

The influence of spatial errors in species occurrence data used in distribution models

Catherine H. Graham^{1*}, Jane Elith², Robert J. Hijmans^{3†}, Antoine Guisan⁴,
A. Townsend Peterson⁵, Bette A. Loiselle⁶ and The Nceas Predicting Species
Distributions Working Group‡

¹Department of Ecology and Evolution, Stony Brook University, New York, 11794, USA; ²School of Botany, University of Melbourne, Parkville, Victoria, 3010, Australia; ³Museum of Vertebrate Zoology, University of California, 3101 Valley Life Sciences Building, Berkeley, CA, USA; ⁴Department of Ecology and Evolution, University of Lausanne, CH-1015 Lausanne, Switzerland; ⁵Natural History Museum and Biodiversity Research Center and Department of Ecology and Evolutionary Biology, University of Kansas, Lawrence, KA, USA; and ⁶Department of Biology, University of Missouri, St Louis, USA

Summary

1. Species distribution modelling is used increasingly in both applied and theoretical research to predict how species are distributed and to understand attributes of species' environmental requirements. In species distribution modelling, various statistical methods are used that combine species occurrence data with environmental spatial data layers to predict the suitability of any site for that species. While the number of data sharing initiatives involving species' occurrences in the scientific community has increased dramatically over the past few years, various data quality and methodological concerns related to using these data for species distribution modelling have not been addressed adequately.

2. We evaluated how uncertainty in georeferences and associated locational error in occurrences influence species distribution modelling using two treatments: (1) a control treatment where models were calibrated with original, accurate data and (2) an error treatment where data were first degraded spatially to simulate locational error. To incorporate error into the coordinates, we moved each coordinate with a random number drawn from the normal distribution with a mean of zero and a standard deviation of 5 km. We evaluated the influence of error on the performance of 10 commonly used distributional modelling techniques applied to 40 species in four distinct geographical regions.

3. Locational error in occurrences reduced model performance in three of these regions; relatively accurate predictions of species distributions were possible for most species, even with degraded occurrences. Two species distribution modelling techniques, boosted regression trees and maximum entropy, were the best performing models in the face of locational errors. The results obtained with boosted regression trees were only slightly degraded by errors in location, and the results obtained with the maximum entropy approach were not affected by such errors.

4. *Synthesis and applications.* To use the vast array of occurrence data that exists currently for research and management relating to the geographical ranges of species, modellers need to know the influence of locational error on model quality and whether some modelling techniques are particularly robust to error. We show that certain modelling techniques are particularly robust to a moderate level of locational error and that useful predictions of species distributions can be made even when occurrence data include some error.

Key-words: error, geo-referencing, locality points, predictive modelling algorithms, species distribution model, uncertainty

*Correspondence author. E-mail: cgraham@life.bio.sunysb.edu

†Present address: International Rice Research Institute, Los Baños, Laguna, Philippines.

‡The Working Group participants are listed at the end of the paper.

Introduction

Delineating the range of a species is important for various applications in applied and theoretical biology. Species distribution modelling is used frequently to predict the full geographical range of species, particularly when an estimate of the probability of occurrence or the relative suitability at a given site is required (Guisan & Thuiller 2005). For example, species distribution modelling has been used to manage species of conservation concern (Gaubert, Papes & Peterson 2006), create richness maps for conservation planning (Loiselle *et al.* 2003; Rissler *et al.* 2006) and predict the geographical spread of invasive species (Peterson 2003). Species distribution modelling requires species occurrence data and environmental spatial data layers, which are combined to create a predictive model describing the suitability of any site for the species. Species occurrence data are increasingly available as a result of data generation, improvement and sharing initiatives in the natural history museum, broader scientific and conservation communities (Graham *et al.* 2004; Suarez & Tsutsui 2004). While species occurrence data provide the basis for much research, their use can be problematic because absence data are generally unavailable, sample sizes are often small and geographical bias and spatial error in the data are generally unexplored or unknown (Hijmans *et al.* 2000; Graham *et al.* 2004; Wieczorek, Guo & Hijmans 2004; Rowe 2005). Elith *et al.* (2006) showed that accurate models can be made with presence-only data and a series of studies have explored how the number of occurrences and bias in data influences model performance (i.e. McPherson, Jetz & Rogers 2004; Barry & Elith 2006; Hernandez *et al.* 2006). Here, we conduct experimental manipulations to evaluate how error influences the accuracy of distribution models.

Errors in occurrence data are caused by a variety of factors, including mistakes in transferral of data from field sheets to electronic databases, rounding errors, failure to specify the geographical datum (the size and shape of the earth and the origin and orientation of the coordinate systems used to map the earth) used to measure geographical location and retrospective georeferencing of imprecise locality descriptions (Barry & Elith 2006; Wieczorek, Guo & Hijmans 2004). Error can be introduced into retrospective georeferencing of textual descriptions of the location where the species was observed. Localities are sometimes described as named places (e.g. 'Berkeley, California') that can have a significant geographical extent. When an offset is used (e.g. '10 miles N of'), the uncertainty caused by the extent is compounded by the uncertainty in the distance (the precision of 10), direction (exactly north?) and the path (straight-line distance, or along a road?) (Wieczorek, Guo & Hijmans 2004). Also, collectors accumulate specimens from a sometimes-broad radius around a field camp, but often use a single locality descriptor for that collecting site (Anderson, Gomez-Laverde & Peterson 2002). Despite these sources of uncertainty, locations within electronic databases are typically georeferenced as a single point. This is because alternatives, such as describing a location by an area encompassing all possible localities, are usually too cumbersome

and do not always facilitate analytical use of the data in species distribution modelling and other applications.

