

**TITLE:**

A statistical explanation of MaxEnt for ecologists

**AUTHORS:**

Jane Elith<sup>1</sup>, Steven Phillips<sup>2</sup>, Trevor Hastie<sup>3</sup>, Miroslav Dudík<sup>4</sup>, Yung En Chee<sup>1</sup>, Colin Yates<sup>5</sup>

**ADDRESSES:**

<sup>1</sup>School of Botany, The University of Melbourne, Parkville, VIC 3010 Australia

<sup>2</sup> AT&T Labs – Research, 180 Park Avenue, Florham Park, NJ 07932, U.S.A.

<sup>3</sup> Department of Statistics, Stanford University, CA 94305, USA

<sup>4</sup> Yahoo! Labs, 111 West 40th Street (17th Floor). New York, NY 10018

<sup>5</sup> Science Division, Western Australian Department of Environment and Conservation, LMB 104, Bentley Delivery Centre, WA6983, Australia.

**KEYWORDS:**

species distribution model, ecological niche, presence-only, absence, machine learning, river, banksias

**SHORT RUNNING TITLE:**

Statistical explanation of MaxEnt

## 1 ABSTRACT:

2 MaxEnt is a program for modelling species distributions from presence-only species records. This  
3 paper is written for ecologists and describes the MaxEnt model from a statistical perspective,  
4 making explicit links between the structure of the model, decisions required in producing a  
5 modelled distribution, and knowledge about the species and the data that might affect those  
6 decisions. To begin we discuss the characteristics of presence-only data, highlighting implications  
7 for modelling distributions. We particularly focus on the problems of sample bias and lack of  
8 information on species prevalence. The keystone of the paper is a new statistical explanation of  
9 MaxEnt which shows that the model minimizes the relative entropy between two probability  
10 densities (one estimated from the presence data and one, from the landscape) defined in covariate  
11 space. For many users, this viewpoint is likely to be a more accessible way to understand the model  
12 than previous ones that rely on machine learning concepts. We then step through a detailed  
13 explanation of MaxEnt describing key components (e.g. covariates and features, and definition of  
14 the landscape extent), the mechanics of model fitting (e.g. feature selection, constraints and  
15 regularization) and outputs. Using case studies for a *Banksia* species native to south-west Australia  
16 and a riverine fish, we fit models and interpret them, exploring why certain choices affect the result  
17 and what this means. The fish example illustrates use of the model with vector data for linear river  
18 segments rather than raster (gridded) data. Appropriate treatments for survey bias, unprojected  
19 data, locally restricted species, and predicting to environments outside the range of the training  
20 data are demonstrated, and new capabilities discussed. Online appendices include additional details  
21 of the model and the mathematical links between previous explanations and this one, example code  
22 and data, and further information on the case studies.

## 23 INTRODUCTION

24 Species distribution models (SDMs) estimate the relationship between species records at sites and  
25 the environmental and/or spatial characteristics of those sites (Franklin, 2009). They are widely  
26 used for many purposes in biogeography, conservation biology and ecology (Elith & Leathwick,  
27 2009a; Table 1). In the last two decades there have been many developments in the field of species  
28 distribution modelling, and multiple methods are now available. A major distinction among  
29 methods is the kind of species data they use. Where species data have been collected systematically  
30 – for instance, in formal biological surveys in which a set of sites are surveyed and the  
31 presence/absence or abundance of species at each site are recorded - regression methods familiar  
32 to most ecologists (e.g. generalized linear or additive models, GLMs or GAMs; or ensembles of  
33 regression trees: random forests or boosted regression trees, BRT) are used.

34  
35 However for most regions systematic biological survey data tend to be sparse and/or limited in  
36 coverage. Species records are available, though, in the form of presence-only records in herbarium  
37 and museum databases. Many of these databases represent well over a century of public and  
38 private investment in biological science and are a hugely important resource of species occurrence  
39 data. The desire to maximize the utility of such resources has spawned an array of SDM methods for  
40 modelling presence-only data. MaxEnt (Phillips *et al.*, 2006; Phillips & Dudík, 2008) is one such  
41 method and is the focus of this paper.

42  
43 MaxEnt's predictive performance is consistently competitive with the highest performing methods  
44 (Elith *et al.*, 2006). Since becoming available in 2004, it has been utilized extensively for modelling  
45 species distributions. Published examples cover diverse aims (finding correlates of species  
46 occurrences, mapping current distributions, and predicting to new times and places) across many  
47 ecological, evolutionary, conservation and biosecurity applications (Table 1). Government and non-  
48 government organisations have also adopted MaxEnt for large-scale, real-world biodiversity  
49 mapping applications, including the Point Reyes Bird Observatory online application  
50 (<http://www.prbo.org/>) and the Atlas of Living Australia (<http://www.ala.org.au/>). JE and SJP's  
51 involvement in such programs identified a need for an ecologically-accessible explanation of

52 MaxEnt. Existing descriptions include concepts from machine learning that tend to be outside the  
53 common experience of many ecologists.

54  
55 In this paper we explain the MaxEnt modelling method with emphases on a statistical explanation  
56 of the method, on what it assumes, and on the impacts of choices made in the modelling process.  
57 We use two case studies to examine the effects of background selection and model settings, and to  
58 illustrate the applicability of the model for exploring ecological relationships with fine-scale, vector-  
59 based environmental data. Our aim is to promote understanding of the method and recommend  
60 useful approaches to data preparation and model fitting and interpretation.

## 61 PREAMBLE: WHAT IS SPECIAL ABOUT THE PRESENCE-ONLY CASE?

62 Expanding use of presence-only data for modelling species distributions has prompted wide  
63 discussion about the sorts of distributions (e.g., potential vs realized) that can be modelled with  
64 presence-only data in contrast to presence-absence data (e.g., Soberón & Peterson, 2005; Chefaoui  
65 & Lobo, 2007; Hirzel & Le Lay, 2008; Jiménez-Valverde *et al.*, 2008; Soberón & Nakamura, 2009;  
66 Lobo *et al.*, 2010). As mentioned in several of these papers, the subject is complex due to the  
67 interplay of data quality (amount and accuracy of species data; ecological relevance of predictor  
68 variables; availability of information on disturbances, dispersal limitations and biotic interactions),  
69 modelling method and scale of analysis. A comprehensive review of the issues would be useful, but  
70 here we restrict ourselves to key points important for this paper.

71  
72 Some of the published discussion suggests that presence-only data in some sense release us from  
73 the problems of unreliable absence records (e.g., Jiménez-Valverde *et al.*, 2008), particularly  
74 emphasising that absences bear such strong imprints of biotic interactions, dispersal constraints  
75 and disturbances that they may preclude modelling of potential distributions (*sensu* Svenning &  
76 Skov, 2004). However, the presence records are also imprinted by many of the factors affecting  
77 absences. If a species is absent from an environmentally suitable area because, say, past  
78 disturbances have caused local extinctions, the signal of that absence will be found in the  
79 distribution of presence records: there will be no presence records in the disturbed area.  
80 Regardless of whether absences are used in modelling, the pattern in the presence records will  
81 suggest the area is unsuitable, and the model will be affected by this patterning. Similarly, if the  
82 detectability of a particular species varies from site to site, then not only does this result in some  
83 false absences in presence-absence data, it also affects the pattern of presences in presence-only  
84 data. This leads naturally to the conclusion that dispensing with absences does not address the  
85 limitations often attributed to absence data, such as the fact that species are not perfectly  
86 detectable and may not occupy all suitable habitat. This thinking means that we will approach the  
87 description of the presence-only modelling problem as one that is trying to model the same  
88 quantity that is modelled with presence-absence data, that is: the probability of presence of a  
89 species (to be defined more carefully below).

90  
91 From here on, we assume that the data available to the modeler are presence-only, i.e. a set of  
92 locations within  $L$ , the landscape of interest, where the species has been observed. Let  $y=1$  denote  
93 presence,  $y=0$  denote absence,  $\mathbf{z}$  denote a vector of environmental covariates, and background be  
94 defined as all locations within  $L$  (or a random sample thereof). Assume the environmental variables  
95 or covariates  $\mathbf{z}$  (representing environmental conditions) are available landscape wide. Define  $f(\mathbf{z})$   
96 to be the probability density of covariates across  $L$ ,  $f_1(\mathbf{z})$  to be the probability density of covariates  
97 across locations within  $L$  where the species is present, and similarly,  $f_0(\mathbf{z})$  where the species is  
98 absent. (Densities – or probability density functions – describe the relative likelihood of random  
99 variables over their range, and can be univariate or multivariate). The quantity that we wish to  
100 estimate is, as with presence-absence data, the probability of presence of the species, conditioned  
101 on environment:  $\Pr(y=1|\mathbf{z})$ . Strictly presence-only data only allows us to model  $f_1(\mathbf{z})$ , which on its  
102 own cannot approximate probability of presence. Presence/background data allows us to model  
103 both  $f_1(\mathbf{z})$  and  $f(\mathbf{z})$ , and this gets to within a constant of  $\Pr(y=1|\mathbf{z})$ , because Bayes' rule gives:

104

$$105 \quad \Pr(y=1|\mathbf{z}) = f_1(\mathbf{z})\Pr(y=1) / f(\mathbf{z}). \quad \dots\dots\dots (1)$$

106 The only quantity that is lacking is the second term,  $\Pr(y=1)$ , i.e. the prevalence of the species  
 107 (proportion of occupied sites) in the landscape. Formally, we say that prevalence is not identifiable  
 108 from presence-only data (Ward *et al.* 2009). This means that it cannot be exactly determined,  
 109 regardless of the sample size; this is a fundamental limitation of presence-only data. As an aside we  
 110 note, however, that absence data are plagued by issues of detection probability (Wintle *et al.*, 2004;  
 111 MacKenzie, 2005) so that even presence-absence data may not yield a good estimate of prevalence.  
 112 A second fundamental limitation of presence-only data is that sample selection bias (whereby some  
 113 areas in the landscape are sampled more intensively than others) has a much stronger effect on  
 114 presence-only models than on presence-absence models (Phillips *et al.*, 2009). Imagine that  $f_1(\mathbf{z})$  is  
 115 contaminated by a sample selection bias  $s(\mathbf{z})$ . This bias will most commonly occur in geographic  
 116 space (e.g. close to roads) but could be environmentally based (e.g. visiting wet gullies) but,  
 117 regardless, will map into covariate space. Under biased sampling, a presence-only model gives an  
 118 estimate of  $f_1(\mathbf{z})s(\mathbf{z})$  rather than  $f_1(\mathbf{z})$ . That is, we get a model that combines the species distribution  
 119 with the distribution of sampling effort (Soberón & Nakamura, 2009). In contrast, for presence-  
 120 absence models, sample selection bias affects both presence and absence records, and the effect of  
 121 the bias cancels out (under reasonable assumptions, see Zadrozny, 2004).

