

# Chapter 13 1

## Modelling Species' Distributions 2

Carsten F. Dormann 3

**Abstract** Species distribution models have become a commonplace exercise over 4  
the last 10 years, however, analyses vary due to different traditions, aims of applica- 5  
tions and statistical backgrounds. In this chapter, I lay out what I consider to be the 6  
most crucial steps in a species distribution analysis: data pre-processing and visual- 7  
isation, dimensional reduction (including collinearity), model formulation, model 8  
simplification, model type, assessment of model performance (incl. spatial autocor- 9  
relation) and model interpretation. For each step, the most relevant considerations are 10  
discussed, mainly illustrated with Generalised Linear Models and Boosted Regres- 11  
sion Trees as the two most contrasting methods. In the second section, I draw 12  
attention to the three most challenging problems in species distribution modelling: 13  
identifying (and incorporating into the model) the factors that limit a species range; 14  
separating the fundamental, realised and potential niche; and niche evolution. 15

### 13.1 Introduction 16

As more species types undergo rapid human-induced extinction, understanding why 17  
species occur where they do is becoming a highly relevant, pressing and potentially 18  
life-saving topic. Conservation actions, such as establishing protected site networks, 19  
adapting land use, providing stepping-stone habitats all require an idea of how the 20  
target species will respond. Furthermore, using organisms as a “bioassay technique”, 21  
i.e. indicators of environmental trends (such as climate change, air pollution, over- 22  
fishing) demands an intimate knowledge of the organism’s niche. Species distribution 23  
modelling (SDM) attempts to identify the probable causes of species whereabouts. 24  
We seek to delineate the realized niche of an organism based on its current distribu- 25  
tion with respect to the environment (see Elith and Leathwick 2009b for definitions 26  
and concepts, Guisan and Thuiller 2005; Kearney 2006; Soberón 2007). 27

---

C.F. Dormann

Department of Computational Landscape Ecology, Helmholtz Centre for Environmental  
Research – UFZ, Permoserstr.15, 04318 Leipzig, Germany  
e-mail: carsten.dormann@ufz.de

## 28 *Why Species Distribution Modelling?*

29 There are several fundamental challenges to this approach (e.g. first and foremost  
30 that it is correlative; see Vaughan and Ormerod 2005 and Dormann 2007b for a  
31 recent critique), and before jumping into the analysis, it is worth considering  
32 whether SDMs are actually useful and fit for the purpose of your specific problem.  
33 For example, at very small spatial scales, differences in environmental conditions  
34 may be too small to be of predictive value and biotic interactions (competition,  
35 predation) may be of crucial importance. In contrast, at the global scale, data  
36 become so coarse that we “only” model the climate niche and specific habitat  
37 requirements cannot be detected.

38 On the other hand, SDMs try to extract ecological information from a species  
39 occurrence pattern when and where it matters. Expert knowledge usually cannot  
40 inform us which trait or limitation will be relevant for our problem at hand. We may  
41 know that a palm tree does not survive sub-zero temperatures, but the observed  
42 distribution will tell you that even  $10 \times 10$  km grid squares with minimum  
43 temperatures well below  $0^{\circ}\text{C}$  harbour this species because of microclimatically  
44 suitable places. Thus, at the spatial resolution under investigation, the physiological  
45 threshold can be misleading even though it may be true. Overall, SDMs are useful  
46 for complementing existing approaches in at least these five areas of research:

- 47 1. Small-extent, decision-support for conservation biology (such as Biological  
48 Action Plans: Zabel et al. 2003, and numerous others)
- 49 2. Testing specific hypotheses, e.g. on the spatial scale of habitat selection (Graf  
50 et al. 2005; Mackey and Lindenmayer 2001), the species-energy hypothesis  
51 (Lennon et al. 2000) or range-size effects on diversity pattern (Jetz and Rahbek  
52 2002)
- 53 3. Generating hypotheses, e.g. on correlation of species traits with environmental  
54 variables (Kühn et al. 2006), which can then be tested experimentally
- 55 4. Identifying hierarchies of environmental drivers (Bjorholm et al. 2005; Borcard  
56 and Legendre 2002; Pearson et al. 2004)
- 57 5. Prospective design of surveys, e.g. optimizing sampling schemes for rare species  
58 (Guisan et al. 2006)

59 Now, we shall focus on the technical side, and assume that you know what you  
60 are doing, ecologically speaking.

61 Analysing the geographic distribution of species' occurrence, abundance or  
62 diversity is, essentially, a statistical task. As such, the fundamental ideas and  
63 principles of good statistics apply (and can be found in the excellent but advanced  
64 book of Hastie et al. 2008). There are at least three reasons why methods for  
65 describing or modelling these patterns have reached a higher level of sophistica-  
66 tion than many other fields in ecology. Firstly, biogeographical data sets are  
67 nowadays large (both in terms of number of data points and potentially explanatory  
68 variables), necessitating the use of new statistical strategies. Secondly, species  
69 distribution data typically carry a largish bunch of common intrinsic statistical

problems and accordingly several solutions have been tailored to these problems 70  
 (presence-only data, low information content of binary data, spatial autocorrela- 71  
 tion, multi-collinearity, model unidentifiability). Thirdly, species distribution 72  
 modelling (SDM) is “sexy”. As habitat of many species is continually lost, as 73  
 climate changes and as environmental management becomes a matter of human 74  
 survival, scientists, decision makers and the general public look for information 75  
 and predictions of possible future scenarios. Consequently, substantial funding (at 76  
 least for ecological topics) over the last decade has enabled talented scientists to 77  
 make a career from SDMs. 78

### *Aims of This Chapter* 79

Recent developments have made the field of SDM somewhat complex, diverse and 80  
 confusing for the newcomer. The aim of this chapter is thus to (1) provide a recipe 81  
 for SDM; (2) briefly discuss a few selected “hot” topics; and (3) give an overview of 82  
 challenges of a more ecological modelling type (dispersal, occupancy, biotic inter- 83  
 actions, functional variables, evolution, changing limiting resources). I shall restrict 84  
 citations to fundamental or specific methodological papers and will therefore have 85  
 to ignore the vast amount of good ecological papers that “only” did it right. On the 86  
 other hand, I am not aware of any paper on species distribution modelling that could 87  
 tick all elements of the recipe below. 88