Techniques have been developed to detect, quantify and document uncertainty in occurrence data (Wieczorek, Guo & Hijmans 2004). Mistakes in transferring information from data sheets often results in completely incorrect georeferences (e.g. using the wrong locality with a particular species), but fortunately some of these errors can be found in data cleaning (Hijmans *et al.* 1999). Some current georeferencing initiatives provide assessments of potential error or geographical footprints associated with a given locality (Guralnick *et al.* 2006). Using this information, highly uncertain locations can be removed from the analysis but this would reduce sample size, which could influence model performance negatively (McPherson, Jetz & Rogers 2004; Hernandez *et al.* 2006).

In order to use museum data in modelling studies and to take full advantage of uncertainty estimates when they accompany georeferences, researchers need to know the extent to which error in occurrence data influences model performance. Therefore, we designed an experiment in which we compared model output obtained with the most accurate location data available, with model output based on deliberately degraded location data. Degraded location data were derived from the original data by adding random error. We conducted this experiment for 10 species in each of four geographical regions, a subset of the evaluation data described elsewhere (Elith *et al.* 2006). We compared a series of modelling techniques to determine whether there was differential performance on degraded occurrence data as part of a broader suite of comparisons and experimental manipulations to evaluate factors important in species distribution modelling. Techniques included those that have been used commonly in species distributional modelling for presence-only data, such as BIOCLIM (Busby 1991) and DOMAIN (Walker & Cocks 1991; Carpenter, Gillison & Winter 1993); regression-based techniques which use background (also referred to as pseudo-absence) information, such as generalized linear and additive models (Yee & Mitchell 1991; Pearce & Ferrier 2000; Austin 2002; Lehmann, Overton & Leathwick 2003; Wintle, Elith & Potts 2005) and multiple regression splines (Leathwick *et al.* 2005); and genetic algorithms (Stockwell & Peters 1999). We also used novel techniques for species distribution modelling, including maximum entropy (Phillips, Anderson & Schapire *et al.* 2006) and boosted regression trees (Friedman *et al.* 2000; Schapire 2003) that are now used commonly in machine-learning statistical research. These newer techniques performed particularly well in a broad comparative study (Elith *et al.* 2006).

Methods

DATA SOURCES

We used data from four regions of the world: the Australian Wet Tropics (AWT), north-east New South Wales, Australia (NSW), New Zealand (NZ) and Switzerland (SWI). For each region, we used 11–13 environmental spatial data layers that were typical of

Table 1. Summary of occurrence data for model building (training) and evaluation (testing) by region

	Training presences			Test presences			Test absences		
	Mean	Min	Max	Mean	Min	Max	Mean	Min	Max
AWT	28	22	45	105	81	169	184	126	237
NSW	61	18	137	231	69	517	939	407	1 989
NZ	461	45	1228	1735	171	4619	16 934	13 284	18 915
SWI	192	22	944	722	82	3553	9 632	6 049	10 442
All regions	186			698			6 922		

what is used generally for species distribution modelling. To describe spatial heterogeneity of each layer we computed the minimum, maximum, mean and standard deviation (see Appendix S1 in Supplementary material). We also quantified spatial heterogeneity by calculating the mean correlation between each pixel and a pixel approximately 5 km away in each of eight directions (N, NE, E, SE, S, SW, W, NW). In all regions, the spatial resolution of the environmental layers was about 100 m.

We used data for 10 species from each of the four regions including birds, frogs and reptiles of the AWT; plants, birds, mammals and reptiles in NSW; and plants of NZ and SWI (see Supplementary material, Appendix S2). We chose species to represent a range of geographical distributions, habitat specialization and (where possible) biological groups/life forms. Each species had > 80 occurrence (presence) records before experimental manipulation of the data, as well as a considerable number of absence records (Table 1). All presence records used in modelling and presence/absence records used in evaluation were obtained from planned surveys that used global positioning system (GPS) receivers to record locations. Therefore, we are confident that the georeferences of these records are very accurate (error mainly < 30 m). For more information on layers and data please see Elith *et al.* (2006).

EXPERIMENTAL MANIPULATION OF DATA

The aim of our experiment was to determine the effect of georeferencing error on the predictive performance of species distribution models in each of the four regions described above. Our focus was to test how error influenced modelling techniques when only occurrence data (i.e. presence-only) were used. Because the true museum data in these regions have variable and sometimes unknown accuracy (Elith *et al.* 2006), they are not a reliable test of the effect of error. Therefore, we derived presence-only samples from accurate presence-absence data. First, for each species we chose randomly four-fifths of the presence records to withhold for evaluation, together with all the absence records. Note that none of the 'observed' absences was used for modelling. With the remaining presences we generated two data sets for fitting the models: one where the locations were degraded to simulate error (error treatment) and a control where the data were not manipulated. This resulted in a mean of 186 (range 22–1228) occurrences for modelling and a robust evaluation data set (Table 1). While this split sample approach does not provide a truly independent test (Araujo *et al.* 2005), it remains a reasonable approach for model evaluation when such data are not available. For the error treatment, x and y coordinates were shifted in a random direction by an amount sampled from a normal distribution with a mean of zero and a standard deviation of 5 km. Our assumption is that this is representative of the error associated with museum data, although error in museum data will vary. If shifting the sites resulted in them moving

outside the study region we recalculated the shift (as above) until the point remained within the bounds of the region.

MODELLING TECHNIQUES

We describe briefly the 10 modelling techniques used in our comparison; detailed descriptions can be found in Elith *et al.* (2006). Three modelling techniques (BIOCLIM, DOMAIN and LIVES) use only presence data in their formulation. All the rest require absence or background data. We used 10 000 random samples from the study region (sometimes referred to as 'pseudo-absences') as absence data for these techniques.

1. BIOCLIM is an envelope-style technique that uses species presence records to create a hyper-space which summarize how these records are distributed with respect to environmental variables (Busby 1991). This envelope specifies the model in terms of percentiles or upper and lower tolerances. BIOCLIM was implemented within DIVA-GIS (Elith *et al.* 2006), a different implementation of the original method reported by Busby (1991).