122  
 123 So far we have treated presence or absence as a binary event, but in reality defining the response  
 124 variable is not straightforward, and in this regard presence-only data are quite different from  
 125 presence-absence data (Pearce & Boyce, 2006). Presence or absence of a species is dependent on  
 126 the time frame and spatial scale -- for example, a vagile species (such as a bird) may be present at  
 127 some times but not others, while a plant species will be more likely to be found in a large plot with  
 128 given environmental conditions than in a small plot with the same conditions. Absence of a plant  
 129 species from a 1km<sup>2</sup> quadrat around a point implies absence in a 1m<sup>2</sup> quadrat around that point,  
 130 but not vice versa. With presence-absence data, it is not hard to incorporate these complexities in  
 131 the formulation of the response variable ( i.e., the specification of what constitutes a sample), or via  
 132 sampling covariates in the model, provided survey details are available (Leathwick, 1998;  
 133 MacKenzie & Royle, 2005; Schulman *et al.*, 2007; Ward, 2007b). However, with presence-only data,  
 134 we typically have occurrence data that do not have any associated temporal or spatial scale. The  
 135 record is usually simply a record of the species at a location, with no information on search area or  
 136 time.

137  
 138 With presence-absence data, the definition of the response variable should naturally be consistent  
 139 with the sampling method: for example, if the available data are surveys of 1m<sup>2</sup> quadrats, then  $y=1$   
 140 should correspond to the species being present in a 1m<sup>2</sup> quadrat. With presence-only data, the  
 141 available data do not usually describe the survey method, so the modeller has considerable leeway  
 142 in defining the response variable. A common approach is to implicitly assume a sampling unit of  
 143 size equal to the grain size of available environmental data (see Elith & Leathwick 2009a for  
 144 discussion of grain).

145  
 146 To summarize, we posit that with presence and background data, we can model the same quantity  
 147 as with presence-absence data, up to the constant  $\Pr(y=1)$ . However, if presence-absence survey  
 148 data are available, we believe it is generally advisable to use a presence-absence modelling method,  
 149 since in that case the models are less susceptible to problems of sample selection bias, the survey  
 150 method will often be known and can be used to appropriately define the response variable for  
 151 modelling, and we take advantage of all information in the data. In particular, presence-absence  
 152 data give us much better information about prevalence than presence-only, because – even though  
 153 there may be some difficulties due to imperfect detection - they solve the major problem of non-  
 154 identifiability. We will come back to this when we discuss the logistic output of MaxEnt.

155

## 156 EXPLANATION OF MAXENT

157 Here for the first time we describe MaxEnt using statistical terminology and notation, providing a  
 158 break from the machine learning terminology in previous papers. As we describe the model we will  
 159 highlight possibilities for – and implications of - modelling choices and defaults, and consider how  
 160 MaxEnt addresses the limitations of presence-only data identified above. We relegate the more  
 161 technical considerations to boxes and online appendices, to avoid interrupting the flow of the  
 162 explanation.

### 163 COVARIATES AND FEATURES

164 Most ecologists, following the statistical literature, call the independent variables in a model the  
 165 covariates, predictors or inputs. In SDMs these include environmental factors that are relevant to  
 166 habitat suitability (e.g. estimates of climate, topography, and soil for plants; temperature, salinity  
 167 and prey abundance for marine fishes). Since species' responses to these tend to be complex, it is  
 168 usually desirable to fit non-linear functions (Austin, 2002). In regression, this can be achieved by  
 169 applying transformations to the covariates – for instance, creating basis functions for polynomials  
 170 and splines, including piecewise linear functions. Complex models are fitted as linear combinations  
 171 of these basis functions in methods including GLMs and GAMs (Hastie *et al.* 2009, Chapter 5). In  
 172 machine learning, basis functions and other transformations of available data are termed features –  
 173 i.e., features are an expanded set of transformations of the original covariates.

174  
 175 In MaxEnt selected features are formed “behind the scenes”, in the same way as in regression,  
 176 where the model matrix is augmented by terms specified in the model (e.g. polynomials,  
 177 interactions). The MaxEnt fitted function is usually defined over many features, meaning that in  
 178 most models there will be more features than covariates. MaxEnt currently has six feature classes:  
 179 linear, product, quadratic, hinge, threshold and categorical (further details in Online Appendix 1).  
 180 Products are products of all possible pair-wise combinations of covariates, allowing simple  
 181 interactions to be fitted. Threshold features allow a “step” in the fitted function; hinge features are  
 182 similar except they allow a change in gradient of the response. Many threshold or hinge features  
 183 can be fitted for one covariate, giving a potentially complex function. Hinge features (which are  
 184 basis functions for piecewise linear splines), if used alone, allow a model rather like a generalized  
 185 additive model (GAM): an additive model, with non-linear fitted functions of varying complexity but  
 186 without the sudden steps of the threshold features. MaxEnt’s default is to allow all feature types  
 187 (conditional on sufficient species data being available), but it is worth considering simpler models,  
 188 as discussed later under implications for modelling.

### 189 THE MAXENT MODEL – A SHORT OVERVIEW

190 Previous papers have described MaxEnt as estimating a distribution across geographic space  
 191 (Phillips *et al.*, 2006; Phillips & Dudík, 2008). Here we give a different (but equivalent)  
 192 characterization that focuses on comparing probability densities in covariate space (Figure 1). In  
 193 doing so we rely strongly on the PhD research of TH's past student, Gill Ward (Ward, 2007b), and  
 194 acknowledge her contribution. Equation 1 shows that if we know the conditional density of the  
 195 covariates at the presence sites,  $f_1(\mathbf{z})$ , and the marginal (i.e. unconditional) density of covariates  
 196 across the study area  $f(\mathbf{z})$ , we then only need knowledge of the prevalence  $\Pr(y=1)$ , to calculate  
 197 conditional probability of occurrence. MaxEnt first makes an estimate of the ratio  $f_1(\mathbf{z})/f(\mathbf{z})$ ,  
 198 referred to as MaxEnt’s “raw” output. This is the core of the MaxEnt model output, giving insight  
 199 about what features are important, and estimating the relative suitability of one place vs another.  
 200 Because the required information on prevalence is not available for calculating conditional  
 201 probability of occurrence, a work-around has been implemented (termed MaxEnt’s “logistic”  
 202 output). This treats the log of the output ---  $\eta(\mathbf{z})=\log(f_1(\mathbf{z})/f(\mathbf{z}))$  --- as a logit score, and calibrates  
 203 the intercept so that the implied probability of presence at sites with “typical” conditions for the  
 204 species (i.e. where  $\eta(\mathbf{z})$  = the average value of  $\eta(\mathbf{z})$  under  $f_1$ ) is a parameter  $\tau$ . Knowledge of  $\tau$  would  
 205 solve the non-identifiability of prevalence, and in the absence of that knowledge MaxEnt arbitrarily

206 sets  $\tau$  to equal 0.5. This logistic transformation is monotone (order preserving) with the raw  
 207 output. We work through each part of the MaxEnt model in the following sections, showing how the  
 208 choice of landscape, species data, and selected settings influence the results.

## 209 THE LANDSCAPE AND SPECIES RECORDS

210 The landscape of interest ( $L$ ) is a geographic area suggested by the problem and defined by the  
 211 ecologist. It might, for instance, be limited by geographic boundaries, or by an understanding of  
 212 how far the focal species could have dispersed. We then define  $L_1$  as the subset of  $L$  where the  
 213 species is present.

214  
 215 The distribution of covariates in the landscape is conveyed by a finite sample – a collection of points  
 216 from  $L$  with associated covariates, typically called a background sample. These data may be  
 217 supplied in the form of grids of covariates covering a pixelation of the landscape; as a default  
 218 MaxEnt randomly samples 10,000 background locations from covariate grids, but the background  
 219 data points can also be specified (see Yates *et al.*, 2010, and case studies below) and grids are not  
 220 essential (case study 2). Note that the background sample does not take any account of the  
 221 presence locations – it is simply a sample of  $L$ , and could by chance include presence locations.  
 222 Using a random background sample implies a belief that the sample of presence records is also a  
 223 random sample from  $L_1$ . We deal later with the case of biased samples.

## 224 DESCRIPTION OF THE MODEL

225 MaxEnt uses the covariate data from the occurrence records and the background sample to  
 226 estimate the ratio  $f_i(\mathbf{z})/f(\mathbf{z})$ . It does this by making an estimate of  $f_i(\mathbf{z})$  that is consistent with the  
 227 occurrence data; many such distributions are possible, but it chooses the one that is closest to  $f(\mathbf{z})$ .  
 228 Minimizing distance from  $f(\mathbf{z})$  is sensible, because  $f(\mathbf{z})$  is a null model for  $f_i(\mathbf{z})$ : without any  
 229 occurrence data, we would have no reason to expect the species to prefer any particular  
 230 environmental conditions over any others, so we could do no better than predict that the species  
 231 occupies environmental conditions proportionally to their availability in the landscape. In MaxEnt,  
 232 this distance from  $f(\mathbf{z})$  is taken to be the relative entropy of  $f_i(\mathbf{z})$  with respect to  $f(\mathbf{z})$  (also known as  
 233 the Kullback-Leibler divergence).

234 Using background data informs the model about  $f(\mathbf{z})$ , the density of covariates in the region, and  
 235 provides the basis for comparison with the density of covariates occupied by the species – i.e.  $f_i(\mathbf{z})$ .  
 236 Constraints are imposed so that the solution is one that reflects information from the presence  
 237 records. For example, if one covariate is summer rainfall, then constraints ensure that the mean  
 238 summer rainfall for the estimate of  $f_i(\mathbf{z})$  is close to its mean across the locations with observed  
 239 presences. The species' distribution is thus estimated by minimizing the distance between  $f_i(\mathbf{z})$  and  
 240  $f(\mathbf{z})$  subject to constraining the mean summer rainfall estimated by  $f_i$  (and the means of other  
 241 covariates) to be close to the mean across presence locations.