## **13.2 A Species Distribution Modelling Recipe** 89

A good cook needs no recipe. Alas, we are trained more in ecology than statistics. 90  
 Moreover, without the right ingredients (a.k.a. data) and tools (software), no dish 91  
 will be tasty. Also, I should mention other recipes along this line: see Harrell (2001) 92  
 for a generic statistical recipe, and Pearson (2007) and Elith and Leathwick (2009a, 93  
 b) for a specific one on SDMs. As for “cooking tools”, I highly recommend using 94  
 code-based software so that each step of the analysis is documented and easily 95  
 reproducible. The functions mentioned in this chapter are all from the free R 96  
 environment for statistical programming (R Development Core Team 2008). 97

The recipe falls into three sections: pre-processing, modelling and model inter- 98  
 pretation (Fig. 13.1). These sections are somewhat arbitrary, but are useful to 99  
 structure the whole endeavour. We shall assume that you have your ingredients 100  
 well prepared: The observed data are as good as we need them, the explanatory 101  
 variables are ecologically relevant and at the same resolution and your statistical 102  
 tools are laid out in front of you. A worked example is available at [http://www.](http://www.mced-ecology.org) 103  
[mced-ecology.org](http://www.mced-ecology.org) (Where's the sperm whale?), which follows the recipe and 104  
 provides example data and R-code. 105

	SDM task		typical example
pre-processing	response	transformation	log, Box-Cox
	predictor	transformation	square-root
		imputation	multiple imputation
		standardisation	centering, scaling
	ecol. dimensional reduction		resource/direct/proxy
	stat. dimensional reduction		PCA, univariate scans
	collinearity		$ r  < 0.7$ , sequential regression
modelling	model	formulation	splines, interactions
		simplification	BIC, pooling of levels
		type	GLM, GAM, BRT
	spatial autocorrelation		GLMM, SEVM
	diagnostics		Influential points, dispersion
	model performance		AUC, cross-validation
interpretation	functional relationship plots		shape, error margins
	Importances, significances		partial $R^2$ , p-values
	ecological interpretation		plausibility, cf. literature

**Fig. 13.1** Overview of the species distribution modelling workflow. The three phases contain various tasks, for which typical examples are given in the *right* column

106 **13.2.1 Pre-processing and Visualization**

107 **The Response Variable**

108 When the data are presence–absence (i.e. binary) no further preparation is needed.  
 109 When data are counts or continuous, we have to make sure that assumptions of the  
 110 modelling approach are met. For parametric modelling approaches (regressions by  
 111 means of GLM or GAM), count data are usually assumed to be Poisson distributed  
 112 but all too often are not. Continuous responses are generally assumed to be  
 113 normally distributed. These assumptions can be checked only after modelling,  
 114 because we need to look at the residuals or compare log-likelihoods of different  
 115 distributions. Generally, if too many zeros have been observed, the data are over-  
 116 dispersed and we have to resort to one of three alternative approaches: a quasi-  
 117 Poisson distribution (where over-dispersion is explicitly modelled); a negative  
 118 binomial distribution (where a clumping parameter is fitted); or a separate analysis  
 119 of zeros and non-zeros (as in zero-inflated or other mixed distribution models:

Bolker 2008). Sometimes people log-transform count data (more precisely:  $y' = \log(y+1)$ ), and find the new  $y'$  to be normally distributed. 120 121

Normally (Gaussian) distributed data show a normal distribution in the model residuals and a straight 1:1 relationship in a QQ-plot of these residuals. Deviations need to be accounted for, e.g. by transforming the data (any good introductory textbook, such as Quinn and Keough (2002), will feature a section on transformations, including useful ones such as the Box-Cox<sup>1</sup> transformation). 122 123 124 125 126

When we have presence-only data (i.e. only locations where a species occurs but no information where it does not), two alternative approaches are available. We could use purpose-built presence-only methods, or we could use all locations without a presence and call them absences (pseudo-absences). Both approaches have their difficulties (Brotons et al. 2004; Pearce and Boyce 2006). The first suffers from a lack of sound methods (in fact, following e.g. Tsoar et al. 2007 and Elith and Graham 2009), I would currently only recommend MaxEnt<sup>2</sup> in this direction and hope for the approach of Ward et al. (2009) to become publically available). The second approach lacks simulation tests on how to select pseudo-absences and how to weight them (see Phillips et al. 2009 for the cutting edge in this field), although it has been argued that the pseudo-absence approach can be as good or better than the purpose-built presence-only methods (Zuo et al. 2008). In what follows, I only consider presence-(pseudo)absence data. 127 128 129 130 131 132 133 134 135 136 137 138 139

## The Explanatory Variables 140

Explanatory variables may also require transforming! Consider a relevant explanatory variable which is highly skewed (e.g. log-normally distributed), as is commonly the case for land-use proportions. Few high-value data points may completely dominate the regression fitted. To give a more balanced influence to all data points, we want the values of the predictors to be uniformly distributed over their range. This will rarely be achievable, and researchers mostly settle for a more or less symmetric distribution of the predictor. Note, however, that ideally we want most data points where they help most. For a linear regression, the mean is always best described, so we would want most data points at the lowest and highest end of the range. For a non-linear function, for example a Michaelis–Menton-like saturation curve, we want most data points in the steep increase, while there is little gained from many points at the high end, once the maximum is reached. As a rule of thumb we need many data points where a curve is changing its slope. 141 142 143 144 145 146 147 148 149 150 151 152 153

Transformation of explanatory variables is particularly needed for regression-type modelling approaches such as GLM and GAM (see below for explanation). Regression trees (used, e.g. in Boosted Regression Trees, BRT, or randomForest) are far less sensitive, if at all (Hastie et al. 2008). It is a good custom to make a 154 155 156 157

<sup>1</sup>boxcox in **MASS** (typewriter and **bold** are used to refer to a function and its **R-package**)

<sup>2</sup>Phillips et al. (2006b): <http://www.cs.princeton.edu/~schapire/maxent/>

158 histogram of each explanatory variable before entering it into an analysis! Trans-  
159 formation options are the same as for the response.

