

WhyWhere2.0: An R package for modeling of species distributions on big environmental data

David R.B. Stockwell*

May 30, 2015

Abstract

Previous studies have indicated that multi-interval discretization (segmentation) of continuous-valued attributes for classification learning might provide a robust machine learning approach to modelling species distributions. Here we apply a segmentation model to the *Bradypus variegatus* – the brown-throated three-toed sloth – using the species occurrence and climatic data sets provided in the niche modelling R package `dismo` and a set of 940 global data sets of mixed type on the Global Ecosystems Database. The primary measure of performance was the area under the curve of the receiver operating characteristic (AUC) on a k-fold validation of predictions of the segmented model and a third order generalized linear model (GLM). This paper also presents further advances in the **WhyWhere** algorithm available as an R package from the development site at <http://github.com/davids99us/whywhere>.

1 Introduction

Ecological niche models such as BIOCLIM [10] and GLMs [1] have been the mainstay of species distribution modelling (SDM). These models are most effective with climatic variables when the response of the species to its environment is continuous and unimodal. Symmetry in the response has been associated with equilibrium of the species with its environment. However non-equilibrium due to ecosystem or climate disturbance and/or species mobility means the a simple and symmetrical niche shape is most likely the exception and not the rule [9].

In contrast to climatic-based SDMs, applied conservation studies tend to focus on habitat features, or habitat suitability indices (HSIs), based on the structure of the environment such as availability of nesting sites, physical barriers and resource sources. HSIs are expressive of finer scale forest and aquatic environments where climatic variables do not down-scale. HSIs are proximal causes of occurrence of species than more distal climate envelopes. Habitat

*All correspondence to be addressed to the author at david.r.stockwell@cqu.edu.au, Adjunct Researcher at Central Queensland University.

variables are often categorical factors and there are very many possible habitat factors that could determine each species. Integration of continuous and categorical models into a single model can be a challenge.

On the basis of these and other studies we hypothesized in [18] that segmented models could be used to fit arbitrary responses, would be at least as accurate as continuous models, and be of greater accuracy when mining large environmental data sets that contain mixed habitat and climatic variables. The second hypothesis was confirmed using the segmented modelling method called **WhyWhere** [18, 17] (WW1) and both are again addressed in this new version (WW2). Moreover, we present theoretical and practical advances of the algorithm enabling the goal of incorporating species distribution modelling into an artificial intelligence environment [2].

2 Methods

A package for species modelling in R called **dismo** includes a number of popular methods including MaxEnt, Bioclim, Domain, GLM, GAM, and RandomForest. This package implements elements of the SDM work flow originally developed in the GARP machine learning system [14], such as pseudo-absences, whereby a random sample of the environment is in lieu of true absence values. The **dismo** package contains bioclimatic variables from the WorldClim database [8] and terrestrial biome data on terrestrial ecoregions [11].

Other R packages that have made **WhyWhere** more accessible are the **raster** package for handling gridded spatial data [19] and a new package called **data.table** for fast aggregation of large data sets [3].

Any species modelling method must address three main stages, and problems arising at any of these stages can lead to poor results: (1) getting the environmental data into a uniform form for analysis (2) determining the best type of model to use to represent the response of the species (3) the interpretation of the results.

2.0.1 Environmental variables

As geographic variables come in varying extents, projections and resolutions, they must generally be unified by coregistration before modelling. But this entails additional processing and inefficiencies. Smaller resolution sets must be enlarged redundantly and information lost when contracted. Memory needs burgeon when coregistering many variables to the finest scale. An ideal method would utilize each data set at its native projection and resolution.

After coregistration of geographic layers, most statistical models input a 'wide' file format, i.e. with variables in columns and locations in rows. This adds an intermediate processing step where all environmental data must be then be held in memory, which can hit memory limitations at the prediction stage when the models are applied across the geographic space.

73 The approach in WW1 was to transform geographic variables into a compact
 74 image format and then to combine at most three variables into the image in the
 75 red, green and blue channels. While fast image processing packages could then
 76 be used for segmentation, the range of the variables were scaled between 0-255,
 77 and also did not obviate the need for coregistration of geographic data. One
 78 advance in WW2 is the manner of building multidimensional models. The key
 79 insight is that the evaluation of a combination of two or more variables can
 80 be performed on the predictions of the models, instead of being performed in
 81 the high-dimensional model model space; a higher dimensional model is not
 82 required.