242  
 243 We note that previous papers describing MaxEnt focused on a location-based definition over a finite  
 244 landscape (typically a grid of pixels). We will call this a definition based in geographic space, and  
 245 compare it with our new description, which focuses on environmental (covariate) space. Note,  
 246 though, that we are not implying by this wording that in either definition there is any consideration  
 247 of the geographic proximity of locations unless geographic predictors are used. In the original  
 248 definition (Phillips *et al.* 2006), the target was  $\pi(x) = p(x|y=1)$ , which was a probability distribution  
 249 over pixels (or locations)  $x$ . This was called the “raw” distribution (Phillips *et al.* 2006), and gave  
 250 the probability, given the species is present, that it is found at pixel  $x$ . Maximising the entropy of the  
 251 raw distribution is equivalent to minimizing the relative entropy of  $f_i(\mathbf{z})$  relative to  $f(\mathbf{z})$ , so the two  
 252 formulations are equivalent (see online Appendix 2 for equations showing the transition from the  
 253 geographic to environmental definitions). The null model for the raw distribution was the uniform  
 254 distribution over the landscape, since without any data we would have no reason to think the  
 255 species would prefer any location to any other. In environmental space the equivalent null model  
 256 for  $\mathbf{z}$  is  $f(\mathbf{z})$ , since without any data, we have no reason to think the species prefers any particular

257 environmental conditions, and therefore occupies environmental conditions in proportion to how  
 258 prevalent they are in the landscape.

259  
 260 Constraints were described above in reference to covariates, but – as explained in the section on  
 261 covariates and features - MaxEnt actually fits the model on features that are transformations of the  
 262 covariates. These allow potentially complex relationships to be modelled. The constraints are  
 263 extended from being constraints on the means of covariates to being constraints on the means of  
 264 the features. We will call the vector of features  $h(\mathbf{z})$  and the vector of coefficients  $\beta$  (note, this  
 265 notation is different to previous papers: Table 2). As explained in Phillips *et al.* (2006), minimizing  
 266 relative entropy results in a Gibbs distribution (Della Pietra *et al.*, 1997) which is an exponential-  
 267 family model:

$$268 \quad f_i(\mathbf{z}) = f(\mathbf{z}) e^{\eta(\mathbf{z})} \quad \dots\dots\dots (2)$$

270 where  $\eta(\mathbf{z}) = \alpha + \beta \cdot h(\mathbf{z})$   
 271 and  $\alpha$  is a normalizing constant that ensures that  $f_i(\mathbf{z})$  integrates (sums) to 1  
 272

273 From this it is clear that the target of a MaxEnt model is  $\eta(\mathbf{z})$ , which estimates the ratio  $f_i(\mathbf{z})/f(\mathbf{z})$ . It  
 274 is a log-linear model, similar in form to a GLM, and depends on both the presence samples and the  
 275 background samples that are used in forming the estimate. Hence the definition of the landscape is  
 276 intimately linked to the solution that is given.

## 277 MECHANICS OF THE SOLUTION

278 In coming to a solution MaxEnt needs to find coefficients (betas) that will result in the constraints  
 279 being satisfied but not match them so closely that it overfits and produces a model with limited  
 280 generalization. MaxEnt handles the issue by setting an error bound, or maximum allowed deviation  
 281 from the sample (empirical) feature means. MaxEnt first automatically rescales all features to have  
 282 the range 0 to 1. Then an error bound ( $\lambda_j$  in equation 3) is calculated for each feature (again note the  
 283 change in notation from previous papers, Table 2). It will reflect the variation in sample values for  
 284 that feature, adjusted by a tuned (pre-set) parameter for the feature class (Phillips and Dudík 2008,  
 285 and equation 3). MaxEnt *could* estimate feature error bounds only from the data, for example using  
 286 cross-validation, but to simplify model fitting and because the data are often biased, it uses feature  
 287 class-specific tuned parameters based on a large international data set (Phillips & Dudík, 2008).  
 288 That dataset covers 226 species, 6 regions of the world, sample sizes ranging from 2 to 5822, and  
 289 11-13 predictors per region (Elith *et al.* 2006). It is possible that the tuning may not work well for  
 290 very different datasets – e.g. if there are many more predictors. The tuned parameters can be  
 291 changed by the user if desired. The pre-tuning also includes restrictions to the set of feature classes  
 292 that will be considered for small samples.

$$293 \quad \lambda_j = \lambda \sqrt{\frac{s^2[h_j]}{m}} \quad \dots\dots\dots (3)$$

295 where  $\lambda_j$  is the regularization parameter for feature  $h_j$ . This feature's variance is  $s^2$  over the  
 296  $m$  presence sites, and its feature class has a tuning parameter  $\lambda$ . Conceptually  $\lambda_j$   
 297 corresponds to the width of the confidence interval, therefore it takes the form of the  
 298 standard error (the square root expression) multiplied by the parameter  $\lambda$  according to the  
 299 desired confidence level.

300  
 301 The lambdas in equation 3 allow regularization – i.e. smoothing the distribution, making it more  
 302 regular. These error bounds are a specific form of regularization called L1-regularization  
 303 (Tibshirani, 1996) that gives sparse solutions (ones with many zeros, i.e. many features removed).  
 304 Regularization is not specific to MaxEnt; it is a common modern approach to model selection. It can

305 be thought of as a way of shrinking the coefficients (the betas) – i.e. penalizing them - to values that  
 306 balance fit and complexity, allowing both accurate prediction and generality. In MaxEnt, the fit of  
 307 the model is measured at the occurrence sites, using a log likelihood (Box 1). A highly complex  
 308 model will have a high log likelihood, but may not generalize well. The aim of regularization is to  
 309 trade off model fit (the first term in equation 4 below) and model complexity (the second term in  
 310 equation 4). In this sense MaxEnt fits a penalized maximum likelihood model (Phillips and Dudík  
 311 2008; equation 4) closely related to other penalties for complexity such as Akaike's Information  
 312 Criterion (AIC, Akaike, 1974). Maximising the penalized log likelihood is equivalent to minimising  
 313 the relative entropy subject to the error bound constraints.

$$\max_{\alpha, \beta} \frac{1}{m} \sum_{i=1}^m \ln(f(z_i) e^{\eta(z_i)}) - \sum_{j=1}^n \lambda_j |\beta_j|$$

315 ..... equation 4

316 subject to  $\int_L f(z) e^{\eta(z)} dz = 1$

317

318 Where  $z_i$  is the feature vector for occurrence point  $i$  of  $m$  sites,  
 319 and for  $j = 1 \dots n$  features

320

321 ----- **Box 1 - Log likelihood** -----

322 In statistics a log likelihood describes the log of the probability of an observed outcome. It  
 323 varies from 0 [ $\ln(1)$ ] to negative infinity [ $\ln(0)$ ]. If the space of outcomes is continuous, we  
 324 measure the probability density at the observed outcome, rather than probability. With  
 325 presence-only data the only known outcomes are presences, so when measuring  
 326 likelihoods, the calculation is simply done at the presence sites (in comparison to logistic  
 327 regression where they are calculated at presence and absence sites). For a set of  
 328 observations the average log likelihood is estimated. When fitting a MaxEnt model from the  
 329 software interface, a gain bar is shown that reports the improvement in penalized average  
 330 log likelihood compared to a null model.

331 ----- **end of Box 1** -----

332

333

### 334 MAXENT'S LOGISTIC OUTPUT

335 MaxEnt (from version 3 onwards) gives a logistic output as its default. It is an attempt to get as  
 336 close as we can to an estimate of the probability that the species is present, given the environment,  
 337  $\Pr(y = 1 | \mathbf{z})$ . This is a post-transformation of the MaxEnt raw output that makes certain assumptions  
 338 about prevalence and sampling effort (Box 2 and Online Appendix 3). These two output types of  
 339 MaxEnt (raw and logistic) are monotonically related, so if the purpose of a study is to rank sites  
 340 according to suitability, it does not matter which type is used – both will yield identical ranking and  
 341 hence identical rank-based measures (e.g., AUC values). MaxEnt's logistic transformation is not a  
 342 commonly used statistical procedure, so here we explain the background and the issues.

343

344 From Eqn. 1 we see that a simple approach to estimate  $\Pr(y=1|\mathbf{z})$  would be to simply multiply  $e^{\eta(\mathbf{z})}$   
 345 by a constant that estimates prevalence; this approach has the disadvantage that  $e^{\eta(\mathbf{z})}$  can be  
 346 arbitrarily large, which implies that we may get an estimate of  $\Pr(y = 1 | \mathbf{z})$  that exceeds 1 (Keating &  
 347 Cherry, 2004; Ward, 2007b). Exponential models can be especially badly behaved when applied to  
 348 new data, for instance, when extrapolating to new environments. To avoid these problems, and to  
 349 side-step the non-identifiability of the species prevalence,  $\Pr(y=1)$ , MaxEnt's logistic output  
 350 transforms the model from an exponential family model (Eqn. 2) to a logistic model:



351

352

$$\Pr(y = 1|\mathbf{z}) = \tau e^{\eta(\mathbf{z})-r} / (1-\tau + \tau e^{\eta(\mathbf{z})-r}) \quad \dots \text{Equation 5}$$

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

----- Box 2: Consider the jaguar: reconciling logistic output and sampling effort -----

The jaguar (*Panthera onca*) and the collared peccary (*Pecari tajacu*) have very similar ranges in South and Central America, and MaxEnt models for the two species would therefore be similar using the default  $\tau$ . However, the jaguar is much rarer than the peccary, so how can the outputs be compared? The answer is that probability of presence is only defined relative to a given definition of presence/absence (i.e., the temporal and spatial scale of a sample; see Preamble). For instance, for a rare species like the jaguar a presence record is likely to derive from sampling over a longer time and/or larger area (e.g. using camera traps over months) than it would for the peccary, which is fairly common and easier to observe. Since with presence-only data there is usually no information on sampling effort, this elasticity in definition is largely conceptual – it explains how to think about the meaning of the probabilities across species, when  $\tau$  is 0.5. When  $\tau$  is 0.5 typical presence sites will have a logistic output near 0.5. This is reasonable as long as we can interpret logistic output as corresponding to a temporal and spatial scale of sampling that results in a 50% chance of the species being present in suitable areas. See Online Appendix 3 for more information.