160 Missing data are a (very) special case of transformation. Although generally  
161 disliked by many analysts, imputation (replacement of missing data) is often a good  
162 idea (see Harrell 2001), particularly if missing data are scattered through the data  
163 set (i.e. across several variables!) and we would lose many data points if we simply  
164 omitted every data point with missing values. Standard imputation uses the other  
165 explanatory variables to interpolate a likely value for the missing one.<sup>3</sup> Replace-  
166 ment by the mean is not an option!

167 “Outliers” are (in general) a red herring: If there is no methodological reason  
168 why a data point is extremely high (e.g. one data set being recorded in winter, while  
169 all other data points are from the summer), then this datum should also be included  
170 in the analysis. Otherwise the data set may be poorly sampled, but the “outlier”  
171 would still represent a (potentially) valuable datum. It would be good practice to  
172 omit it later on and see if the results are robust to this omission. Furthermore, in  
173 multi-dimensional data sets (i.e. those with several explanatory variables), a datum  
174 might be an “outlier” in one dimension, but an ordinary data point in all others: why  
175 delete it?

176 Finally, all continuous variables should be standardized before the analysis.<sup>4</sup>  
177 This reduces collinearity, particularly with interactions (Quinn and Keough 2002).  
178 As a convenient side effect, regression coefficients are now directly comparable:  
179 the larger their absolute value, the more important this term is in the model (they  
180 become standardized regression coefficients).

## 181 Collinearity

182 Collinearity refers to the existence of correlated explanatory variables. Some  
183 predictors are only proxies for an underlying, latent variable. For example, consider  
184 temperature and rainfall, which are largely governed by distance to ocean (ocean-  
185 ity), altitude and regional terrain. Collinear predictors can lead to biased models due  
186 to inflated variances (Quinn and Keough 2002). There are many cures to this  
187 ailment, but no remedy. Logically speaking, if two predictors are tightly linked to  
188 an underlying (but elusive) causal variable, there is no way to find out which is the  
189 “correct” predictor for our analysis. We may choose precipitation over temperature  
190 when modelling plants (or the other way around for insects), but there is no  
191 guarantee that this choice allows us a sensible extrapolation of our model. Further-  
192 more, using Principal Component Analysis (PCA) or any of the other tailor-made  
193 methods for collinearity (Partial Least Squares, penalized regression, latent root  
194 regression, sequential regression, and many others) will not solve the ecological  
195 problem, only the statistical. These methods will produce either a new data set of

---

<sup>3</sup>transcan and aregImpute in **Hmisc**

<sup>4</sup>scale

uncorrelated variables, or “consider” the correlation when estimating model parameters. 196  
197

So where is the problem? Imagine an organism whose distribution is governed entirely by its sensitivity to frost. When we combine our climate variables into one or more principal components, model the species' distribution, and then predict to a climate change scenario, the fact that both rainfall and mean summer temperature are correlated with number of frost days will dilute its impact in the model. The total effect of “frost” is distributed over all correlated variables. As a consequence, any climate prediction will underestimate the effect of frost and hence yield a “wrong” expected future distribution. If we don't know the true underlying causal mechanism, no statistics can help us here (or at least very little). Any ecological knowledge used in variable pre-selection, however, will lead to a smaller bias in scenario projections! 198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208

**Dimensional Reduction** 209

Often we may have dozens or even hundreds of potential explanatory variables (e.g. from multispectral remote sensing or landscape metrics). We should try to reduce this set to as few as possible for two reasons: (1) The more variables we have, the more they will be correlated. (2) The more variables we have, the more likely one of them will spuriously contribute to our model (type I error). For SDMs, Austin (2002) and Guisan and Thuiller (2005) argue that we should choose “resource” over “direct” and “direct” over “indirect” variables. For example, the abundance of prey (hardly ever available) or nesting opportunities will be a resource variable when analysing the distribution pattern of a bird of prey. Temperature or human disturbance could be direct variables, impacting on the bird without moderation by other variables. Indirect variables would be altitude or length of road in a grid cell, which are substitutes, surrogates or proxies for other, more directly acting variables. These indirect variables are often not immediately perceivable by the organism (such as altitude by a plant or length of road verges by a rodent). So if we have two (correlated) variables, we should discard the one “further away” from the species' ecology. 210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225

If we are unable to reduce the data set sufficiently (i.e.  $k \gg N$ ), we should use dimensional reduction techniques, such as Principal Component Analysis<sup>5</sup> or its more sophisticated variants that also allow categorical variables (nMDS<sup>6</sup>). The scores for the most important axes in this new parameter hyperspace can be used as explanatory variables. Note that interpretation is often extremely impaired by automatic dimensional reductions. It is thus always advisable to use ecological understanding rather than statistical functions at this step! 226  
227  
228  
229  
230  
231  
232

---

<sup>5</sup>prcomp

<sup>6</sup>isoMDS in MASS or, more conveniently, metaMDS in vegan

233 An alternative is to “filter” the data by importance. We can use a robust and able  
234 technique to tell us which variables are important. Next, we use only those 5 or 12  
235 variables filtered from the initial pool of variables, and continue. Regression-tree  
236 based methods are very useful for this, and I recommend randomForest and Boosted  
237 Regression Trees. If you plan to model your data with BRT anyway, there is little  
238 point in reducing the data before.

239 Finally, be aware that any model can only find correlations with the variables  
240 provided. Of course, we know that our hypothetical bird of prey depends on specific  
241 prey. Without this information, we may actually be modelling the niche of the prey,  
242 not of the predator!