2. DOMAIN is a distance-based technique that estimates the environmental similarity, using the Gower distance metric, between a site of interest and the nearest presence record in environmental space (Carpenter, Gillison & Winter 1993). It uses presence-only data and the resulting predictions range between 0 and 100.

3. LIVES uses a limiting factor method that postulates that the occurrence of a species is determined only by the environmental factor that most limits its distribution. The limiting factor of the species is defined as the environmental factor that has the minimum similarity (based on a modified Gower metric) among the environmental factors considered in the mode. The suitability index of a grid cell to the species is the maximum of similarities of the target grid to each of the training grids (Carpenter, Gillison & Winter 1993; Li & Hilbert, in press).

4. Generalized linear models (GLMs) used the occurrence data for each species and 10 000 random samples from the background as the dependent variable and environmental variables as independent variables. The 10 000 background samples were weighted so that the total weight for presence was equal to the total weight for absence (also performed for generalized additive models (GAMs), below). We fitted parametric terms, usually some combination of linear, quadratic and/or cubic terms, using a stepwise procedure where successively simpler fits were evaluated with Akaike's information criterion (AIC).

5. Implementation of GAMs was similar to GLMs but this method uses non-parametric, data-defined smoothers to fit non-linear functions. Because of their greater flexibility, GAMs are considered more capable of modelling complex ecological response shapes than GLMs (Yee & Mitchell 1991; Guisan, Edwards & Hastie 2002; Wintle, Elith & Potts 2005).

6. Multivariate adaptive regression splines (MARS) provide an alternative regression-based technique for fitting non-linear

responses. MARS uses piecewise linear fits rather than smooth functions and a fitting procedure that makes them much faster to implement than GAMs (Leathwick *et al.* 2005). Our implementations of the regression techniques presented here did not attempt to model interactions (but see Elith *et al.* 2006).

7. GARP uses a genetic algorithm to select a set of rules (e.g. adaptations of regression and range specifications) that predicts most accurately the species distribution (Stockwell & Peters 1999). We used two versions of GARP: the desktop version that has been used widely for modelling data from natural history collections (DK-GARP) and an improved version of GARP that has not yet been released (OM-GARP). This new version has updated algorithms for developing rule sets. DK-GARP was not run on the NZ data set because the data set was too large; therefore DK-GARP was excluded from some analyses.

8. Maximum entropy models (MAXENT) have been developed within the machine community and use presence and background absence data. The approach is to model a species' distribution as an unknown probability distribution over the set of background points (i.e. the values at the background points are non-negative and sum to 1). The environmental variable values at the presence localities impose some constraints on the unknown distribution. For example, the mean and variance of the environmental variables should be close to their empirical values. The maximum entropy approach then approximates the unknown distribution using the distribution over the background points that maximizes entropy, subject to these constraints. Because entropy is a convex function it has no local maxima, so the unique maximum (the maximum entropy distribution) can be found by a variety of optimization methods. We implemented MAXENT so that it fitted interactions and refer to it as MAXENT-T to be consistent with Elith *et al.* (2006).

9. Boosted regression trees (BRT) combines two algorithms: the boosting algorithm iteratively calls the regression-tree algorithm to construct a combination or 'ensemble' of trees. The regression trees are fitted sequentially, and use a gradient descent algorithm to model iteratively the residuals that reflect the lack of fit from the previous set of trees. Rather than using the stepwise processes in GLMs and GAMs that are based on model comparisons using AICs, our version of BRT uses cross-validation to compare model fit with predictive ability (Schapire 2003; Elith *et al.* 2006).

EVALUATION

The modelled distributions were evaluated for predictive performance using the partitioned data described above. We used the area under the receiver operating characteristic (ROC) curve to assess the agreement between the presence-absence records and the model predictions (Fielding & Bell 1997). A ROC curve plot is created by plotting the sensitivity values, the true-positive fraction, against 1 – specificity, the false-positive fraction, for all available probability thresholds. Model performance is determined by calculating the area under this curve (AUC) such that a curve that maximizes sensitivity for low values of the false-positive fraction is considered a good model (Hanley & McNeil 1982; Fielding & Bell 1997). Hence, AUC is a threshold independent measure of the ability of a model to discriminate between sites where a species is present, vs those where it is absent (Hanley & McNeil 1982). AUC ranges from 0 to 1, where a score of 1 indicates perfect discrimination, a score of 0.5 implies random predictive discrimination and values less than 0.5 indicate performance worse than random. We did not use other evaluation statistics, such as correlation or maximum kappa, because these gave qualitatively similar results in a previous large

model comparison study that used the same data (Elith *et al.* 2006). We emphasize that the approaches to model evaluation used here balance correct predictions at both presence and absence sites, and as such focus on current prediction of distributions rather than characterization of ecological niches as described by some authors (e.g. Peterson 2007).

DATA ANALYSES

We report AUC values for predictions made with accurate location data and for predictions made with degraded location data. This allowed us to assess the predictive performance across techniques and to determine if predictive performance decreases when error is added. We used Wilcoxon's signed-rank tests for matched pairs to determine if addition of error influenced model performance; that is, we compared AUC of models built with original vs. degraded data. One analysis was conducted for each of the four regions based on the mean model performance of the nine to 10 modelling techniques (DK-GARP was not run on NZ) for each of the 10 species modelled in each region. A second set of analyses was conducted for each technique using data from each species and region ($n = 40$). We applied Bonferroni corrections because multiple comparisons were conducted.