Alternatively, if the value of  $\tau$  is available for a given level of sampling effort, it could be used instead of the default, and then the predictions for the two species would be directly comparable. Tau measures a form of rarity (Rabinowitz *et al.* 1986). The jaguar has very low local abundance even in suitable areas within its range, so a very small value  $\tau$  is appropriate for all but the most intensive sampling schemes. The estimate of  $\tau$  could come from expert knowledge or targeted surveys. While  $\tau$  is determined by prevalence, and vice versa,  $\tau$  is arguably more ecologically intuitive, as it is more a property of the species while prevalence strongly depends on the choice of study area.

----- end of Box 2 -----

395

## IMPLICATIONS FOR MODELLING

396

397

398

399

400

401

402

These properties of the MaxEnt model have several implications for how it should be used.

MaxEnt relies on an unbiased sample (as do all species modelling methods), so efforts in collecting a comprehensive set of presence records (cleaned for duplicates and errors) and dealing with biases are critical (Newbold, 2010). Methods are implemented for dealing with biased species data (see case study 1, and Dudík *et al.*, 2006; Phillips *et al.*, 2009; Elith *et al.*, 2010 in press). The main alternatives are to provide background data with similar biases to those in the presence data (e.g. by using sites surveyed for other species in the same biological group), or to use a bias grid that

403 indicates the biases in the survey data (see tutorial provided with MaxEnt for an example). All the  
404 values in this grid should be positive (or specified as no data), and should be scaled to represent  
405 relative survey effort across the landscape  $L$ . There is one additional important consideration. If the  
406 covariate grids are unprojected (i.e. latitude and longitude in degrees, for instance WorldClim data  
407 - <http://www.worldclim.org/>), any region covering a non-trivial range in latitude (say, more than  
408 200km, especially away from the equator) will have grid cells of varying area. For instance, in  
409 Australia cells in the north are approximately 1.3 times the area of cells in the south. MaxEnt  
410 randomly samples cells, implicitly assuming equal area cells. Solutions are to project the grids to an  
411 equal area projection, create a grid showing the variations in cell area, that can then be used as a  
412 bias grid, or create your own background sample with appropriate sampling weights (case study 1).  
413

414 The MaxEnt solution is affected by the landscape (region) used for the background sample, as  
415 demonstrated by VanDerWal *et al.* (2009). Conceptually, that landscape should include the full  
416 environmental range of the species, and exclude areas that definitely have not been searched  
417 (unless the reason for no searching is that there is unambiguous knowledge that the species does  
418 not occur there). A local endemic that is, for instance, likely to be geographically restricted due to  
419 barriers to dispersal, should be modelled with background selected from areas into which it might  
420 have dispersed. Cleared areas that would not be surveyed because there is no remaining habitat for  
421 the species should be excluded. Excluding areas from the background sample can be achieved  
422 through use of masks, as explained in the online tutorial for MaxEnt (and see Table 2). Predictions  
423 can still be made to excluded areas, if required, by using the projection facilities. We will discuss  
424 some caveats to these general concepts for background selection in the first case study.  
425

426 MaxEnt includes a range of feature types, and subsets of these can be used to simplify the solution.  
427 By default, the program restricts the model to simple features if few samples are available (linear is  
428 always used; quadratic with at least 10 samples; hinge with at least 15; threshold and product with  
429 at least 80) because – as for any modelling method – few samples provide limited information for  
430 determining the relationships between the species and its environment (Barry & Elith, 2006;  
431 Pearson *et al.*, 2007). In such cases, it is also a good idea to first reduce the candidate predictor set  
432 using ecological understanding of the species (Elith & Leathwick, 2009b). Hinge features tend to  
433 make linear and threshold features redundant, and one way to form a model with relatively smooth  
434 fitted functions, more like a GAM, is to use only hinge features (e.g. Elith *et al.*, 2010 in press, and  
435 case study 1). Excluding product features creates an additive model that is easier to interpret,  
436 though less able to model complex interactions.  
437

438 MaxEnt has an inbuilt method for regularization (L1-regularization) that is reliable and known to  
439 perform well (Hastie *et al.*, 2009). It implicitly deals with feature selection (relegating some  
440 coefficients to zero) and is unlikely to be improved - and more likely, degraded - by procedures that  
441 use other modelling methods to pre-select variables (e.g. Wollan *et al.*, 2008). In particular, it is  
442 more stable in the face of correlated variables than stepwise regression, so there is less need to  
443 remove correlated variables (unless some of them are known to be ecologically irrelevant), or  
444 preprocess covariates by using PCA and selecting a few dominant axes. Note, though, that since  
445 there are often many variables available, some expert pre-selection of a candidate set is often a  
446 good idea; Elith and Leathwick 2009b. Selecting proximal variables is likely to be particularly  
447 important when models are to be used in different regions or climates. If smoother models are  
448 required, regularization parameters can be increased by the user (e.g., see Elith *et al.* 2010 in  
449 press).  
450

451 If comparing models for different species some care is needed in use of the logistic outputs because  
452 probability of presence is only defined relative to a given level of sampling effort, which as a default  
453 is assumed to be one that results in a 50% chance of observing the species in suitable areas (Box 2).  
454 The implied sampling effort therefore depends on the species. This presents some challenges for  
455 cross-species comparisons of habitable areas, but these are a direct result of using presence-only  
456 data, and is not a unique problem to MaxEnt. Some users may in fact see the species-specific scaling

457 as an opportunity, since the literature on favourability functions (e.g., Real *et al.*, 2006) claims that  
 458 probability of presence is itself hard to work with.

## 459 USING MAXENT

### 460 CASE STUDY 1: MODELLING CURRENT AND FUTURE DISTRIBUTIONS OF A PLANT

461 This analysis predicts the current distribution of *Banksia prionotes*, then uses the model to identify  
 462 where suitable environments for the species are likely to occur under climate change. In it we  
 463 highlight the importance of choice of landscape and dealing with survey bias, debiasing background  
 464 samples from unprojected grids, use of a reduced set of feature types for a smoother model, and  
 465 tools for assessing the environments in new times or places.

466 *B. prionotes* is a woody shrub to small tree native to south west Western Australia (WA). It is widely  
 467 distributed across its range, and shows a preference for deep sandy soils. Often a dominant plant in  
 468 scrubland and low woodlands, it is an important nectar source for honeyeaters, and an outstanding  
 469 ornamental species for cut flowers.

470  
 471 Methods: Here we use species data from the Banksia Atlas (Taylor & Hopper, 1988; Yates *et al.*,  
 472 2010), with 361 records for *B. prionotes* from the 4631 sites across the South West Australia  
 473 Floristic Region (SWAFR) that were surveyed for *Banksia* and for which we had complete  
 474 environmental data. The atlas is the result of a community science project, and records could either  
 475 be interpreted as presence-only or presence-absence data, depending on what assumptions are  
 476 made about the search patterns of contributors. Here we treat them as presence-only data, but use  
 477 the full set of locations as one “background” treatment. To demonstrate the effect of this choice, two  
 478 alternative backgrounds (i.e. landscape definitions) were evaluated: a sample of 10000 sites within  
 479 the SWAFR (Yates *et al.* 2010, and Figure 2), and a sample of 20000 sites across the whole of  
 480 Australia. The larger number of sites across Australia was used to ensure good representation of all  
 481 environments, based on previous tests of the effects of background sample size on model structure  
 482 for these predictors (Elith unpubl.). Because the covariate data for this study are unprojected, these  
 483 samples were weighted according to cell area (see methods in Online Appendix 4) but otherwise  
 484 random.

485  
 486 Using random sites within the floristic region implies that the presence records are a random  
 487 sample from all locations where the species is present in the region which is unlikely because  
 488 records were from extant vegetation patches in likely suitable environments (the region has been  
 489 extensively cleared for agriculture, and some of the more inland areas are too arid for many *Banksia*  
 490 species). Using random sites across Australia implies the species could have dispersed anywhere  
 491 across the continent, and the whole continent considered available for sampling. This is  
 492 questionable because the desert areas to the north and east of the inhabited area are likely barriers  
 493 to dispersal. We will come back to implications of this later.

494  
 495 Yates *et al.* (2010) identified important climatic drivers for plants of southwest Western Australia.  
 496 We base our candidate set of predictors on their study, but use a different data source so we can  
 497 train and predict over the whole of Australia. Described in online Appendix 4, our covariates (all  
 498 unprojected, at 0.01 degree or ~ 1km grid resolution) included five climate variables: isothermality  
 499 (ISOTHERM), mean temperature of the wettest quarter (TEMPWETQ), mean temperature of the  
 500 warmest quarter (TEMPWARMQ), annual precipitation (RAIN) and precipitation of the driest  
 501 quarter (RAINDRYQ), and an estimate of the solum plant-available water holding capacity (solwhc).  
 502 We present this as a demonstration study only, and recognize that, for rigorous application in this  
 503 region, better soils data and predictors representing land transformation are needed for more  
 504 precise predictions (Yates *et al.* 2010). The future environment was represented by changes  
 505 predicted under the A1FI scenario for 2070 estimated over the ensemble of 23 GCMs in IPCC AR4  
 506 (Solomon *et al.* 2007); the solwhc was assumed to remain as it is now.

507 Models were fitted and projected to both current and future climates (Figure 3) using only hinge  
508 features, with default regularization parameters (see Appendix 5 for model details, and for a  
509 comparison with models fitted with all feature types). We fitted all models on the full data sets but  
510 also used 10-fold cross validation to estimate errors around fitted functions and predictive  
511 performance on held-out data. The latter is a good test for each model but – given the different  
512 backgrounds – not comparable across models. Note also that the AUC in this case is calculated on  
513 presence vs background data (Phillips *et al.* 2006). We also divided the atlas data into training and  
514 testing sets for a manual 5-fold cross-validation, testing each model on identical withheld data via  
515 two test statistics (area under the receiver operating characteristic curve (AUC), and correlation,  
516 COR; details in online Appendix 4). Example code for running such analyses are available online  
517 (Appendix 4).