### 243 Exploratory Data Plotting

244 Can we finally start? No! It is both good practice and highly advisable to look at the  
245 data by plotting them in any reasonable combination conceivable (see, e.g., Bolker  
246 2008). Plot thematically related explanatory variables as scatterplot<sup>7</sup> to detect  
247 collinearity. Plot each explanatory variable against the response (henceforth called  
248 X and y, respectively) and look for nonlinear effects. Plot y against two Xs  
249 (Fig. 13.2) and a hull-polygon around the data to see that 40% or so of the parameter  
250 space is not in your data set. This is the area outside the convex hull in Fig. 13.2.  
251 The more variables (and hence dimensions) your data set has, the more severe this  
252 problem becomes. It is so prominent among statisticians (though not among  
253 ecologists) that it is referred to as the “curse of dimensionality” (Bellman 1957;  
254 Hastie et al. 2008). Repeat this plotting for any number of variables. Getting a  
255 feeling for the data is crucial, and many later errors can be avoided. Every minute  
256 invested at this stage saves hours later on.

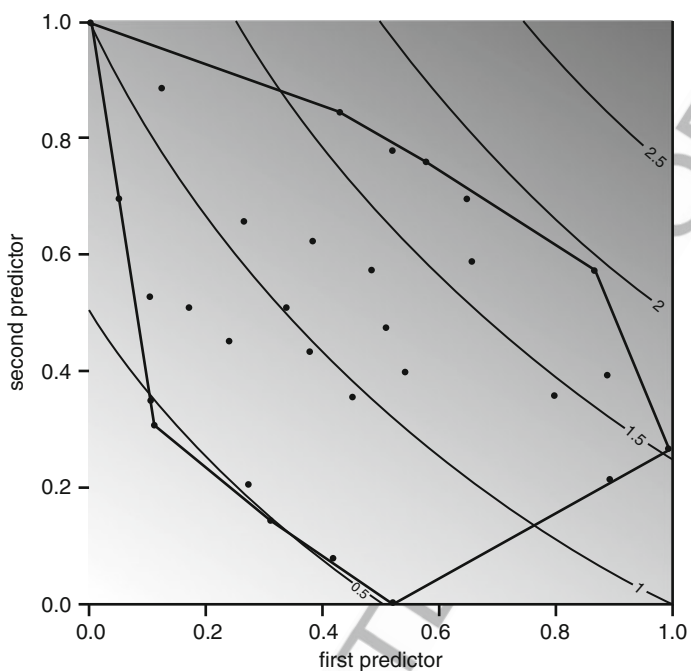
### 257 13.2.2 Modelling

258 Here, we arbitrarily divide the process of deriving a “usable” model into two steps.  
259 The first, model building, selects the variables to be included, the type of non-  
260 linearity and order of interactions considered, and the criteria for selecting the final  
261 complexity of a model. The second step, model parameterisation, performs the final  
262 step of using the data to calculate the best estimates for variable effects. It is this  
263 model that we want to use for interpolation, hypothesis testing or extrapolation. Note  
264 that in some methods these two steps are implicitly taken care of and that there is no  
265 two-step process (mainly machine learning, where model selection is done internally  
266 through cross-validation in order to prevent models from being “unreasonably”

---

<sup>7</sup>pairs





**Fig. 13.2** Visualizing the parameter space supported by data. In this case, the *top-right* and *bottom-left* corner of the parameter space of the two predictors has not actually been sampled by data (despite a low correlation of  $r = -0.26$ ). The parameter space actually sampled is indicated by the convex hull covering 57% of the area, declines dramatically with the number of dimensions (“curse of dimensionality”). In other words: we have few data points to look at interactions of higher order

large: e.g. Hastie et al. 2008). For more traditional approaches (and here I am 267  
 thinking of GLMs), we may want to have these steps functionally separated. 268

**Model Formulation**

269

We have reduced our data set to a moderate number of predictors in the step 270  
 “Dimensional reduction” above. Now we still need to specify in which functional 271  
 form the predictors are allowed to correlate with the response. In early years, both 272  
 non-linear and interactive model terms were neglected, making many of their 273  
 findings less trustworthy. Modern methods (such as BRT) will automatically have 274  
 non-linearity and interactions build-in. It is still important to understand the rele- 275  
 vance of non-linearity and interactions, even when using the tree-based methods, 276  
 because we still have to be able to interpret the results. The information on the 277  
 importance of a variable often returned by machine-learning algorithms does not 278  
 allow us to see *how* the variables act. As shown in the case study at [http://www.](http://www.mced-ecology.org) 279  
[mced-ecology.org](http://www.mced-ecology.org) (Where’s the sperm whale?), the functional relationship must be 280

281 plotted to gauge its shape. For interactions we need to plot each variable at each  
282 level of the other variable, thus visualizing synergistic or compensatory effects of  
283 the two variables.

284 The key idea behind SDM, i.e. the environmental niche of a species, implies a  
285 hump-shaped relationship between any environmental predictor and a species'  
286 occurrence: there are lower and upper limits. Hence, we *must* allow the model to  
287 be nonlinear. If we happen to only sample a part of the entire gradient, we also  
288 need to consider saturation curves, which are again non-linear. The simplest, and  
289 generally sufficient, way to include non-linearity is by generating a new, squared  
290 dummy variable for each continuous predictor.<sup>8</sup> This represents the third element  
291 of a Taylor series (which can be expanded to represent any function). When using  
292 GAM or other spline-based approaches, non-linearity is governed by the smoothing  
293 function used. Here the issue is not so much how to model non-linearity, but  
294 rather how much non-linearity we allow for. Reducing the “wiggleness” of splines  
295 (either by stepwise model selection for the number of knots in each predictor<sup>9</sup> or  
296 by shrinkage of spline fits<sup>10</sup>) prevents over-fitting and should be the standard  
297 approach.