To test whether techniques were sensitive to error to differing degrees we used ranked data because they provide the position of a given technique relative to other techniques for a given modelling challenge. This allowed us to determine if rank changed with error, i.e. if some techniques were more sensitive to error than others. Ranks are preferable to actual values because they are not sensitive to absolute values of AUC, which is important because interpretation of relative changes in AUC is somewhat problematic. Specifically, when a model has a very low AUC (i.e. close to 0.5 or random), adding error to localities would be less likely to influence model performance than adding error to a highly accurate model. Further, an incremental change in model performance of 0.1 for random predictions (e.g. AUC of 0.5–0.6) is not equivalent to the same change of 0.1 for a high-performing model. We therefore compared model performance (ranked performance relative to other models) using ranked AUC scores across all technique types using a Friedman's analysis of variance (ANOVA) for multiple dependent samples. We ran separate tests for original and degraded data. We followed these tests with a series of *post-hoc* Wilcoxon's paired tests to determine differences in model performance. We did not run all possible combinations of tests, but chose only those that would allow us to determine significant differences among techniques (i.e. if technique A was better than B, and B better than C, we did not compare A to C).

Results

INFLUENCE OF ERROR ON MODEL PERFORMANCE ACROSS REGIONS

Model predictions with the original data (control treatment) tended to have higher mean AUC scores than those obtained with the degraded data (error treatment; Table 2), although in all regions there were some instances in which model performance was higher when using degraded points (Fig. 1). In NZ and SWI, model performance was significantly worse overall with degraded data (Wilcoxon's paired test, $Z = 2.70$, $n = 10$, $P = 0.007$; $Z = 2.49$, $n = 10$, $P = 0.012$, respectively), while in NSW, performance using original data was only

Table 2. Mean and standard deviation (SD) of area under the curve (AUC) scores for models built on control (no error added) data calculated across all species and regions for each modelling method. Diff refers to the mean absolute difference between AUC scores of models built with control and degraded data. *Z*-values and *P*-values are associated with Wilcoxon's paired test; *indicates that result is significant at 0.05 after Bonferroni correction for multiple comparisons

Modelling method	Mean	SD	Difference	<i>Z</i> -value	<i>P</i> -value
BIOCLIM	0.66	0.09	0.056	3.04	0.0023*
BRT	0.76	0.10	0.036	3.47	0.0005*
DOMAIN	0.74	0.10	0.073	4.93	< 0.0001*
LIVES	0.70	0.09	0.056	4.65	< 0.0001*
GAM	0.73	0.10	0.053	2.18	0.0290
GARP-DK	0.67	0.11	0.055	1.88	0.0598
GARP-OM	0.71	0.11	0.060	2.97	0.003*
GLM	0.73	0.11	0.047	2.12	0.0249
MARS	0.71	0.11	0.059	2.00	0.0452
MAXENT	0.73	0.11	0.049	1.88	0.0600

marginally better than degraded data ($Z = 1.78$, $n = 10$, $P = 0.074$). In contrast, in AWT, model performance with the error treatment was higher than with the control treatment, although not significantly ($Z = 0.36$, $n = 10$, $P = 0.72$). The relatively minor effect of location error on the predictive performance of models, particularly in AWT, may have been due to the somewhat low average performance of control models. If model performance was not much better than random (note AUC values around and below 0.5 for AWT in Fig. 1), then decline in model performance with degraded data cannot be expected. Despite these results, mean AUC scores remained relatively high overall for predictions made with degraded data (Table 2); that is, even when differences were statistically significant, they were not particularly large.

INFLUENCE OF ERROR ON MODEL PERFORMANCE ACROSS MODELS

Not all model techniques were influenced equally by error. The quality of the predictions made with BIOCLIM, BRT, DOMAIN, LIVES and OM-GARP diminished significantly when degraded data were used (Table 2, Fig. 2). However, there was no evidence of decline in performance of regression-based techniques (GLM, GAM and MARS), MAXENT-T and DK-GARP (note that DK-GARP was run on only three of the four regions). This lack of a significant decline may be a result of several species which had higher AUC scores for models built with degraded data (Fig. 2).

To compare model performance across techniques we used ranked data (see Methods for justification). Model performance varied across the different techniques on both original and degraded data (Friedman's ANOVA, $n = 40$, d.f. = 8, $\chi^2 = 67.55$, $P < 0.0001$ and $\chi^2 = 77.59$, $P < 0.0001$, respectively). The rank of different methods varied slightly as a result of the error treatment; however, the best-performing and worst-performing models were consistent (Fig. 3). BRT had a significantly higher performance than other techniques

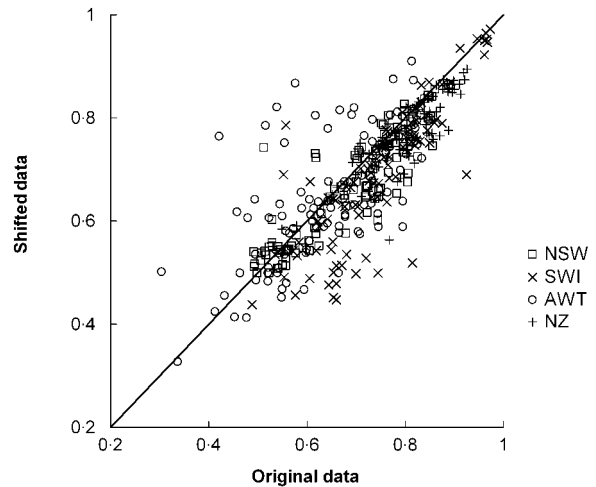


Fig. 1. Area under the curve values for model predictions made with the original data vs. predictions made with the degraded data. Line is $x = y$ (points below the line indicate worse performance with degraded locality data).

on the original data and a small but consistent decline in model performance with the error treatment (Table 1, Fig. 2). BRT was not different from MAXENT-T when models were built on degraded data (Fig. 3). A few MAXENT-T models had very low performance in AWT, especially on the original data (Fig. 2). These models were improved significantly when reconducted with the newest version of MAXENT (Steven Phillips, personal communication); however, determining what caused these low values led to improvement of the MAXENT model, so rerunning these models with a newer version would have biased our results (i.e. tuning the model based on evaluation with independent data). DOMAIN and BRT had the largest decline in performance with degraded data.