518  
519 Results: Atlas background (model 1) produced a mapped distribution in the inhabited region with  
520 more of an eastward emphasis compared with other background treatments (Figure 3). The  
521 coastward (westerly) bias in the distribution of survey sites (Figure 2) affected the distributions  
522 predicted by models 2 and 3 (random background across SWAFR or Australia) but was factored out  
523 by using atlas background (model 1). The more easterly distribution is more consistent with the  
524 known ecology of the species, and with the observed distribution (Taylor and Hopper 1988).  
525 Variable importance varies with data set, with TEMPWETQ being much more prominent when  
526 using an all-Australia background than when restricted to the south-west. Similarly, shapes of fitted  
527 functions vary across data sets (Appendix 5). This is to be expected, because each data set implies a  
528 different modelling question (e.g. the all-Australia background asks: why is this species only in  
529 environments occurring in the southwest?).

530  
531 An increasing number of SDM applications involve prediction to new environments (e.g. to new  
532 places or times; Elith & Leathwick, 2009a). These are contentious applications, making strong  
533 assumptions (Dormann, 2007) and usually requiring prediction to environments not sampled by  
534 the training data. MaxEnt has been extended to include new capabilities to inform users about  
535 predicting to novel environments (Elith *et al.*, 2010 in press). MaxEnt already provides mapped  
536 information on the effect of model "clamping" – i.e. the process by which features are constrained to  
537 remain within the range of values in the training data. This identifies locations where predictions  
538 are uncertain due to the method of extrapolation, by showing where clamping substantially affects  
539 the predicted value. We feel that extreme care should be taken whenever extrapolating outside the  
540 training, so new calculations ("MESS maps", i.e. multivariate environmental similarity surfaces)  
541 display differences between the training and prediction environments (Figure 3). In this case they  
542 show that, compared with environments at the atlas sites, the northern parts of the SWAFR will  
543 experience novel climates in 2070 (Figure 3 model #1). Models based on random background  
544 across SWAFR or the continent (models 2 and 3) require less extrapolation (because wider  
545 sampling of background points brings with it wider sampling of environments) but, given the  
546 problems with the realism of these treatments, we do not view the result as a necessary advantage  
547 for future predictions.

548  
549 Online Appendices 5 and 6 include further information on how these models predict across the  
550 continent, for both current and future climates. They provide interesting insights into model  
551 variation across scales, regions, and datasets, and emphasize the importance of choice of  
552 background (see commentary, Appendix 5). In particular, it is interesting that model 3 restricts  
553 predictions to the correct general area, and has the highest 10-fold cross-validated AUC (Table 3),  
554 yet has the poorest ecological justification for its choice of background and is least likely to be  
555 useful for managing the species locally. The advantage of limiting background to local, reachable  
556 areas (models 1 and 2) is that contrasts between occupied and unoccupied environments in the  
557 local area are the model focus, and – particularly with fine-scale environmental data –  
558 differentiation useful at the management scale might be achievable. It is also likely to be the most  
559 ecologically realistic choice for many locally restricted species. On the other hand, if models are to  
560 be projected well outside the local geographic area, use of local backgrounds brings with it the

561 penalty that prediction to other areas is likely to involve considerable extrapolation. Some trade-off  
562 is clearly required.

## 563 CASE STUDY 2: MODELLING THE DISTRIBUTIONS OF FISH IN RIVERS.

564 This analysis predicts the current distribution of *Gadopsis bispinosus*, the two-spined blackfish, in  
565 rivers of south-eastern Australia. In the preamble we make a case that with presence and  
566 background data, we can model the same quantity as with presence-absence data, up to the  
567 constant  $\Pr(y=1)$ . One implication of that is that we should be able to use the same types of data,  
568 including fine-scale, detailed information, to model ecological relationships – i.e. we need not be  
569 restricted to coarse grid cells and basic climate variables. Here we use detailed ecological  
570 information at the river segment scale to model the distribution of a native fish species. To our  
571 knowledge it is the first example using MaxEnt with vector (river segment) data.

572  
573 *G. bispinosus* is a native freshwater fish endemic to south-eastern Australia. It occurs in cool, clear  
574 upland or montane streams with abundant in-stream cover. It is most common in medium to large  
575 streams that are deep enough for reduced stream velocities, and in forested catchments with  
576 relatively small sediment inputs (Lintermans 2000).

### 577 Methods:

578 The species data are from surveys (described further in online Appendix 7) of the inland-draining  
579 rivers of northwest Victoria, Australia. In this area there are ten major river systems grouped into  
580 four regions that start in hilly to mountainous terrain and drain northwards. *G. bispinosus* was  
581 recorded at 255 sites. We use covariate data from the 255 capture sites as our sample of  $L_1$  and a  
582 random sample of 10000 of the ~240000 river segments for our sample of  $L$ , the background data.  
583  
584

585 The candidate predictor set comprised 20 variables summarizing information across three  
586 hierarchically nested spatial scales (segment, immediate watershed and entire upstream catchment  
587 area) and also downstream to the large river system draining to the ocean. The environmental  
588 variables estimate climate, river slope, riparian vegetation and catchment characteristics (Online  
589 Appendix Table S7.1). River system was also included to quantify spatial variation in land  
590 characteristics and disturbances not covered by the environmental predictor set.

591  
592 These segment-based (non-gridded) data are modelled using the SWD (samples-with-data) format  
593 in MaxEnt – this involves presenting spreadsheet-like summaries of environments at both presence  
594 and background sites. All environmental variables were continuous except the categorical river  
595 system covariate. Default settings for features and regularization were used for model training, and  
596 10-fold cross-validation used to obtain out-of-sample estimates of predictive performance and  
597 estimates of uncertainty around fitted functions. For mapping, the model was projected to a  
598 selected area in the Goulburn-Broken catchment. Technically, this was achieved by projecting to  
599 SWD format data, then linking the predictions to the relevant river segments in a GIS. Online  
600 Appendix 8 includes data and code for replicating this case study, including information on how to  
601 run MaxEnt from batch files.

602  
603 Results: Consistent with ecological knowledge about the species, the model predicts *G. bispinosus*  
604 will most frequently occur in the larger streams of montane areas (Figure 4). These locations are  
605 identified as those whose upstream catchments have relatively more precipitation in the warmest  
606 quarter and steeper maximum stream slopes. Amongst these, emphasis on segments with warmer  
607 summer maximum temperatures served to exclude the higher elevation cold streams (Figure 5).  
608 Jackknife tests of variable importance help to identify those with important individual effects; the  
609 three most important single predictors were the summed length of all upstream links  
610 (TOTLENGTH\_UCA), the upstream maximum slope (US\_MAXSLOPE) and the amount of riparian  
611 tree cover upstream (UC\_RIP\_TRECOV); and the predictor with the most information not present in  
612 the other variables is the segment-based maximum temperature of the warmest month

613 (MAXWARMP\_TEMP). Many predictors had small to minimal impacts in the final model. The model  
614 shows strong discrimination on held out data, with a cross-validated AUC of 0.97.

615  
616 Extensions / alternatives: Since records on one river system might share a more similar  
617 environment than those on different systems, an alternative approach to cross-validation would be  
618 to test the predictions iteratively on held-out rivers. We chose not to do it in this case, because  
619 presence records were concentrated in relatively few river systems, so the training sets would be  
620 substantially reduced, and the test sets, relatively few.

## 621 CONCLUSIONS

622 Here we have described MaxEnt from a statistical viewpoint, showing that the model minimizes the  
623 relative entropy between two probability densities defined in feature space. An understanding of  
624 the model leads naturally to recommendations for implementation, and ours included the  
625 importance of providing appropriate background samples, of dealing with sample biases, and of  
626 tuning the model – through feature type selection and regularization settings - to suit the data and  
627 application. Presence-only data are a valuable resource and potentially can be used to model the  
628 same ecological relationships as with presence-absence data, provided that biases can be dealt with  
629 and except for the non-identifiability of prevalence.

630  
631 MaxEnt is regularly updated, usually to include new capabilities to suit the expanding applications,  
632 and also sometimes to change the program defaults to those most often used in practice. Recent  
633 new capabilities include the cross-validation and MESS maps (i.e. estimates of how the  
634 environmental space in predicted times and places compares with that of the training data)  
635 demonstrated in case study 1. In addition, new clickable maps allow users to interrogate  
636 predictions spatially, providing information for any grid cell on the components of the prediction  
637 (i.e. what contributes to its particular value) and where the environmental conditions “sit” on the  
638 fitted functions. Maps of limiting factors show the variable most influencing the prediction for every  
639 grid cell. For further details see Elith *et al.* (2010 in press) and the most recent online tutorial  
640 (<http://www.cs.princeton.edu/~schapire/maxent/>). SDMs can provide useful information for  
641 exploring and predicting species distributions, and we are keen to see their continued development  
642 and use for learning about and conserving the world's biodiversity.

## 643 ACKNOWLEDGEMENTS

644 JE was supported by an Australian Research Council grant, FT0991640 and by an early consultancy  
645 that raised the question of how to explain MaxEnt to end-users (Jeff Tranter, Environmental  
646 Resources Information Network, Canberra, AUS). TH was partially supported by grant DMS-  
647 1007719 from the U.S. National Science Foundation. Simon Ferrier, John Baumgartner and Tord  
648 Snäll provided useful feedback on ideas and/or the manuscript. Robert Hijmans provided the  
649 method for taking samples proportional to area. Thanks to the three reviewers including Mark  
650 Robertson and Janet Franklin for generous and constructive comments and good ideas.