298 Interactions are similarly relevant. Statistically, an interaction is the product of  
299 the participating main effects. Ecologically, it means that we need to know the  
300 value of all variables included in the interaction, not only the main effects. Because  
301 this is highly relevant and often difficult for the beginner, let me briefly give an  
302 example. Assume that global patterns of plant diversity are well-predicted by the  
303 predictors “annual precipitation” and “mean annual temperature” – and their  
304 interaction. For the main effects, wet or hot means more species, but not necessar-  
305 ily. When a site is hot, it needs to *also* be wet to have high species richness;  
306 otherwise it may well be a barren desert. But when cold, a site will never support  
307 many plant species, independent of precipitation. In this example, neither tempera-  
308 ture nor rainfall alone is sufficient to predict species richness at any site, but we  
309 need to interpret them in concert.

310 Classification and regression trees (CARTs) embrace non-linearity and interac-  
311 tions in an elegant and natural way. Their boosted (BRT) or bagging (randomFor-  
312 est) extensions hence do not require specification of non-linearity and interactions.

### 313 Model Simplification

314 One of the fundamental problems in building statistical models is the trade-off  
315 between the variance explained by the model, and the bias it produces when

---

<sup>8</sup>This can be done either manually ( $X1.2 <- X1^2$ ) or as part of the model formula ( $y \sim X1 + I(X1^2)$ ); higher-order polynomials should be specified using `poly` ( $y \sim \text{poly}(X1, \text{degree}=3)$ ), which calculates orthogonal polynomials.

<sup>9</sup>As proposed for the function `gam` in package **gam**: see `?gam::step.gam`

<sup>10</sup>As proposed for the function `gam` in package **mgcv**: see `?mgcv::step.gam`.

validating it on a new, independent data set (variance-bias-trade-off: Hastie et al. 316  
2008). Smaller models are more robust, i.e. less biased, at the expense of being not 317  
very good in explaining variance. The way to derive the “optimal” model size is 318  
through cross-validation (CV). For some modelling approaches this is automati- 319  
cally implemented, but the majority of model types require the user to carry out this 320  
step.<sup>11</sup>  $N$ -fold cross-validation encompasses a random assignment of data points to 321  
the  $N$  subset, with  $N$  usually between 3 and 10. Care should be taken to have equal 322  
prevalence in all subsets, e.g. by randomizing 0s and 1s separately (stratified 323  
randomization). The model is then fitted to  $N-1$  of the  $N$  subsets and evaluated 324  
on (by predicting to) the remaining subset. This is repeated for all  $N$  subsets and 325  
evaluations are averaged. Based on these values, we can select the best modelling 326  
strategy (both model complexity and model type). An alternative approach is to 327  
bootstrap the entire model building process and use bootstrapped measures of 328  
model performance. Since a bootstrap requires several thousand runs, and a CV 329  
only a few, CV is far more common. 330

Information theoretical approaches are based on analytical methods to describe 331  
this CV. Hence Akaike's Information Criterion (AIC) or Schwartz'/Bayesian 332  
Information Criterion (BIC) are implicitly also based on cross-validation. While 333  
it is clear that too large a model will be over-fitting, and that too small a model will 334  
not capture as much of the variation as it should in the data, the “true” model will 335  
always remain elusive, and our “optimal” model will only be a caricature of the 336  
truth. However, here is much to be learned from this caricature! 337

## Model Type

338

At this point we have to choose one (or more) method(s) to do our analysis with. 339  
The good “traditional” approaches comprise Generalised Linear Models (GLM) 340  
and Generalised Additive Models (Guisan and Zimmermann 2000). Discriminant 341  
Analysis has been given up on, as have been Neural Networks and CARTs (Guisan 342  
and Thuiller 2005). “Modern” approaches are often based on either multidimen- 343  
sional extensions of GAMs (such as MARS and SVM) or machine-learning varia- 344  
tions of CART (such as BRT and randomForest: Hastie et al. 2008). Anyone using a 345  
machine-learning method should familiarize himself with this method. The major- 346  
ity of them are performed on real data sets, where the truth is unknown and the 347  
performance of a method was hence assessed by cross-validation. These compar- 348  
isons show, broadly speaking, that model types sometimes differ dramatically in 349  
performance, that each model type can be misused and that both GLM and BRT are 350  
reliable methods when used properly. 351

---

<sup>11</sup>Do not confuse out-of-bag model weights and alike with cross-validation of the entire model. When data are held back during the building of sub-models (e.g. in randomForest or BRT), this does not represent a cross-validation of the entire model (i.e. the average of sub-models).

352 This is not the place to explain the differences between all of them (see Hastie  
353 et al. 2008 for a recent and comprehensive description or Elith and Leathwick  
354 2009a). It has to suffice to make clear the main difference in the machine-learning  
355 approach to “traditional” statistical models. In traditional models (e.g. GLM), we  
356 specify the functional relationship between the response and its predictors. For  
357 example, we decide to include precipitation as a non-linear predictor for plant  
358 species richness. This model proposal is then fitted to the data. In machine learning,  
359 we propose only the set of predictors, but not the model structure. Here, an  
360 algorithm builds a model proposal, fits it to a part of the data set and evaluates its  
361 performance on the other part of the data. It then proposes a modification of the  
362 original model and so forth. Machine-learning algorithms<sup>12</sup> differ in scope, origin,  
363 complexity, and speed, but they all share this validation step which is used to steer  
364 the algorithm towards a better model formulation. There are plenty of studies  
365 comparing different modelling approaches (Guisan et al. 2007; Meynard and  
366 Quinn 2007; Pearson et al. 2006; Segurado and Araújo 2004). Rather, we shall  
367 continue using GLM and BRT as representatives for the two most common good  
368 approaches.

369 The choice of model type has much to do with availability of software, current  
370 fashion and, of course, with the specific aim of the study. Further complications  
371 arise if the design of the survey may require a mixed model approach (e.g. due to  
372 repeated measurements or surveys split across observers), if spatial autocorrelation  
373 needs to be addressed, if zero-inflated distributions have to be employed, and if  
374 corrections for detection probability shall be modeled. The more additional require-  
375 ments are imposed on the model, the more GLMs become the sole possible  
376 method.<sup>13</sup> Alternatively, you may want to go for a Bayesian SDM (see Latimer  
377 et al. 2006, for a primer).