Discussion

Point occurrence data are available increasingly from museums, herbaria, surveys and numerous other sources. However, all these data include some locational or georeferencing error at some scale (Hijmans *et al.* 1999; Rowe 2005). By comparing models run with high-quality and with degraded data, we evaluated what influence error might have on the accuracy of species distribution models across a variety of taxa and regions. For many species, models run with data subject to random locational errors resulted in less accurate models, as predicted initially. None the less, this effect was variable and in at least two of the regions studied here, New South Wales and the Australian Wet Tropics, the error treatment had either a marginal effect or no consistent effect on model performance. Even in regions where the error treatment resulted in decreased model performance – Switzerland and New Zealand – AUC scores remained relatively high. Values of $AUC > 0.75$ are generally considered adequate for use in conservation planning and other applications (Pearce & Ferrier 2000) and values above 0.7 are considered reasonable

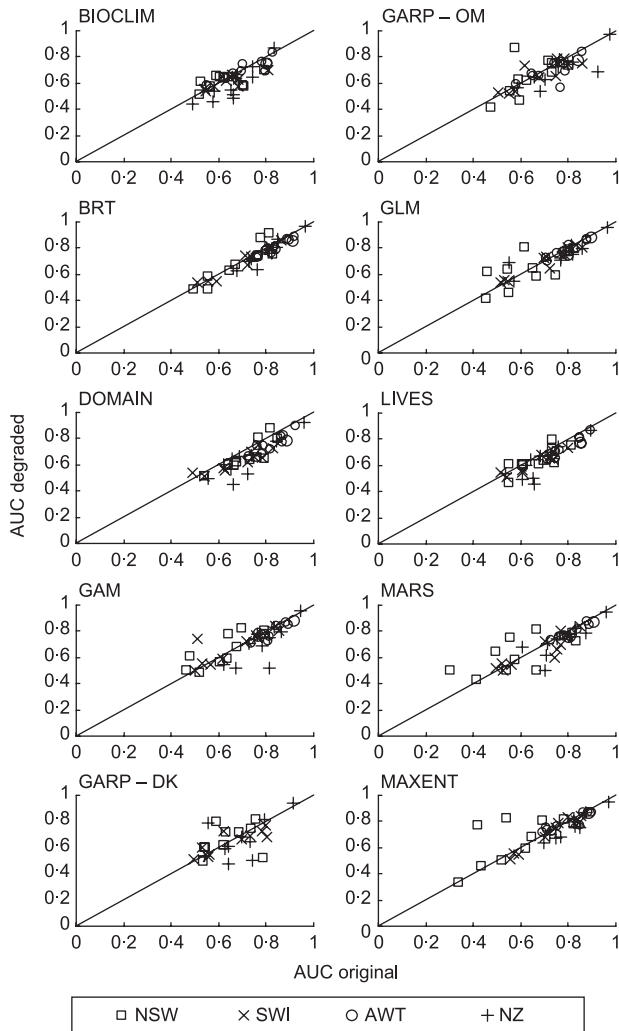


Fig. 2. Area under the curve values for model predictions made with the original data vs. predictions made with the degraded data for each of the 10 techniques. Line is $x = y$ (points below the line indicate worse performance with degraded locality data). Symbols on graphs indicate region where asterisk = New South Wales (NSW), circle = New Zealand (NZ), square = Australian Wet Tropics (AWT) and the plus sign = Switzerland (SWI).

by Swets (1988). Hence, our overall results are encouraging because, from a locational accuracy viewpoint, they suggest that modelling techniques may be robust to some error in data. None the less, these results should be viewed as preliminary, as we included only one error treatment and performance of models might change under different levels of error. Further, one might argue that if error does not influence model performance then predictions must be rather imprecise to begin with. However, the frequently high AUC scores do not support that view, and it appears that modelling approaches are robust to moderate levels of error in the data. The key here is not to over-interpret results of these models. They are useful at certain spatial scales, and should not be used at finer scales where they may be less useful.

Modelling techniques differed both in their overall performance and their sensitivity to error. BRT stood out as the

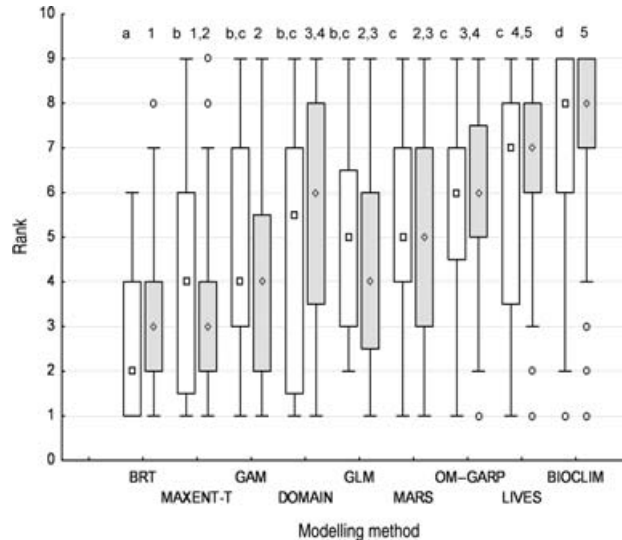


Fig. 3. Box plot of ranked performance (ranked relative to other techniques) of each method across all species/regions using original (clear bars) and degraded (hatched bars) data. The medians (central line), quartiles (box), and non-outlier range (whiskers) and outliers are shown in the plot. The letters, in the case of original data, and numbers, in the case of degraded data, refer to *post-hoc* tests. The same letter or number indicates that there are no significant (at $P < 0.05$) differences among methods and different letters or numbers mean that there is a significant difference.

best technique based on original data, and together with MAXENT-T it ranked highest in the error treatment. MAXENT-T results, on average, were not affected by shifting locations of data. The average performance of BRT has altered more than that of MAXENT on degraded data, but in the final analysis on degraded data the two are comparable. These two techniques also performed well across a comprehensive comparison study that included 16 modelling techniques (or variants), 226 species and six geographical regions (Elith *et al.* 2006), and across complementary analyses of the same data involving distinct analytical treatments (e.g. changing grain size; Guisan *et al.* 2007).