83 For example, the response on a single variable is called a membership func-
 84 tion (for reasons explained later). Prediction assigns a membership value to each
 85 data point in the training set, producing a vector of values in the range [0,1]
 86 which can be evaluated (using the ROC or AUC). The membership vectors for
 87 two single variables can be combined using a fuzzy AND operator to produce a
 88 new membership vector. This membership vector can then be evaluated (using
 89 the ROC or AUC), thus evaluating the performance of the conjunction with-
 90 out calculating the two-dimensional membership function. The combination of
 91 membership functions follows the AND, OR Zadeh operators [20]:

$$AND : x \wedge y = \min(f(x), f(y)) \quad (1)$$

$$OR : x \vee y = \max(f(x), f(y)) \quad (2)$$

92 This eliminates the need to develop and express a higher dimensional model,
 93 enabling a data mining approach where we can explore large databases for a
 94 parsimonious model of the species in analysis. Because of this we here focus on
 95 single variable models.

96 2.0.2 Choice of model structure

97 Ecological theory maintains that the response of species to the environment is
 98 generally 'humped' around an ideal, or restricted to a range of a variable (e.g.
 99 the range of temperature tolerances). Such a model must be at least quadratic
 100 to represent a unimodal response, and order three to incorporate skewness. An
 101 ideal method would be robust to any non-linear response type. Many environ-
 102 mental variables are categorical particularly soil and vegetation type. Ideally a
 103 modeling system integrates both types: continuous (e.g. temperature, rainfall)
 104 or categorical (e.g. biome or soil type) but very few do.

105 MaxEnt solves this problem by providing a range of potential response types
 106 [12]. The approach in **WhyWhere** is cutting up the range of the variable into
 107 discrete categorical factors. WW1 [18] used a color quantization algorithm in the
 108 GIF image format. WW2 segments a single continuous or categorical variables
 109 into open-closed intervals. For example, the output of the R cut function on the
 110 numbers 1..4 into 2 levels as open-closed intervals:

```
> cut(1:4,2)
```

-Submitted-

```
[1] (0.997,2.5] (0.997,2.5] (2.5,4]      (2.5,4]
Levels: (0.997,2.5] (2.5,4]
```

111 The unique 4 categorical values could also be represented in the same open-
112 closed syntax.

```
> cut(1:4,4)
```

```
[1] (0.997,1.75] (1.75,2.5] (2.5,3.25] (3.25,4]
Levels: (0.997,1.75] (1.75,2.5] (2.5,3.25] (3.25,4]
```

113 The cuts do not need to be uniformly spaced. The default method of de-
114 termining the cut locations in WW1 followed the median cut algorithm which
115 assigns equal numbers of points to a limited number of categories [7] in the
116 red, green, blue 3D space. This has been shown to retain good visual appear-
117 ance, but also has a statistical justification of minimizing the variance across the
118 range, by minimizing the variance in each category. There are other methods
119 of multi-interval discretization of continuous-valued attributes for classification
120 learning [4] implemented in the **discretization** R package.

121 The calculation of the intensity of species' response is straightforward on a
122 segmented variable. The locations where the species occurs can be thought of
123 as a sampling of the environmental space with a count of values in the each
124 environmental category (S), and expressed as a density by normalizing the sum
125 over the values in the categories to one. The prior density of values in each
126 environmental category is G . The change in the density from G to S indicates
127 strength of the the response of the species to its environment in that category –
128 the membership function M . A membership function is a fuzzy truth value as
129 a function $f : \mathbb{R} \rightarrow [0, 1]$ from a numeric domain to the real unit interval. The
130 membership function we used in WW1 and in WW2 for each category i in f is:

$$M_i = S_i / (S_i + G_i) \quad (3)$$

131 What we would like is the conditional probability of species being present
132 given the environmental $P(S|G)$ for each category of G . Frequentest calculations
133 give us $P(G|S)$. While $P(S|G)$ could be obtained using Bayes Theorem, it
134 requires an estimate of probability of the occurrence of a species $P(S)$, which in
135 opportunistic data is not well defined, and dependent on season, search effort and
136 other uncontrolled variables. The best we can do is an approximate equivalence
137 that holds under certain conditions (principally independence of variables) and
138 has been shown to be sufficient and useful in modelling such relationships [15,
139 16]. Where β is a normalization factor:

$$P(S|E) = \beta P(E|S) \quad (4)$$

140 One may ask why not calculate S/G and not $S/(S + G)$? Because we are
141 developing a model on a background set, developed from a random sample of
142 points in G , this can result in points in S_i that are not in G_i . That is, the
143 species occurs in environments that are not represented in the background set.