378 If your data and model require an unusual combination of steps (say a combina-  
379 tion of zero-inflated data with nested design and spatial autocorrelation, while  
380 predictors are highly correlated and many values missing), and you develop a  
381 way to cook this dish, then you should do (at least) two things: Firstly, evaluate  
382 your method for its ability to detect an effect that you *know* is there (“power”).  
383 Secondly, evaluate your method for its sensitivity to detect effects that you know  
384 are *not* there (“type I error”). Both evaluations should be amply replicated, should  
385 be based on simulated data (so that you know the truth) and should (finally) confirm  
386 that your new methods is reliable!

---

<sup>12</sup><http://www.machinelearning.org/> is a good place to start exploring this field

<sup>13</sup>Most of these “complications” can be handled by standard extensions of GLMs (see, e.g. Bolker 2008, and various dedicated R-packages). They will, however, make the model less stable, require larger run-times and still rely on getting the distribution right. There is, of course, the alternative of Bayesian implementations. Since these are also fundamentally maximum likelihood approaches, they are similar to sophisticated GLMs. In any case, there is no Bayesian Boosted Regression Tree (not to speak of a combination with spatial terms and mixed effects). It runs against the Bayesian philosophy to use boosting or bagging, and there is no efficient implementation either.

## Spatial Autocorrelation

387

Spatial autocorrelation (SAC) refers to the phenomenon that data points close to each other in space are more alike than those further apart. For example, species richness in a given site is likely to be similar to a site nearby, but very different from sites far away. This is mainly due to the fact that the environment is more similar within a shorter distance. Hence, SAC in the raw data (species richness) is a consequence of SAC in the environment (topography, climate), something Legendre (1993) termed “spatial dependence”. In SDMs, we do not care about SAC per se, but about SAC in the model’s residuals (i.e. unexplained by the environment), because it distorts model coefficients (Bini et al. 2009; Dormann 2007a). To date it is unclear whether this residual SAC is mainly due to model misspecification (omission of non-linearity and interactions), due to variation in sampling coverage, due to omission of important predictors, or due to ecological processes (territoriality, dispersal). Only with respect to some of these problems can *statistical* solution be found. The spatial toolbox is rich in approaches (Beale et al. 2010; Carl et al. 2008; Dormann et al. 2007; Mahecha and Schmidtlein 2008). In any case, SDM residuals should be investigated for spatial autocorrelation, and attempts should be made to correct for it. If spatial models yield similar coefficient estimates (GLM) as non-spatial models, then there seems to be little value in “going spatial”: the ranges of the spatial autocorrelation may or may not be related to the ecological scale of movement or behavioral patterns (Betts et al. 2009; Dormann 2009).

## Tweaking the Model

408

There are several ways in which the quality of the model can be increased (Maggini et al. 2006). One important start is to investigate the model residuals. They indicate whether model assumptions were violated (e.g. when residuals are highly skewed or their variance is not the same throughout the range of fitted values) or if some non-linear relationship went unnoticed (residuals may, e.g. show a hump-shaped trend against fitted values).

Model diagnostics<sup>14</sup> will also indicate outliers, i.e. data points that have a high influence on the model coefficients. We can use weights to decrease an outlier’s impact. Weights are also useful when the balance between presences and absences is very disturbed. Down-weighting the more common category so that model weights sum to the same value for 0s and 1s has been shown to increase the sensitivity of binomial models (Maggini et al. 2006). The same approach is recommended when using pseudo-absences (Elith and Leathwick 2009a).

By including data from other scales or broader geographic coverage, regional or local SDMs can also be improved. Pearson et al. (2004) used European distribution and climate data to fit a niche model for four plant species. Predicted probabilities

---

<sup>14</sup>Diagnosics for GLMs fitted in R are given by plotting the model object.

425 of occurrence from this model were then used as input variable alongside land-  
426 cover variables in the second-step model for the UK. Thereby the authors avoided  
427 the problem that the climate gradient in the UK is much shorter than of the species'  
428 global distribution.

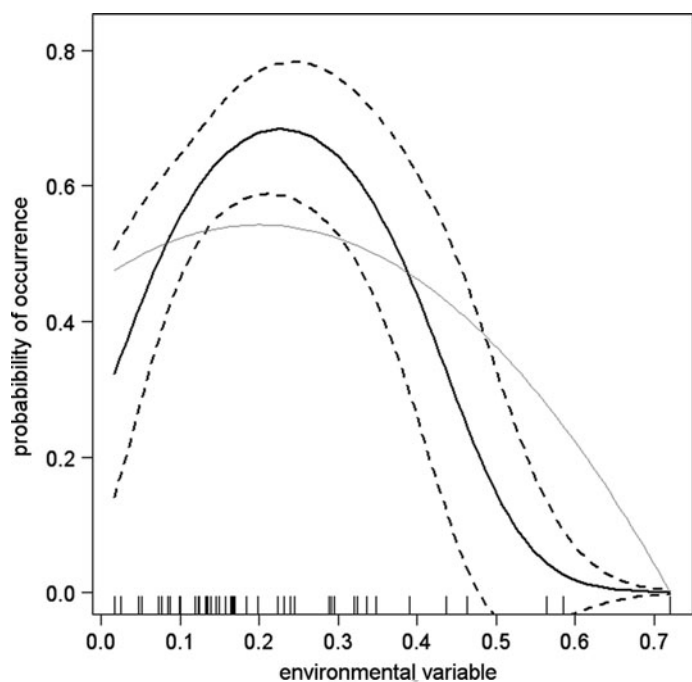
## 429 **Assessing Model Performance**

430 To quantify how well our model fits the data, we compare model predictions with  
431 field data (usually on a hold-out sample; e.g. the subset of a cross-validation).  
432 Traditionally, the probability predictions from the model were converted into  
433 presences and absences and then a confusion matrix could be used to calculate  
434 various parameters of choice (e.g. commission and omission error, kappa, etc:  
435 Fielding 2002). The AUC (“area under curve”) is currently the most commonly  
436 used measure of discriminatory power of a model. Its value (between 0.5 for  
437 random and 1 for perfect) quantifies the ability of the model to put the data points  
438 into the correct class (i.e. presence or absence), independent of the threshold  
439 required by the other measures mentioned. It has recently received justified criti-  
440 cism because its values are not comparable across different prevalences (and the  
441 criticism extends to kappa, too; see Lobo et al. 2008). Currently, misclassification  
442 rates, commission and omission errors are more *en vogue* again, because they can  
443 be intuitively interpreted. Furthermore, by assigning different weights to false  
444 negatives (omission error) and to false positives (commission error), conservation  
445 management can come to more sophisticated and balanced decisions (Rondinini  
446 et al. 2006).