The regression-based models (GAM, GLM and MARS) showed no statistical change in model performance in the error treatment while models built using presence-only data (BIOCLIM, LIVES and DOMAIN) all declined in model performance. This difference can probably be attributed to whether models use absence data. In presence-only modelling techniques, there is no attempt to use absence data or create weighted response functions between occurrence records and environmental variables, which is probably the explanation for this decline in model performance. When absence data are used, a degraded presence point will be evaluated in the context of the absence data and this information will be used to build the model. Because presence-only techniques do not use background information they are more sensitive to inaccurate presences. For example, DOMAIN, which had the largest drop in performance between original and degraded data, uses the Gower distance metric to calculate the distance of a grid cell to the nearest (in climate space) occurrence

point, making it relatively sensitive to the occurrence of new combinations of the environmental variables. The large drop in performance of DOMAIN here is consistent with climate change studies where DOMAIN under-predicted species distributions severely under climate scenarios when compared to physiological distribution models (Hijmans & Graham 2006).

Differences among geographical regions might be influenced by three different factors: the environmental variables used in model building, the type of species modelled and the number of evaluation points. The types of environmental variables and the natural heterogeneity in the variables varied from region to region; however, this seems unlikely to have caused large differences in model performance across regions. While the spatial autocorrelation in SWI variables was the lowest among regions, potentially explaining the large influence of the error treatment in this region, variables for NZ were not more heterogeneous than AWT or NSW. The taxonomic groups modelled varied by region. In AWT and NSW a number of mobile species were modelled, whereas in NZ and SWI, plants which are potentially less mobile were modelled. Mobile species might be less influenced by locational error than sedentary species because their locations, even when accurate, represent their movement in the landscape over which they search for habitat and food resources (Guisan & Thuiller 2005). Finally, NZ and SWI had the largest amount of evaluation data available, so finer differences among models were detected more easily. A similar result was found by Elith *et al.* (2006); more discrimination among modelling methods was afforded by larger amounts of evaluation data.

We assumed that a 5 km shift in geographical coordinates of occurrence data was a reasonable amount to evaluate the influence of locational error on model performance. None the less, this represents a single data degradation treatment, and in some cases larger errors will occur in museum data. While we are not aware of research where data are manipulated explicitly, there are studies that have attempted to build models with imprecise occurrences. These studies have had mixed success. Lloyd & Palmer (1998) determined that they could model accurately three species of South African birds by choosing a random point within a 0.25 degree resolution grid where the species was present, and using this occurrence with environmental data of 0.0833 degrees (5 min) to run a model. McPherson, Jetz & Rogers (2006) followed a similar methodology using Ugandan breeding bird atlas data and obtained very poor results, although the difference in resolution was larger than the South African study (distribution data had a 0.25 degree resolution and the environmental data had a 0.01 degree resolution). Environments in Uganda might also be more spatially heterogeneous within the large grid cells than those of South Africa, but we cannot compare quantitatively the relative amount of spatial heterogeneity in each region. Further, these authors suggest that the Ugandan data are highly biased and this might have affected their results. In the current study, we had the advantage of using high-quality data that likely covered the environmental space of the regions, so that the bias usually encountered in occurrence data (e.g. roadside bias, accessibility, etc.) was unlikely to have

influenced our results. None the less, given the mixed results associated with the influence of locational error on model performance in these atlas studies and the current study, further research is warranted.

There are several avenues for future studies. First, multiple data degradation treatments could be applied to explore further the effects of error on model success. Further, because imprecision might not only be in errors in occurrence records but could also be in environmental variables used as predictors in distribution models, it might be useful to study both types of uncertainty simultaneously. Secondly, we have evaluated how error influences predictive success of models but not model inference. It would be useful to determine how the relative importance and shape of predictor variables changes with error. Thirdly, alternative techniques for dealing with locational uncertainty could be explored. For example, instead of taking the value of environmental layers at a single occurrence point, a mean (or median) value could be taken from a larger area. If uncertainty has been reported, the size of the area could be based on this uncertainty. In this approach, models are parameterized using a higher spatial resolution and predicted to a lower resolution. This approach worked well for British birds (Araujo *et al.* 2005) and otters in Spain and Portugal (Barbosa *et al.* 2003), but poorly for plant species in Great Britain (Collingham *et al.* 2000) and birds in Uganda (McPherson, Jetz & Rogers 2006).

A fourth area of suggested research deals with when and how to use inaccurate data, given that data with a relatively large error footprint/uncertainty may degrade models (Engler, Guisan & Rechsteiner 2004). In many species distribution modelling studies researchers will have a range of different qualities of data and could partition occurrence data by accuracy and then use a hierarchical modelling approach (Pearson & Dawson 2003; Pearson, Dawson & Lui 2004). In this approach environmental data are resampled to correspond to the uncertainty of a given partition of occurrence data and a model is created. This model is then used as an environmental input into a subsequent model built at a finer resolution with higher-quality occurrence data (Pearson & Dawson 2003; Pearson, Dawson & Lui 2004). Alternatively, if the level of imprecision in the occurrence data is unknown or difficult to estimate (i.e. in data for mobile species that could be observed traversing habitat that may or may not be suitable), a multiscale model can be used (Pearson & Dawson 2003). Here, environmental variables sampled at different resolutions are included among explanatory variables for the model to pick from.

In summary, we have explored how error in occurrence data influences model performance. By means of manipulations, in which known amounts of error are introduced into occurrence data sets, we were able to evaluate the sensitivity of different modelling approaches to greater uncertainty in occurrence locations. Our results indicate not only which analytical approaches may be most robust to such error, but also permit the more general conclusion that species distribution modelling approaches in general are fairly robust to locational error, and that usable models can be built even when

occurrence data are imprecise. Our results are especially important for conservation practitioners who require accurate maps of the species they are trying to protect. Finally, we hope that our conclusions will prompt the further development and use of these methods for species conservation and management, including predictions of how species distribution will be altered by land use and climate change.

