penicillata*. *Journal of*  
716 *Applied Ecology*, **46**, 842–851.
- 717 Newbold, T. (2010) Applications and limitations of museum data for con-  
718 servation and ecology, with particular attention to species distribution  
719 models. *Progress in Physical Geography*, **34**, 3–22.
- 720 Patil, A. (2007) *Bayesian Nonparametrics for Inference of Ecological Dy-*  
721 *namics*. Ph.D. thesis, University of California Santa Cruz.
- 722 Pearson, R.G., Raxworthy, C.J., Nakamura, M. & Townsend Peterson, A.  
723 (2006) Predicting species distributions from small numbers of occurrence  
724 records: a test case using cryptic geckos in Madagascar. *Journal of Bio-*  
725 *geography*, **34**, 102–117.
- 726 Phillips, S.J., Anderson, R.P. & Schapire, Robert, E. (2006) Maximum en-

727     tropy modeling of species geographic distributions. *Ecological Modelling*,  
728     **190**, 231–259.

729     Phillips, S.J., Dudík, M., Elith, J., Graham, C.H., Lehmann, A., Leathwick,  
730     J.R. & Ferrier, S. (2009) Sample selection bias and presence-only dis-  
731     tribution models: implications for background and pseudo-absence data.  
732     *Ecological Applications*, **19**, 181–97.

733     Phillips, S.J. & Elith, J. (2013) On Estimating Probability of Presence from  
734     Use-Availability or Presence-Background Data. *Ecology*.

735     Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D. & R Core Team (2012) *nlme*:  
736     *Linear and Nonlinear Mixed Effects Models*. R package version 3.1-106.

737     Preston, C., Pearman, D. & Dines, T. (2002) *New atlas of the British and*  
738     *Irish flora: an atlas of the vascular plants of Britain, Ireland, the Isle of*  
739     *Man and the Channel Islands*. Oxford University Press, Oxford.

740     R Development Core Team (2012) *R: A Language and Environment for*  
741     *Statistical Computing*. R Foundation for Statistical Computing, Vienna,  
742     Austria. ISBN 3-900051-07-0.

743     Rasmussen, C. & Williams, C. (2006) *Gaussian processes for machine learn-*  
744     *ing*. MIT Press.

- 745 Ridgeway, G. (2013) *gbm: Generalized Boosted Regression Models*. R pack-  
746 age version 2.1.
- 747 Robin, X., Turck, N., Hainard, A., Tiberti, N., Lisacek, F., Sanchez, J.C.  
748 & Mller, M. (2011) proc: an open-source package for r and s+ to analyze  
749 and compare roc curves. *BMC Bioinformatics*, **12**, 77.
- 750 Rue, H., Martino, S. & Chopin, N. (2009) Approximate Bayesian inference  
751 for latent Gaussian models by using integrated nested Laplace approx-  
752 imations. *Journal of the Royal Statistical Society: Series B (Statistical*  
753 *Methodology)*, **71**, 319–392.
- 754 Schmidberger, M., Tierney, L. & Mansmann, U. (2009) State of the Art in  
755 Parallel Computing with R. *Journal of Statistical Software*, **31**.
- 756 Sigourney, D.B., Munch, S.B. & Letcher, B.H. (2012) Combining a Bayesian  
757 nonparametric method with a hierarchical framework to estimate individ-  
758 ual and temporal variation in growth. *Ecological Modelling*, **247**, 125–134.
- 759 Sinclair, S.J.S., White, M.M.D. & Newell, G.R. (2010) How useful are species  
760 distribution models for managing biodiversity under future climates. *Ecol-*  
761 *ogy and Society*, **15**.
- 762 Sinka, M.E., Bangs, M.J., Manguin, S., Coetzee, M., Mbogo, C.M., Hem-

ingway, J., Patil, A.P., Temperley, W.H., Gething, P.W., Kabaria, C.W.,  
 Okara, R.M., Van Boeckel, T.P., Godfray, H.C.J., Harbach, R.E. & Hay,  
 S.I. (2010) The dominant Anopheles vectors of human malaria in Africa,  
 Europe and the Middle East: occurrence data, distribution maps and  
 bionomic precis. *Parasites & Vectors*, **3**, 117.