447 Only rarely will a second set of data be available to investigate the quality of our  
448 model(s) through external validation. A different recording strategy, another time  
449 slice or data from a different geographic location represent really independent data,  
450 and could thus be considered an external validation. The internal validation  
451 (described above as cross-validation) is an optimistic assessment of model quality.  
452 When using SDMs to infer underlying mechanisms, external validation is less of  
453 an issue than when using them to extrapolate to a future climate or other sites.  
454 Because the cross-validated models are optimistic, they give narrower error bands  
455 than they should.

### 456 **13.2.3 Interpretation**

457 Once we have arrived at what we regard as a final model, we should make every  
458 effort to understand what it means. A first and most relevant step is to visualize the  
459 functional relationships within the model. The plot of how occurrence probability  
460 is related to, say, annual precipitation should be accompanied by a confidence band  
461 around this line. It may be useful to plot the data as rug (ticks on the axis representing



**Fig. 13.3** Functional relationship between an environmental variable and a binary response. Rug (ticks on lower axis) indicate for which  $x$ -values data were available. Lines represent a quadratic fit (solid) and its standard deviation (dashed). Thin grey line is the true, underlying, data-generating function. Note the few data points upon which the declining half of the function is based (6 of a total of 50 have a value  $>0.35$ )

positions of the data values) into this figure to visualise the support at each point in parameter space (Fig. 13.3).

For interactions, visualization becomes more difficult. Two-way interactions can still be plotted (e.g. as a 3-D plot or as a contour plot). No confidence bands can be included, though. Here it is again very important to indicate the position of the samples to identify regions of the parameter space that have not been sampled. For higher dimensions, or for a model that averages across many sub-models, we can do the same plots (called marginal plots for main effects because they represent the marginal changes to a predictor, averaging across all other predictors). We can also slice through higher dimensions, i.e. calculate a marginal plot for specific values of other predictors (often their median).

Spending time plotting is again well invested. We will detect errors in the model, scratch our head over inexplicable (and hence overly complex) patterns, and be forced to extract the main conclusions from it. It is this phase where the traditional GLM is superior to the BRT, because variable interpretation is easier. It is, however, also this phase where we may realize that BRTs are superior to GLMs because they can model step-changes and thresholds much better. Personally,

479 I think we should not publish patterns we do not understand. There are, as the  
480 previous steps have shown, several decisions that could generate artifacts and their  
481 publication cannot be seen as progress.

### 482 **13.3 Beyond Recipes: New Challenges for Species** 483 **Distribution Models**

484 The above recipe can be used to derive a static description of environmental  
485 correlates with distribution data. But they often leave the analyst unsatisfied.  
486 Many assumptions might be suspected to be violated (Dormann 2007b), such as  
487 stationarity, unbiased coverage, or equilibrium with environment.<sup>15</sup> There are three  
488 key challenges with the following problems of current SDMs: (1) A niche model  
489 describes the *current* niche, and it is unclear which factors will be limiting elsewhere  
490 or in the future. (2) Climate change projections delimit only the *potential* future  
491 distribution, and it is unclear whether the species will ever fill this new range; e.g.,  
492 due to dispersal constraints. And, (3), our understanding of the adaptive potential of  
493 a species is currently very poor. This sets, in part, the research agenda for species  
494 distribution models. Let us look at these challenges in more detail.

#### 495 ***What Limits a Species' Range?***

496 An organism is constrained in its population dynamics by resources, competitors,  
497 predators and diseases, density dependence, reproductive opportunity, mutualists,  
498 environmental stochasticity and so forth (e.g. Krebs 2002). The same holds true for  
499 its spatial distribution (e.g. Gaston 2009; Holt and Barfield 2009), but additionally  
500 spatial constraints come into play (e.g. distance between habitat fragments, mini-  
501 mum territory size, Allee effects due to low population density). With an SDM, we  
502 are usually only able to quantify some of these limitations, and, accordingly, SDMs  
503 often do not transfer very well to other sites (Schröder and Richter 1999; Randin  
504 et al. 2006; Duncan et al. 2009, but see Herborg et al. 2007). In particular, biotic  
505 interactions are hardly ever quantified explicitly within SDMs (although some of  
506 them will also correlate with the environmental data used). But they matter in real  
507 data (e.g. Preston et al. 2008; Schweiger et al. 2008), and they impact model  
508 performance and predictive ability (as shown, in a simulation study, by Zurell  
509 et al. 2009). It is thus a key challenge to incorporate biotic interactions into SDMs,

---

<sup>15</sup>Actually, the term “equilibrium” is a bit misleading. What is meant is that the entire width of its niche is filled. Within this niche, there may well be unoccupied sites, e.g. due to metapopulation dynamics. A problem arises, when a species does not occupy say the dry end of its soil moisture niche for historic reasons. Then the estimate of this end of the niche will be biased.



but to date such attempts are few and far between (e.g. Bjornstad et al. 2002). A recent review (Thuiller et al. 2008) indicates some avenues to do so, but all of them are based on the explicit modelling of populations within cells, if not individuals. It is unclear, how to derive a more general dynamic SDM without spending years per species on incorporating detailed ecological knowledge.