Acknowledgements

Data were kindly provided by organizations for whom a number of authors worked. We also thank the Missouri Botanical Garden, especially Robert Magill and Trisha Consiglio, for access to TROPICOS and Gentry transect databases; and Andrew Ford, CSIRO Atherton, for AWT PA plant records. We would like to thank T. Wohlgenuth and U. B. Braendi from WSL Switzerland for access to the Swiss data sets. This research was initiated in a working group at the National Center for Ecological Analysis and Synthesis (NCEAS), Santa Barbara, USA: 'Testing Alternative Methodologies for Modelling Species' Ecological Niches and Predicting Geographic Distributions', conceived of and led by Peterson and Moritz. C. Graham was supported partly by a NASA New Investigator Grant.

Working group participants

Robert P. Anderson, Miroslav Dudík, Jane Elith, Simon Ferrier, Catherine H. Graham, Antoine Guisan, Robert J. Hijmans, Falk Huettmann, John Leathwick, Anthony Lehmann, Jin Li, Lucia Lohmann, Bette Loiselle, Glenn Manion, Craig Moritz, Miguel Nakamura, Yoshinori Nakazawa, Jake Overton, A. Townsend Peterson, Steven Phillips, Karen Richardson, Ricardo Scachetti Pereira, Robert Schapire, Jorge Soberón, Stephen Williams, Mary Wisz, Niklaus Zimmermann.

References

Anderson, R.P., Gomez-Laverde, M. & Peterson, A.T. (2002) Geographical distributions of spiny pocket mice in South America: insights from predictive models. *Global Ecology and Biogeography*, **11**, 131–141.

Araujo, M.B., Thuiller, W., Williams, P.H. & Reginster, I. (2005) Downscaling European species atlas distributions to a finer resolution: implications for conservation planning. *Global Ecology and Biogeography*, **14**, 17–30.

Austin, M.P. (2002) Spatial prediction of species distribution: an interface between ecological theory and statistical modelling. *Ecological Modelling*, **157**, 101–118.

Barbosa, A.M., Real, R., Olivero, J. & Vargas, J.M. (2003) Otter (*Lutra lutra*) distribution modeling at two resolution scales suited to conservation planning in the Iberian Peninsula. *Biological Conservation*, **114**, 377–387.

Barry, S. & Elith, J. (2006) Error and uncertainty in habitat models. *Journal of Applied Ecology*, **43**, 413–423.

Busby, J.R. (1991) BIOCLIM – a bioclimate analysis and prediction system. *Nature Conservation: Cost Effective Biological Surveys and Data Analysis* (eds C.R. Margules & M.P. Austin), pp. 64–68. CSIRO, Canberra, Australia.

Carpenter, G., Gillison, A.N. & Winter, J. (1993) DOMAIN: a flexible modelling procedure for mapping potential distributions of plants and animals. *Biodiversity and Conservation*, **2**, 667–680.

Collingham, Y.C., Wadsworth, R.A., Huntley, B. & Hulme, P.E. (2000) Predicting the spatial distribution of non-indigenous riparian weeds: issues of spatial scale and extent. *Journal of Applied Ecology*, **37**, 13–27.

Elith, J., Graham, C.H., Anderson, R.P., Dudík, M., Ferrier, S., Guisan, A., Hijmans, R.J., Huettmann, F., Leathwick, J.R., Lehmann, A., Li, J., Lohmann, L.G., Loiselle, B.A., Manion, G., Moritz, C., Nakamura, M., Nakazawa, Y., Overton, J.M., Peterson, A.T., Phillips, S.J., Richardson, K., Scachetti-Pereira, R., Schapire, R.E., Soberón, J., Williams, S., Wisz, M.S. & Zimmermann, N.E. (2006) Novel methods improve prediction of species' distributions from occurrence data. *Ecography*, **29**, 129–151.

Engler, R., Guisan, A. & Rechsteiner, L. (2004) An improved approach for predicting the distribution of rare and endangered species from occurrence and pseudo-absence data. *Journal of Applied Ecology*, **41**, 263d274.

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

Friedman, J.H., Hastie, T. & Tibshirani, R. (2000) Additive logistic regression: a statistical view of boosting. *Annals of Statistics*, **28**, 337–407.

Gaubert, P., Papes, M. & Peterson, A.T. (2006) Natural history collections and the conservation of poorly known taxa: ecological niche modeling in central African rainforest genets (*Genetta* spp.). *Biological Conservation*, **130**, 106–117.

Graham, C.H., Ferrier, S., Huettmann, F., Moritz, C. & Peterson, A.T. (2004) New developments in museum-based informatics and applications in biodiversity analysis. *Trends in Ecology and Evolution*, **19**, 497–503.

Guisan, A., Edwards, T.C. & Hastie, T. (2002) Generalized linear and generalized additive models in studies of species distributions: setting the scene. *Ecological Modelling*, **157**, 89–100.

Guisan, A., Graham, C.H., Elith, J., Huettmann, F. & NCEAS Species Distribution Modelling Group (2007) Sensitivity of predictive species distribution models to change in grain size. *Diversity and Distributions*, **13**, 332–340.

Guisan, A. & Thuiller, W. (2005) Predicting species distribution: offering more than simple habitat models. *Ecology Letters*, **8**, 993–1009.

Guralnick, R.P., Wiecek, J., Beaman, R. & Hijmans, R.J. and the Bio-Geomancer working Group (2006) BioGeomancer: automated georeferencing to map the world's biodiversity data. *Plos Biology*, **4**, 1908–1909.

Hanley, J.A. & McNeil, B.J. (1982) The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology*, **143**, 29–36.

Hernandez, P., Graham, C.H., Master, L.L. & Albert, D.L. (2006) A comparison of the performance of species distribution models using a range of species' occurrences. *Ecography*, **29**, 773–785.