Spiegelhalter, D.J., Best, N.G., Carlin, B.P. & van der Linde, A. (2002)  
 Bayesian measures of model complexity and fit. *Journal of the Royal  
 Statistical Society: Series B (Statistical Methodology)*, **64**, 583–639.

Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *Jour-  
 nal of the Royal Statistical Society: Series B (Statistical Methodology)*, **58**,  
 267–288.

Vanhatalo, J., Pietiläinen, V. & Vehtari, A. (2010) Approximate infer-  
 ence for disease mapping with sparse Gaussian processes. *Statistics in  
 Medicine*, **29**, 1580–607.

Vanhatalo, J., Veneranta, L. & Hudd, R. (2012) Species distribution mod-  
 eling with Gaussian processes: A case study with the youngest stages  
 of sea spawning whitefish (*Coregonus lavaretus* L. s.l.) larvae. *Ecological  
 Modelling*, **228**, 49–58.

- 781 Ward, G., Hastie, T., Barry, S., Elith, J. & Leathwick, J.R. (2009) Presence-  
782 only data and the em algorithm. *Biometrics*, **65**, 554–63.
- 783 Wenger, S.J. & Olden, J.D. (2012) Assessing transferability of ecological  
784 models: an underappreciated aspect of statistical validation. *Methods in*  
785 *Ecology and Evolution*, **3**, 260–267.
- 786 Wisz, M.S., Pottier, J., Kissling, W.D., Pellissier, L., Lenoir, J., Damgaard,  
787 C.F., Dormann, C.F., Forchhammer, M.C., Grytnes, J.A., Guisan, A.,  
788 Heikkinen, R.K., Høye, T.T., Kühn, I., Luoto, M., Maiorano, L., Nils-  
789 son, M.C., Normand, S., Ockinger, E., Schmidt, N.M., Termansen, M.,  
790 Timmermann, A., Wardle, D.A., Aastrup, P. & Svenning, J.C. (2013) The  
791 role of biotic interactions in shaping distributions and realised assemblages  
792 of species: implications for species distribution modelling. *Biological Re-*  
793 *views of the Cambridge Philosophical Society*, **88**, 15–30.
- 794 Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal  
795 likelihood estimation of semiparametric generalized linear models. *Journal*  
796 *of the Royal Statistical Society (B)*, **73**, 3–36.
- 797 Yackulic, C.B., Chandler, R.B., Zipkin, E.F., Royle, J.A., Nichols, J.D.,  
798 Campbell Grant, E.H. & Veran, S. (2012) Presence-only modelling using

799     MAXENT: when can we trust the inferences? *Methods in Ecology and*  
800     *Evolution*, **4**, 236–243.

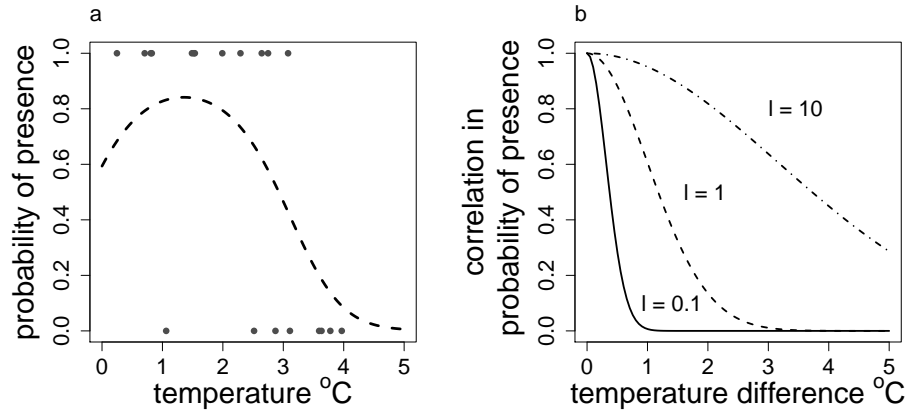


Fig. 1: Illustration of the covariance function using synthetic data: (a) observed presence-absence data (points) and the true underlying probability of presence as a function of temperature (dashed line); (b) correlation between probability of presence at different sites, calculated from temperature difference between these sites using the covariance function with three different lengthscale parameters (discussed in the text). Models fitted to the data using these three lengthscales are illustrated in Fig. 2.