### ***Fundamental, Potential and Realized Niches***

The reason why many biogeographers refer to SDMs as “Species Distribution Models” and not as “niche models” is because they do not believe that we model the niche of the target species. In fact, as Jiménez-Valverde et al. (2008) argue, because we do not know *why* a species is absent in some sites, we are in the dark about its niche. The discussion of what a niche is, and what we are modelling, has sparked several interesting and not always compatible publications (e.g. Kearney 2006; Soberón 2007). Hence, Araújo and Guisan (2006) have named the “clarification of the niche concept” the first of five challenges for SDMs. While we cannot resolve this issue here, it is important to realize that the “niche” based on the correlation between geographic distributions and environmental conditions is quite a bit more vague than the niche discussed in evolutionary ecology, where resources and other causal drivers are envisaged (see, e.g. Losos 2008).

More to the point in this context is the challenge to quantify how much of the fundamental niche is actually covered by the realized niche as extracted from SDMs. If, on one extreme, the realized niche is pretty much also the fundamental niche (i.e. there are no biotic interactions alike to constrain the distribution at the scale we are analyzing), then we can merrily predict future distributions of this species (e.g. under climate or land-use change). At worst, we are overestimating the future, “potential” distribution (if in the future biotic interactions may become limiting or if species do not reach the sites). At the other extreme, if the fundamental niche is considerably wider than what we model, any projection can be fundamentally flawed (Dormann et al. 2010). I am not aware of any study assessing the overlap of realized and fundamental niche for geographic distributions (see also Nogués-Bravo 2009). It could require transplant experiments into areas beyond the current range and the manipulation of biotic interactions there. The few studies going into this direction point at a large discrepancy between fundamental and realized niche. Battisti et al. (2006), for example report on a range shift after a particularly warm summer, which was not reverted afterwards, indicating that it was dispersal limitation that prevented a filling of the niche. Similarly, several studies point at the importance of dispersal limitation (Nekola 1999; Ozinga et al. 2005; Samu et al. 1999; Svenning and Skov 2004), leading to both a bias in the modeled environment-occurrence relationship as well as the width of the niche itself. There is, as yet, no standard way to wed SDMs and dispersal (see Johst et al. 2002; King and With 2002; Lavorel et al. 2000; Lischke et al. 2006; Midgley et al. 2006; Schurr et al. 2007 for attempts, Thuiller 2004).

551 *Niche Evolution*

552 Another important and fast developing field related to species distribution model-  
553 ling is the study of niche evolution. I shall use this term very loosely, as is often  
554 done, to also include micro-evolutionary changes, genetic (and ecological) drift  
555 within species and genotypic plasticity (Pfenninger et al. 2007). Climate change  
556 projections using SDMs rely on the assumption that species are not able to adapt  
557 significantly to altering environmental conditions. This assumption is implicit in the  
558 extrapolation of the fitted niche: if a species was able to adapt rapidly, then the  
559 present niche would not be related to its future niche.

560 The problem is that we have considerable, if patchy, evidence that niches can  
561 rapidly evolve (reviewed in Thompson 1998), change within the fundamental niche  
562 (Dormann et al. 2010) or at least that variability within a species is large enough to  
563 allow it to shift its niche when confronted with novel environments (e.g. Ackerly  
564 et al. 2006; Broennimann et al. 2007; Hajkova et al. 2008; Holt 2003; Holt and  
565 Gaines 1992). There is, as yet, no synthesis of niche evolution nor, to my knowledge,  
566 any mechanistic approach to incorporate geno- and phenotypic plasticity into SDMs  
567 or spatial population models. There is, on the other hand, more than anecdotal  
568 evidence that microevolutionary processes are at play and matter ecologically  
569 (Hampe and Petit 2005; Phillips et al. 2006a). Hence, this field still awaits being  
570 embraced by species distribution models.

571 **13.4 Concluding Remarks**

572 This chapter tries to strike a balance between guidance for novices to species  
573 distribution modelling – by providing a recipe for the most crucial elements of  
574 SDMs – and an embedding of SDMs into the currently most relevant statistical  
575 challenges. Analysing a species' distribution can be a very useful starting point for  
576 further investigations or process-based modelling attempts. The correlative nature  
577 of modelling in general, and species distribution modelling specifically, should  
578 always be remembered. Tempting as it may be to incorporate a lot of ecological  
579 knowledge into mechanistic or statistical models, only little of this information will  
580 actually be relevant *at the focal scale*. The main intellectual challenges that remain  
581 to not over-interpret one's findings and to seek independent corroboration.

582 **Acknowledgments** Over the years, many colleagues helped develop the above recipe. I am  
583 particularly grateful to Boris Schröder, Björn Reineking and Jane Elith, as well as the many  
584 participants of statistical workshops on this topic. I am also grateful to Fred Jopp, Hauke Reuters  
585 and Dietmar Kraft for improving a previous version. Funding by the Helmholtz Association is  
586 acknowledged (VH-NG-247).

# Metadata of the chapter that will be visualized online

---

Series Title		
Chapter Title	Modelling Species' Distributions	
Chapter SubTitle		
Copyright Year	2011	
Copyright Holder	Springer-Verlag Berlin Heidelberg	
Corresponding Author	Family Name	Dormann
	Particle	
	Given Name	<b>Carsten F.</b>
	Suffix	
	Division	Department of Computational Landscape Ecology
	Organization	Helmholtz Centre for Environmental Research – UFZ
	Address	Permoserstr.15, 04318, Leipzig, Germany
	Email	carsten.dormann@ufz.de

---

---

**Abstract** Species distribution models have become a commonplace exercise over the last 10 years, however, analyses vary due to different traditions, aims of applications and statistical backgrounds. In this chapter, I lay out what I consider to be the most crucial steps in a species distribution analysis: data pre-processing and visualisation, dimensional reduction (including collinearity), model formulation, model simplification, model type, assessment of model performance (incl. spatial autocorrelation) and model interpretation. For each step, the most relevant considerations are discussed, mainly illustrated with Generalised Linear Models and Boosted Regression Trees as the two most contrasting methods. In the second section, I draw attention to the three most challenging problems in species distribution modelling: identifying (and incorporating into the model) the factors that limit a species range; separating the fundamental, realised and potential niche; and niche evolution.

---