**Table 1 – Examples of published studies using MaxEnt, showing variation in purpose, scale and organism**

Primary purpose	Scales	Organisms	Refs
Predict current distributions as input for conservation planning, risk assessments or IUCN listing, or new surveys	Andes	Humming-birds	Tinoco <i>et al.</i> (2009)
	Global	Stony corals seamounts	Tittensor <i>et al.</i> (2009)
Understand environmental correlates of species occurrences, groups of species, or other	Norway	Macrofungi	Wollan <i>et al.</i> (2008)
	Portugal	European wildcat	Monterroso <i>et al.</i> (2009)
Predict potential distributions for invasive species, or explore expanding distributions	New Zealand	Ants	Ward (2007a)
	China	Nematode	Wang <i>et al.</i> (2007)
Predict species richness or diversity	California	Amphibians and reptiles	Graham and Hijmans (2006)
	Brazil	Myrtaceae 19 species	Murray-Smith <i>et al.</i> (2009)
Predict current distributions for understanding morphological / genetic diversity (“phylogeography”, “phyloclimatic studies”), endemism and evolutionary niche dynamics	Global	Seaweeds	Verbruggen <i>et al.</i> (2009)
	Andes	Birds	Young <i>et al.</i> (2009)
	Madagascar	Bats	Lamb <i>et al.</i> (2008)
Hindcast distributions to understand patterns of endemism, vicariance etc	NW Europe	Pond snails	Cordelier and Pfenninger (2009)
	Brazilian coast	Forests	Carnaval and Moritz (2008)
Forecast distributions to understand changes with climate change / land transformation; includes retrospective studies	Mediterr’n + surrounds	Cyclamen	Yesson and Culham (2006)
	Regional W. Australia	Banksia	Yates <i>et al.</i> (2010)
	Canada	Butterflies	Kharouba <i>et al.</i> (2009)
Test model performance against other methods	Patagonia	Insects	Tognelli <i>et al.</i> (2009)
	Local region in California	Rare plants	Williams <i>et al.</i> (2009)
	Regional to national	Many species	Elith <i>et al.</i> (2006)

**Table 2:** Terminology used in this paper

Item / concept	Definition	Notation
Background	A sample of points from the landscape	
Entropy	A measure of disperdness. Previous papers <sup>1</sup> described the model as maximizing entropy in geographic space; this paper focuses on minimizing relative entropy in covariate space.	
Features	An expanded set of transformations of the original covariates	
Mask	A gridded layer of 1 / no data used to indicate areas to be included in background sampling (=1) and those to be excluded (=no data). To be included as a predictor. For projecting to the whole region, a grid called mask, but containing any values – say, 1 across the whole region of interest – should be supplied along with all other covariate grids.	
MESS map	Multivariate Environmental Similarity Surface –measures the similarity of any given point to a reference set of points, with respect to the chosen predictor variables. It reports the closeness of the point to the distribution of reference points, gives negative values for dissimilar points and maps these values across the whole prediction region (Elith <i>et al.</i> 2010 in press)	
Prevalence is not identifiable	Prevalence cannot be exactly determined from presence-only data in isolation, regardless of the sample size. This is a fundamental limitation of presence-only data.	
Probability density functions	Describe the relative likelihood of random variables over their range; can be univariate or multivariate.	
Regularization (tuning) parameters	Regularization refers to smoothing the model, making it more regular, so as to avoid fitting too complex a model. In MaxEnt the regularization parameters can be changed if required.	$\beta$ in previous papers <sup>1</sup> , $\lambda$ in this paper
Sampling bias	Some areas in the landscape are sampled more intensively than others. Usually occurs in geographic space but could be environmentally based.	$s(z)$
Weights or coefficients	These are the parameters of the model that weight the contribution of each feature.	$\lambda$ in previous papers <sup>1</sup> , $\beta$ in this paper

<sup>1</sup> Phillips *et al.* (2006), Phillips & Dudík (2008)



**Table 3. Variable importance and evaluation statistics for case study 1.** Variable names and abbreviations for evaluation statistics are consistent with the text.

Model (background)	Variable importance						AUC (10fold CV but varying data sets)	AUC; COR (5fold CV on atlas data)
	RAIN DRYQ	RAIN	TEMP- WARMQ	TEMP- WETQ	ISO- THERM	SOL- PWHC		
1 (atlas)	57.9	30.7	7.9	0.4	1.1	2.0	0.92	0.96; 0.62
2 (southwest)	45.3	35.4	4.7	3.4	9.9	1.4	0.90	0.93; 0.52
3 (Australia)	19.7	17.7	5.3	54.0	3.0	0.3	0.99	0.91; 0.45

## FIGURES

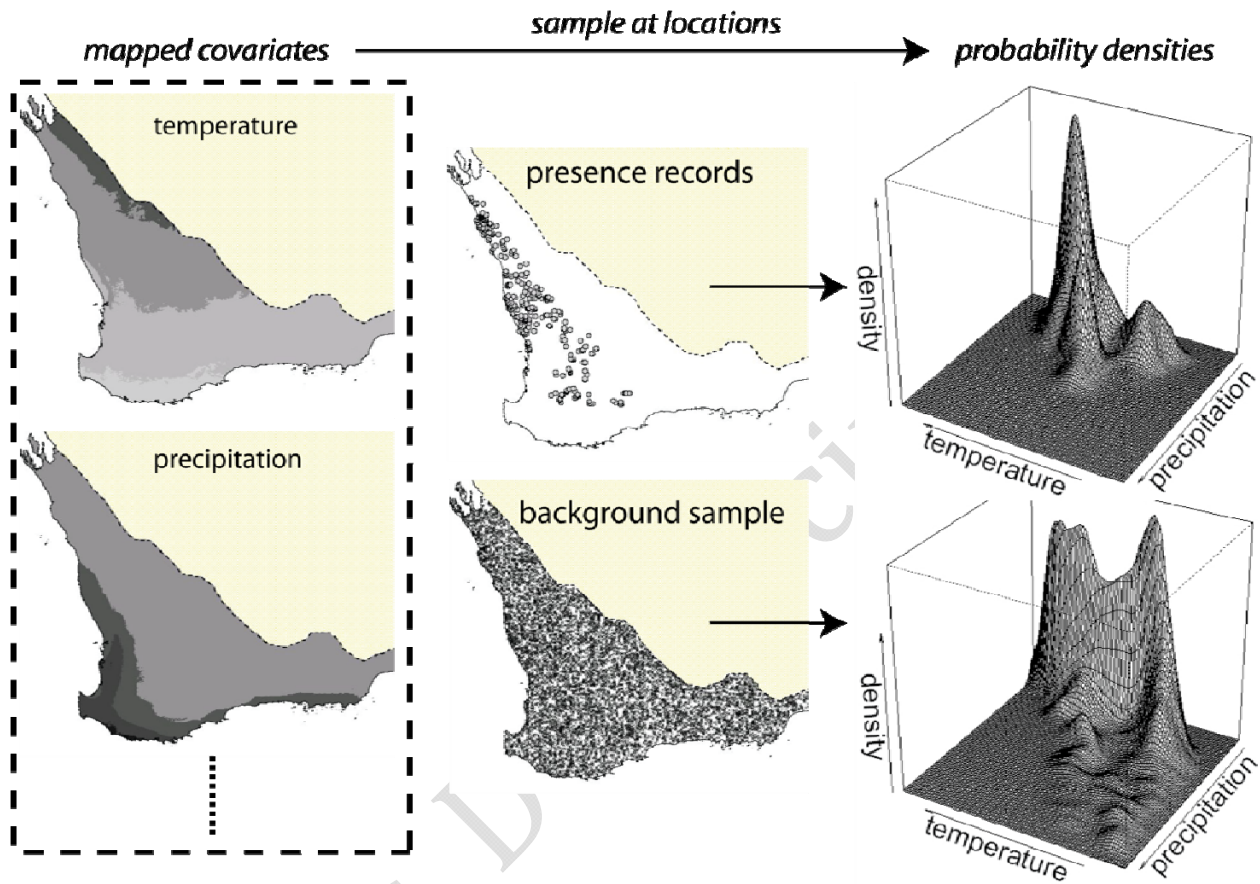


Figure 1 – A diagrammatic representation of the probability densities relevant to our statistical explanation, using data presented in case study 1. The maps on the left are two example mapped covariates (temperature and precipitation). In the centre are the locations of the presence and background samples. The density estimates on the right are not in geographic (map) space, but show the distributions of values in covariate space for the presence (top right) and background (bottom right) samples. These could represent the densities  $f_1(\mathbf{z})$  and  $f(\mathbf{z})$  for a simple model with linear features.

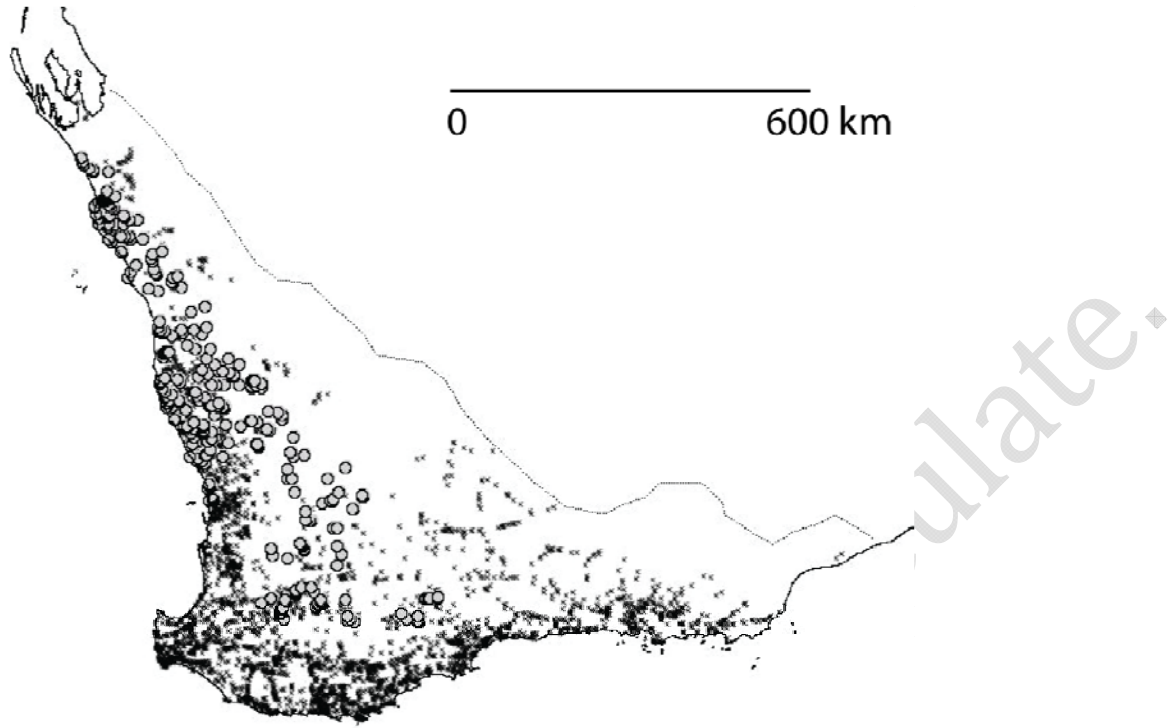


Figure 2: All banksia atlas sites (black) with occurrences of *B. prionotes* in grey circles.