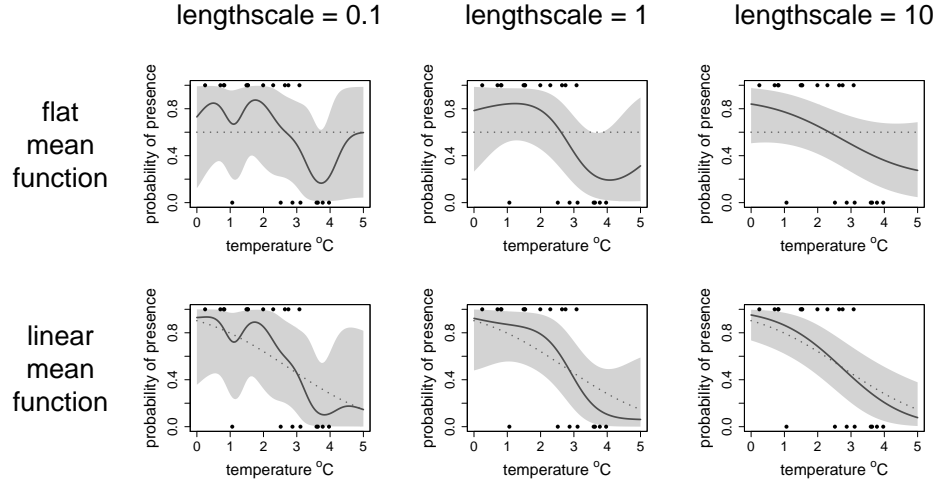


Fig. 2: Effects of the mean function and lengthscale on the fitted GP model.

Shown are the observed data (points), the value of the mean function (dotted line), the probability of presence predicted by the GP model (solid line) and associated 95% credible intervals for this prediction (shaded grey area). Models are fitted with either the default flat mean function at the mean probability of presence (upper row) or a mean function representing some prior knowledge about how probability of presence relates to temperature, as described in the text (lower row).

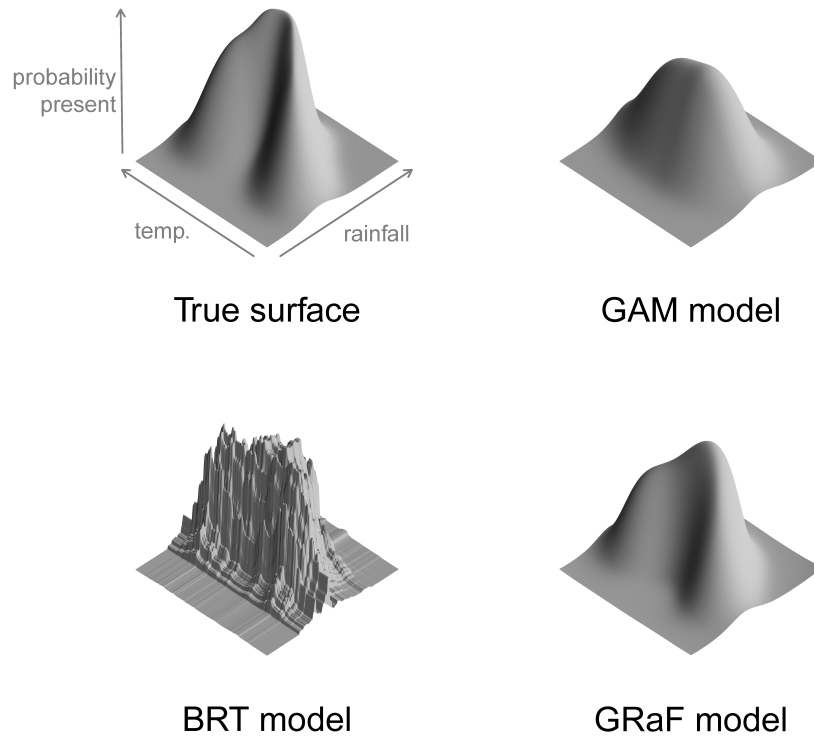


Fig. 3: Predictive surfaces fitted by boosted regression trees (BRT), a generalized additive model with univariate smoothers (GAM) and a GP model to simulated data with a strong non-linear interaction. The true surface represents the probability of presence of a hypothetical species in response to temperature and rainfall. Models were fitted to 1000 random presence/absence observations drawn from the true probability surface (a mixture of Gaussians).