Hijmans, R.J., Garrett, K.A., Huamán, Z., Zhang, D.P., Schreuder, M. & Bonierbale, M. (2000) Assessing the geographic representativeness of genebank collections: the case of Bolivian wild potatoes. *Conservation Biology*, **14**, 1755–1765.

Hijmans, R.J. & Graham, C.H. (2006) Testing the ability of climate envelope models to predict the effect of climate change on species distributions. *Global Change Biology*, **12**, 2272–2281.

Hijmans, R.J., Schreuder, M., de la Cruz, J. & Guarino, L. (1999) Using GIS to check co-ordinates of germplasm accessions. *Genetic Resources and Crop Evolution*, **46**, 291–296.

Leathwick, J.R., Rowe, D., Richardson, J., Elith, J. & Hastie, T. (2005) Using multivariate adaptive regression splines to predict the distributions of New Zealand's freshwater diadromous fish. *Freshwater Biology*, **50**, 2034–2052.

Lehmann, A., Overton, J.M. & Leathwick, J.R. (2003) GRASP: generalized regression analysis and spatial prediction. *Ecological Modelling*, **160**, 165–183.

Li, J. & Hilbert, D. (in press) A new predictive model, LIVES, for the potential distributions and habitats of species using presence-only data. *Ecological Modelling*.

Lloyd, P. & Palmer, A.R. (1998) Abiotic factors as predictors of distribution in southern African bulbuls. *Auk*, **115**, 404–411.

Loiselle, B.A., Howell, C.A., Graham, C.H., Goerck, J.M., Brooks, T., Smith, K.G. & Williams, P.H. (2003) Avoiding pitfalls of using species distribution models in conservation planning. *Conservation Biology*, **17**, 1591–1600.

McPherson, J.M., Jetz, W. & Rogers, D.J. (2004) The effects of species' range sizes on the accuracy of distribution models: ecological phenomenon or statistical artefact? *Journal of Applied Ecology*, **41**, 811–823.

McPherson, J.M., Jetz, W. & Rogers, D.J. (2006) Using coarse-grained occurrence data to predict species distributions at finer spatial resolutions – possibilities and limitations. *Ecological Modelling*, **192**, 499–522.

Pearce, J. & Ferrier, S. (2000) An evaluation of alternative algorithms for fitting species distribution models using logistic regression. *Ecological Modelling*, **128**, 127–147.

Pearson, R.G. & Dawson, T.P. (2003) Predicting the impacts of climate change on the distribution of species: are bioclimate envelope models useful? *Global Ecology and Biogeography*, **12**, 361–371.

Pearson, R.G., Dawson, T.P. & Liu, C. (2004) Modelling species distributions in Britain: a hierarchical integration of climate and land-cover data. *Ecography*, **27**, 285–298.

Peterson, A.T. (2003) Predicting the geography of species' invasions via ecological niche modeling. *Quarterly Review of Biology*, **78**, 419–433.

Peterson, A.T. (2007) Uses and requirements of ecological niche models and related distributional models. *Biodiversity Informatics*, **3**, 59–72.

Phillips, S.J., Anderson, R.P. & Schapire, R.E. (2006) Maximum entropy modeling of species geographic distributions. *Ecological Modelling*, **190**, 231–259.

Rissler, L.J., Hijmans, R.J., Graham, C.H., Moritz, C. & Wake, D.B. (2006) Phylogeographic lineages and species comparisons in conservation analyses: a case study of California herpetofauna. *American Naturalist*, **167**, 655–666.

- Rowe, R.J. (2005) Elevational gradient analyses and the use of historical museum specimens: a cautionary tale. *Journal of Biogeography*, **32**, 1883–1897.
- Schapire, R. (2003) The boosting approach to machine learning – an overview. *MSRI Workshop on Nonlinear Estimation and Classification, 2002* (eds D.D. Denison, M.H. Hansen, C. Holmes, B. Mallick & B. Yu), pp. 1–23. Springer, NY.
- Stockwell, D. & Peters, D. (1999) The GARP modelling system: problems and solutions to automated spatial prediction. *International Journal of Geographic Information Science*, **13**, 143–158.
- Suarez, A.V. & Tsutsui, N.D. (2004) The value of museum collections for research and society. *Bioscience*, **54**, 66–74.
- Swets, J.A. (1988) Measuring the accuracy of diagnostic systems. *Science*, **240**, 1285–1293.
- Walker, P.A. & Cocks, K.D. (1991) HABITAT: a procedure for modelling a disjoint environmental envelope for a plant or animal species. *Global Ecology and Biogeography Letters*, **1**, 108–118.
- Wieczorek, J.R., Guo, Q. & Hijmans, R.J. (2004) The point-radius method for georeferencing point localities and calculating associated uncertainty. *International Journal of Geographic Information Science*, **18**, 745–767.
- Wintle, B.A., Elith, J. & Potts, J. (2005) Fauna habitat modelling and mapping in an urbanising environment; a case study in the Lower Hunter Central Coast region of NSW. *Austral Ecology*, **30**, 729–748.
- Yee, T.W. & Mitchell, N.D. (1991) Generalized additive models in plant ecology. *Journal of Vegetation Science*, **2**, 587–602.

Received 30 October 2006; accepted 14 August 2007
 Handling Editor: Robert Freckleton

Supplementary material

The following supplementary material is available for this article.

Appendix S1. Environmental GIS predictors and descriptions for each data set.

Appendix S2. Species and their occurrence in each data set.

This material is available as part of the online article from: <http://www.blackwell-synergy.com/doi/full/10.1111/j.1365-2664.2007.01408.x>

(This link will take you to the article abstract.)

Please note: Blackwell Publishing is not responsible for the content or functionality of any supplementary material supplied by the authors. Any queries (other than missing material) should be directed to the corresponding author for the article.

Copyright of *Journal of Applied Ecology* is the property of Blackwell Publishing Limited and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.