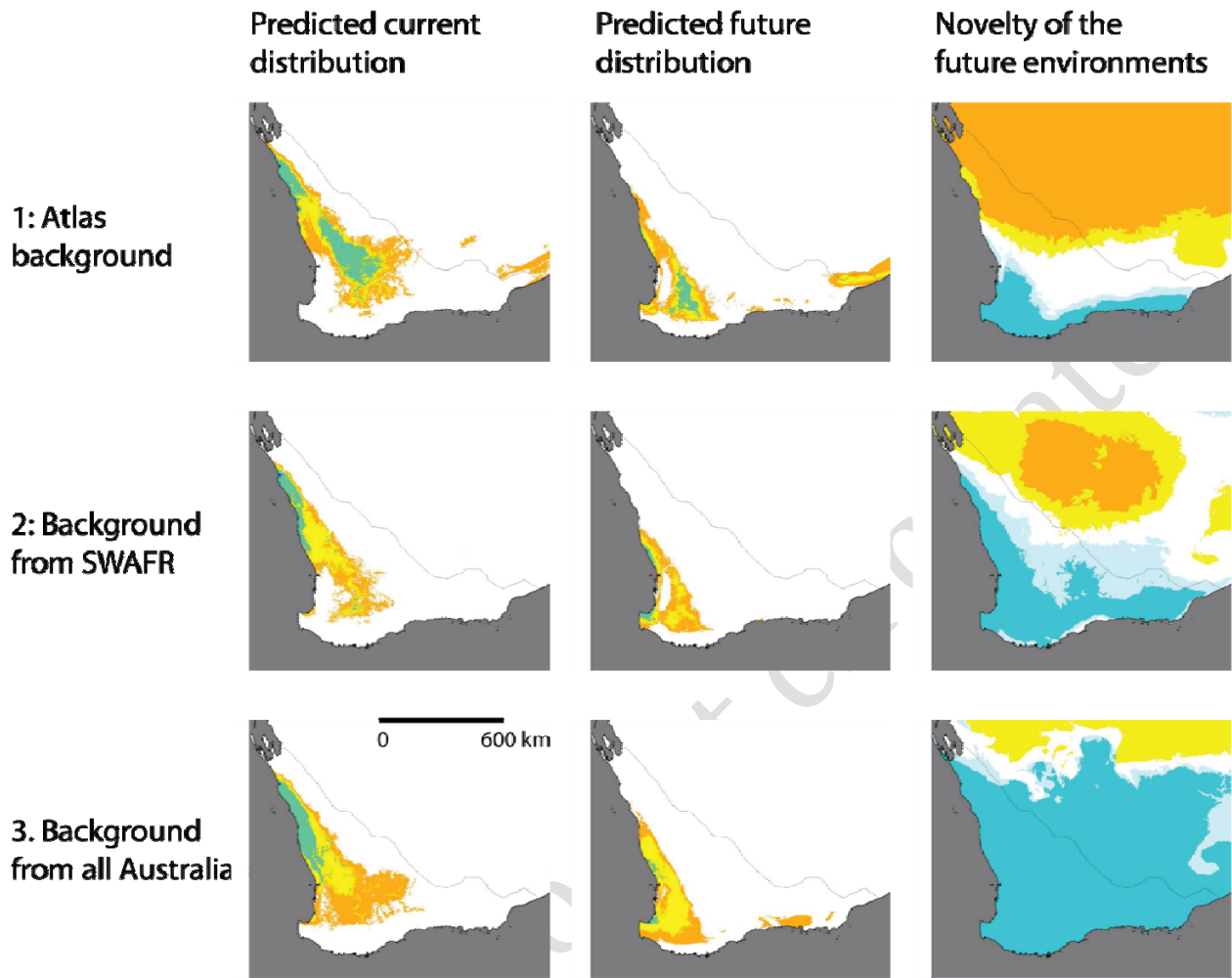


Figure 3. Model results for case study 1, showing for the three data sets (in rows): predicted current and future distributions, and extent of extrapolation compared with the training data. Predicted distributions are logistic outputs, from low values (white, 0 to 0.2) through orange, yellow, green to blue (0.8 to 1.0). For extrapolation maps, warm colours indicate extrapolation is occurring, with orange the most extreme. Grey indicates the ocean.

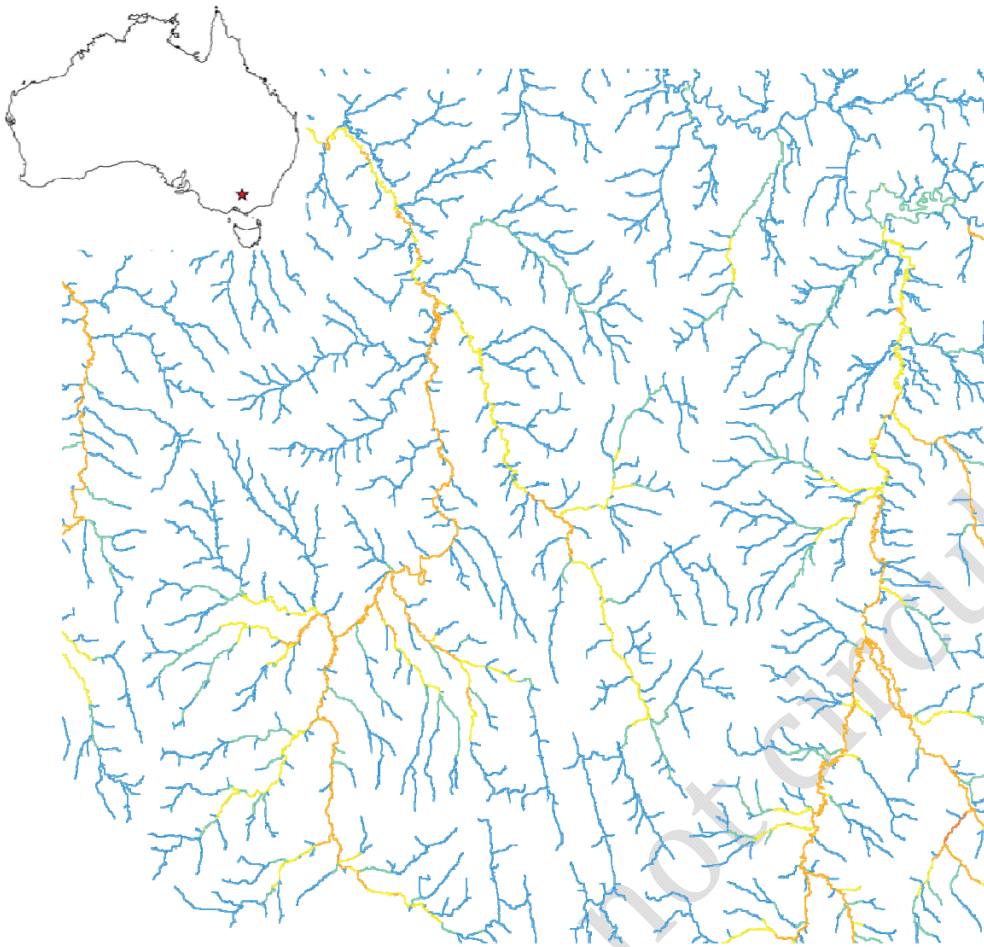


Figure 4 – Predicted distribution of *Gadopsis bispinosus*, showing logistic output predictions from MaxEnt. Legend: predictions in equal intervals from 0 to 1, from blue (low) through green – yellow –orange (high). Scale: east to west the rivers map spans 45km. The star on the inset shows location.

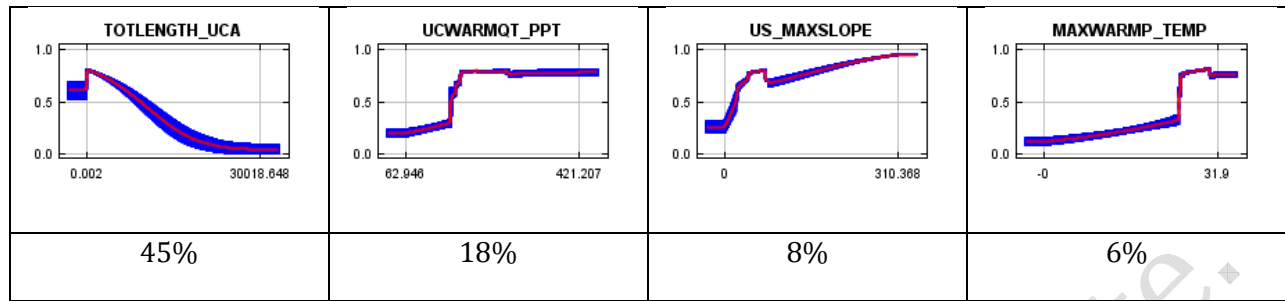


Figure 5: Partial dependence plots showing the marginal response of *Gadopsis bispinosus* to the four most important variables (i.e., for constant values of the other variables), with variable importance below each graph. The y axes indicates logistic output.

In press. Do not circulate.