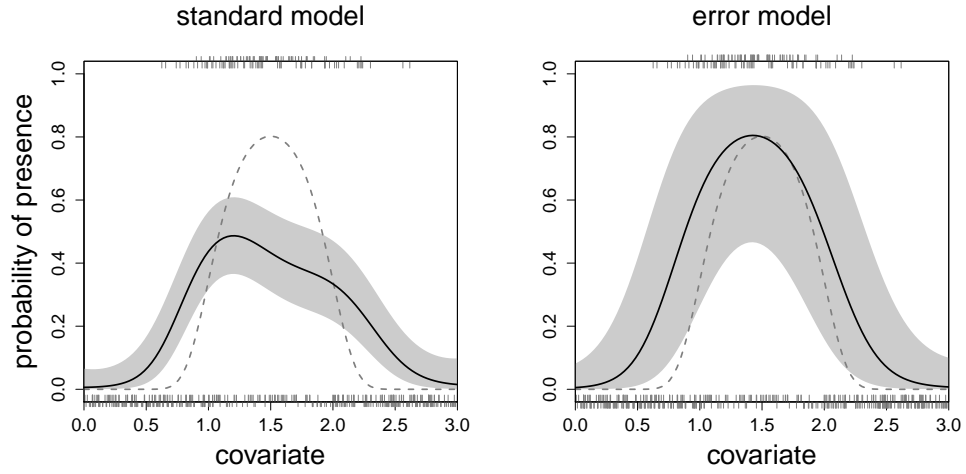


Fig. 4: Comparison of predictive surfaces fitted by GP models to simulated data with the covariate measured under error either ignoring measurement error (standard model) or accounting for it (error model). Solid lines give predictions and shaded regions represent 95% credible intervals. Three hundred presence-absence points were generated from the true model (dashed line) given the correct value of the covariate (tick marks outside box, presence on top line and absence on bottom line). Normally distributed random noise was then added to simulate covariates measured under error (tick marks inside box). Both models were then fitted to these data with measurement error. The error model was provided with the correct standard deviation of the error (0.5) whilst the standard model ignored the error.