## REFERENCES

- Akaike, H. (1974) A new look at statistical model identification. *IEEE Transactions on Automatic Control*, **AU-19**, 716-722.
- Austin, M.P. (2002) Spatial prediction of species distribution: an interface between ecological theory and statistical modelling. *Ecological Modelling*, **157**, 101-118.
- Barry, S.C. & Elith, J. (2006) Error and uncertainty in habitat models. *Journal of Applied Ecology*, **43**, 413-423.
- Carnaval, A.C. & Moritz, C. (2008) Historical climate modelling predicts patterns of current biodiversity in the Brazilian Atlantic forest. *Journal of Biogeography*, **35**, 1187-1201.
- Chefaoui, R.M. & Lobo, J.M. (2007) Assessing the effects of pseudo-absences on predictive distribution model performance. *Ecological Modelling*, **210**, 478-486.
- Cordellier, M. & Pfenninger, M. (2009) Inferring the past to predict the future: climate modelling predictions and phylogeography for the freshwater gastropod *Radix balthica* (Pulmonata, Basommatophora). *Molecular Ecology*, **18**, 534-544.
- Della Pietra, S., Della Pietra, V. & Lafferty, J. (1997) Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **19**, 1-13.
- Dormann, C.F. (2007) Promising the future? Global change projections of species distributions. *Basic and Applied Ecology*, **8** 387-397.
- Dudík, M. & Phillips, S.J. (2008) Generative and discriminative learning with unknown labeling bias. *Advances in Neural Information Processing Systems*, **21**
- Dudík, M., Schapire, R.E. & Phillips, S.J. (2006) Correcting sample selection bias in maximum entropy density estimation. *Advances in Neural Information Processing Systems 18: Proceedings of the 2005 Conference* pp. 323-330. Cambridge, MA.
- Elith, J. & Leathwick, J.R. (2009a) Species distribution models: ecological explanation and prediction across space and time. *Annual Review of Ecology, Evolution and Systematics*, **40**, 677-697.
- Elith, J. & Leathwick, J.R. (2009b) Conservation prioritization using species distribution models. Chapter 6 in: *Spatial conservation prioritization: quantitative methods and computational tools* (ed. by A. Moilanen, K.A. Wilson & H.P. Possingham). Oxford University Press.
- Elith, J., Kearney, M. & Phillips, S.J. (2010 in press) The art of modelling range-shifting species *Methods in Ecology and Evolution*,
- 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.S., 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.
- Franklin, J. (2009) *Mapping Species Distributions: Spatial Inference and Prediction*. Cambridge University Press, Cambridge, UK.
- Graham, C.H. & Hijmans, R.J. (2006) A comparison of methods for mapping species ranges and species richness. *Global Ecology & Biogeography*, **15**, 578.
- Hastie, T., Tibshirani, R. & Friedman, J.H. (2009) *The elements of statistical learning: data mining, inference, and prediction, second edition*, 2nd edn. Springer-Verlag, New York.
- Hirzel, A.H. & Le Lay, G. (2008) Habitat suitability modelling and niche theory. *Journal of Applied Ecology*, **45**, 1372-1381.
- Jiménez-Valverde, A., Lobo, J.M. & Hortal, J. (2008) Not as good as they seem: the importance of concepts in species distribution modelling. *Diversity and Distributions*, **14**, 885-890.
- Keating, K.A. & Cherry, S. (2004) Use and interpretation of logistic regression in habitat selection studies. *Journal of Wildlife Management*, **68**, 774-789.
- Kharouba, H.M., Algar, A.C. & Kerr, J.T. (2009) Historically calibrated predictions of butterfly species' range shift using global change as a pseudo-experiment. *Ecology*, **90**, 2213-2222.
- Lamb, J.M., Ralph, T.M.C., Goodman, S.M., Bogdanowicz, W., Fahr, J., Gajewska, M., Bates, P.J.J., Eger, J., Benda, P. & Taylor, P.J. (2008) Phylogeography and predicted distribution of African-Arabian and Malagasy populations of giant mastiff bats, *Otomops* spp. (Chiroptera : Molossidae). *Acta Chiropterologica*, **10**, 21-40.

- Leathwick, J.R. (1998) Are New Zealand's *Nothofagus* species in equilibrium with their environment? *Journal of Vegetation Science*, **9**, 719-732.
- Lintermans, M. (2000) The Status of Fish in the Australian Capital Territory: a Review of Current Knowledge and Management Requirements. Technical Report No. 15, Environment ACT, Canberra.
- Lobo, J.M., Jiménez-Valverde, A. & Hortal, J. (2010) The uncertain nature of absences and their importance in species distribution modelling. *Ecography*, **33**, 103-114.
- MacKenzie, D.I. (2005) Was it there? Dealing with imperfect detection for species presence/absence data. *Australia and New Zealand Journal of Statistics*, **47**, 65-74.
- MacKenzie, D.I. & Royle, J.A. (2005) Designing efficient occupancy studies: general advice and tips on allocation of survey effort. *Journal of Applied Ecology*, **42**, 1105-1114.
- Monterroso, P., Brito, J.C., Ferreras, P. & Alves, P.C. (2009) Spatial ecology of the European wildcat in a Mediterranean ecosystem: dealing with small radio-tracking datasets in species conservation. *Journal of Zoology*, **279**, 27-35.
- Murray-Smith, C., Brummitt, N.A., Oliveira-Filho, A.T., Bachman, S., Moat, J., Lughadha, E.M.N. & Lucas, E.J. (2009) Plant Diversity Hotspots in the Atlantic Coastal Forests of Brazil. *Conservation Biology*, **23**, 151-163.
- Newbold, T. (2010) Applications and limitations of museum data for conservation and ecology, with particular attention to species distribution models. *Progress in Physical Geography*, **34**, 3-22.
- Pearce, J.L. & Boyce, M.S. (2006) Modelling distribution and abundance with presence-only data. *Journal of Applied Ecology*, **43**, 405-412.
- Pearson, R.G., Raxworthy, C.J., Nakamura, M. & Townsend Peterson, A. (2007) Predicting species distributions from small numbers of occurrence records: a test case using cryptic geckos in Madagascar. *Journal of Biogeography*, **34**, 102-117.
- Phillips, S.J. & Dudík, M. (2008) Modeling of species distributions with Maxent: new extensions and a comprehensive evaluation. *Ecography*, **31**, 161-175.
- Phillips, S.J., Anderson, R.P. & Schapire, R.E. (2006) Maximum entropy modeling of species geographic distributions. *Ecological Modelling*, **190**, 231-259.
- Phillips, S.J., Dudík, M., Elith, J., Graham, C.H., Lehmann, A., Leathwick, J. & Ferrier, S. (2009) Sample selection bias and presence-only distribution models: implications for background and pseudo-absence data. *Ecological Applications*, **19**, 181-197.
- Rabinowitz, D., Cairns, S. & Dillon, T. (1986) Seven forms of rarity and their frequency in the flora of the British Isles. *Conservation Biology: the Science of Scarcity and Diversity* (ed. by M.E. Soulé), pp. 182-204. Sinauer Associates, Sunderland, Massachusetts, USA.
- Real, R., Barbosa, A.M. & Vargas, J.M. (2006) Obtaining environmental favourability functions from logistic regression. *Environmental and Ecological Statistics*, **13**, 237-245.
- Schulman, L., Toivonen, T. & Ruokolainen, K. (2007) Analysing botanical collecting effort in Amazonia and correcting for it in species range estimation. *Journal of Biogeography*, **34**, 1388-1399.
- Soberón, J. & Nakamura, M. (2009) Niches and distributional areas: Concepts, methods, and assumptions. *Proceedings of the National Academy of Sciences*, **106**, 19644-19650.
- Soberón, J.M. & Peterson, A.T. (2005) Interpretation of models of fundamental ecological niches and species' distributional areas. *Biodiversity Informatics*, **2**, 1-10.
- Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K.D., Tignor, M. & Miller, H.L. (eds) (2007) *Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change*.
- Svenning, J.C. & Skov, F. (2004) Limited filling of the potential range in European tree species. *Ecology Letters*, **7**, 565-573.
- Taylor, A. & Hopper, S.D. (1988) The Banksia Atlas. AGPS, Canberra, ACT.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B*, **58**, 267-288.
- Tinoco, B.A., Astudillo, P.X., Latta, S.C. & Graham, C.H. (2009) Distribution, ecology and conservation of an endangered Andean hummingbird: the Violet-throated Metaltail (*Metallura baroni*). *Bird Conservation International*, **19**, 63-76.



- Tittensor, D.P., Baco, A.R., Brewin, P.E., Clark, M.R., Consalvey, M., Hall-Spencer, J., Rowden, A.A., Schlacher, T., Stocks, K.I. & Rogers, A.D. (2009) Predicting global habitat suitability for stony corals on seamounts. *Journal of Biogeography*, **36**, 1111-1128.
- Tognelli, M.F., Roig-Junent, S.A., Marvaldi, A.E., Flores, G.E. & Lobo, J.M. (2009) An evaluation of methods for modelling distribution of Patagonian insects. *Revista Chilena De Historia Natural*, **82**, 347-360.
- VanDerWal, J., Shoo, L.P., Graham, C. & Williams, S.E. (2009) Selecting pseudo-absence data for presence-only distribution modeling: how far should you stray from what you know? *Ecological Modelling*, **220**, 589-594.
- Verbruggen, H., Tyberghein, L., Pauly, K., Vlaeminck, C., Van Nieuwenhuyze, K., Kooistra, W., Leliaert, F. & De Clerck, O. (2009) Macroecology meets macroevolution: evolutionary niche dynamics in the seaweed *Halimeda*. *Global Ecology and Biogeography*, **18**, 393-405.
- Ward, D. (2007a) Modelling the potential geographic distribution of invasive ant species in New Zealand. *Biological Invasions*, **9**, 723-735.
- Ward, G. (2007b) *Statistics in ecological modeling; presence-only data and boosted mars*. Stanford University, Palo Alto.
- Williams, J.N., Seo, C.W., Thorne, J., Nelson, J.K., Erwin, S., O'Brien, J.M. & Schwartz, M.W. (2009) Using species distribution models to predict new occurrences for rare plants. *Diversity and Distributions*, **15**, 565-576.
- Wintle, B.A., McCarthy, M.A., Parris, K.M. & Burgman, M.A. (2004) Precision and bias of methods for estimating point survey detection probabilities. *Ecological Applications*, **14**, 703-712.
- Wollan, A.K., Bakkestuen, V., Kauserud, H., Gulden, G. & Halvorsen, R. (2008) Modelling and predicting fungal distribution patterns using herbarium data. *Journal of Biogeography*, **35**, 2298-2310.
- Yates, C., McNeill, A., Elith, J. & Midgley, G. (2010) Assessing the impacts of climate change and land transformation on *Banksia* in the South West Australian Floristic Region. *Diversity and Distributions*, **16**, 187-201.
- Yesson, C. & Culham, A. (2006) A phyloclimatic study of *Cyclamen*. *BMC Evolutionary Biology*, **6**, 72-95
- Young, B.F., Franke, I., Hernandez, P.A., Herzog, S.K., Paniagua, L., Tovar, C. & Valqui, T. (2009) Using spatial models to predict areas of endemism and gaps in the protection of Andean slope birds. *Auk*, **126**, 554-565.
- YunSheng, W., BingYan, X., FangHao, W., QiMing, X. & LiangYing, D. (2007) The potential geographic distribution of *Radopholus similis* in China. *Agricultural Sciences in China*, **6**, 1444-1449.
- Zadrozny, B. (2004) Learning and evaluating classifiers under sample selection bias. p114 in *Proceedings of the Twenty-First International Conference on Machine Learning*. Association of Machine Learning, New York, New York, USA.