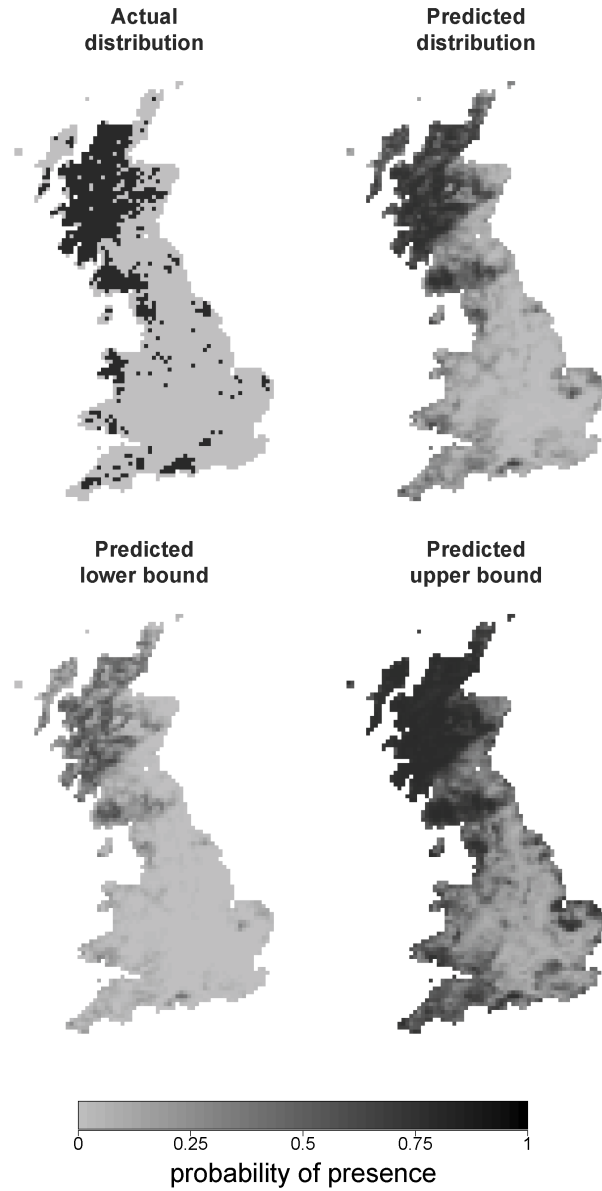


Fig. 5: True and predicted distributions of Bog myrtle (*Myrica gale*; prevalence 0.27) in Great Britain. Shown are the true distribution, the predicted distribution from a GP model (the Maximum *a posteriori* prediction), and lower and upper bounds on this prediction, representing our uncertainty in it (95% credible intervals, automatically generated by the GP). The GP model was fitted to 300 presence-absence data points - around 10% of the dataset.

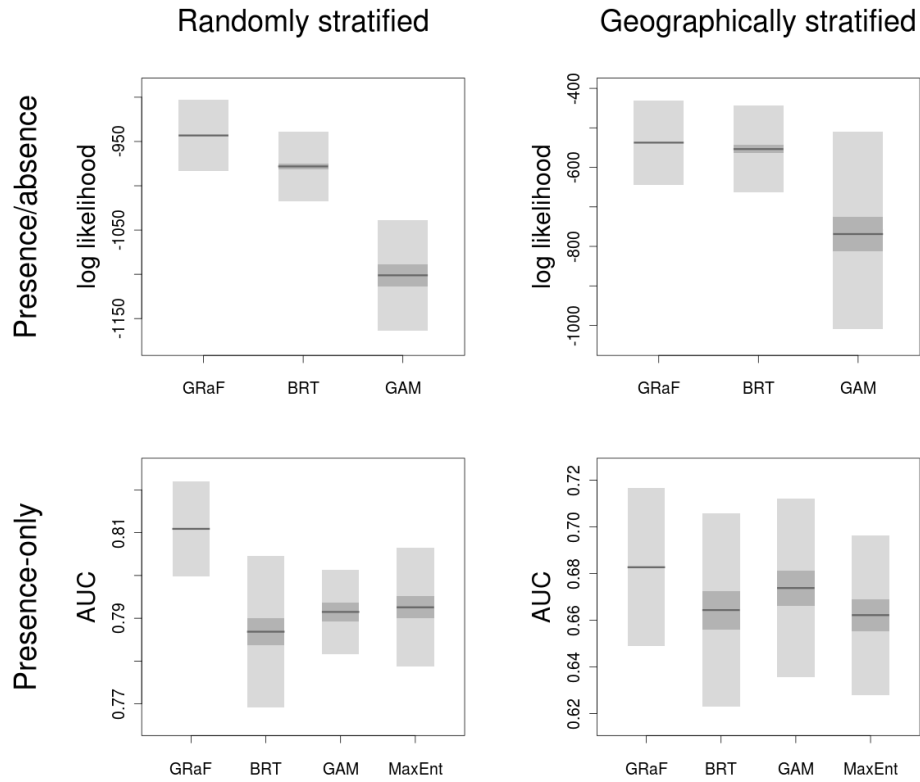


Fig. 6: Marginal validation statistics for model predictions to withheld training sets for presence/absence and presence-only data with two types of cross-validation. Centre lines give the means of the marginal validation statistic and light grey boxes give  $\pm 1$  standard deviation of the marginal statistics, as an indication of the likely differences in performance of each model on an ‘average’ species in the plant dataset. Dark grey boxes give  $\pm 1$  standard deviation of the estimated difference in the mean of the statistic for each model from the statistic for GRaF, as a visual representation of the statistical tests carried out. As GRaF was used as the contrast for these tests, no dark grey box is presented for this model. Higher log-likelihoods and higher AUCs indicate more accurate predictions to the evaluation set.

## 802 **Supporting Information**

803 **Appendix S1.** Statistical explanation and specification

804 **Appendix S2.** Assessment of approximation error

805 **Appendix S3.** Explanation of data for SDM comparison

806 **Appendix S4.** Demonstration of GRaF R package

807 **Appendix S5.** R code and data to compare SDMs and reproduce figures