144 Use of $S/(S + G)$ avoids a divide by zero error. The membership function is a
145 proportional relationship which is sufficient to compare the efficiency of single
146 variables.

147 There are many ways to evaluate skill of a model. The approach to evaluating
148 the strength of the response in WW1 was significance with the Chi-squared test
149 or a K-S test, however we find it more convenient to use the area under the
150 curve of the receiver operating characteristic, or AUC. The receiver operating
151 curve (ROC) is widely used to compare classification models, while the AUC
152 provides a robust measure of skill. It is the probability that a model correctly
153 classifies a random draw of a positive and negative example.

154 2.0.3 Ecological interpretation

155 How does the structure of a multi-dimensional model embody the ecology of the
156 species? Many methods used in species modelling are based in other domains
157 (e.g. linear regression) without clear ecological interpretation. For example,
158 when inconsistent units such temperature and rainfall are combined in a gener-
159 alized linear model, how are they to be interpreted?

160 Ecological principles such as competitive interactions [13] tend to be in the
161 form of logical expressions. The principle of Liebig's Law of the Minimum states
162 that growth is controlled by the scarcest necessary resource. This is logically a
163 fuzzy conjunction of limiting factors – a Zadeh AND. Another law of ecology
164 is Gauss's law of competitive exclusion. This is a proposition that two species
165 competing for the same resource cannot coexist at constant population, due to
166 effect of slight advantages magnified over generations. This is a fuzzy disjunction
167 – a Zadeh OR. Thus fuzzy AND and OR can represent elements of established
168 ecological theory. The approach of modelling with a logical expression of a small
169 number of variables has utility in interpreting as ecological theory.

170 2.1 WhyWhere Algorithm

171 The first step is to 'presample' which when given a set of presence data locations,
172 and a geographic file that serves as a mask, produces a combined list of the
173 presence data and a randomly sampled list of sites of both presence and absence.
174 The mask file defines the geographic extent for sampling the background data.
175 If absences are available then they need not be generated.

176 The inputs to WW2 are: a **data.table** from presample with the longitude
177 and latitude of known locations, and a list of environmental data files that
178 may be read into the **raster** package. This file is then input to the main
179 routine with a list of the geographic files. Parameters include **multi** for searching
180 conjunctions of variables and **segment** to select the form of segmentation.

181 The algorithm proceeds by looping through the environmental variables and
182 creating and evaluating a membership function on each one. A table of the
183 variables so-far is retained. In the current implementation of a multi-dimensional
184 model, only the best variable is combined with each new variable using the fuzzy
185 minimum and the AUC recalculated. Alternative approaches to searching the

space of conjunctions may be implemented in future. It is possible to monitor the progress with the plot option. This plots out the best model so far and prints out a list of the best models considered in a streaming work flow.

Listing 1: Listing of the main algorithm

```

input locations
input a mask file
prepare background points and combine with presence points
for all environmental files do
  develop membership function for file
  insert AUC and variable name into ordered result
  test conjunct of this variable with best so far
  if result is better then insert into ordered result
output table of results

```

We show the results for the brown-throated three-toed sloth (*Bradypus variegatus*) that is documented and modeled in **dismo**. The environmental data consist of 9 environmental files from the WorldClim data set covering the South American continent. Figure 1 shows (A) the best variable, (B) the AUC of the best model, (C) the membership function, and (D) the predicted distribution with the presence points. The highest rated variable in the *Bradypus* data set is *bio7* = the temperature annual range (*bio5* – *bio6*) where *bio5* = maximum temperature of the warmest month and *bio6* = minimum temperature of the coldest month. The result is very similar to the results from GARP and MaxEnt in [12]. Table 2 lists the results for all variables.

Table 1: The AUC of environmental variables on the *Bradypus* data: WW is WhyWhere and GLM is a third order generalised linear model.

	name	WW	GLM
1	bio7	0.77	0.74
2	bio17	0.72	0.68
3	bio12	0.72	0.70
4	bio6	0.71	0.66
5	bio16	0.68	0.66
6	biome	0.66	0.65
7	bio5	0.64	0.62
8	bio1	0.61	0.58
9	bio8	0.59	0.57
10	limit	0.00	0.00

The membership function shown graphically in Figure 1.C is represented internally as a lookup table, shown in Table2. The prior distribution (*G1*) is the blue line, and the distribution of presences (*G2*) is the red line. The membership function is shown by the grey bars ($G1/(G1 + G2)$). Note the almost uniform distribution of background classes in the variable width quantile cuts.

-Submitted-

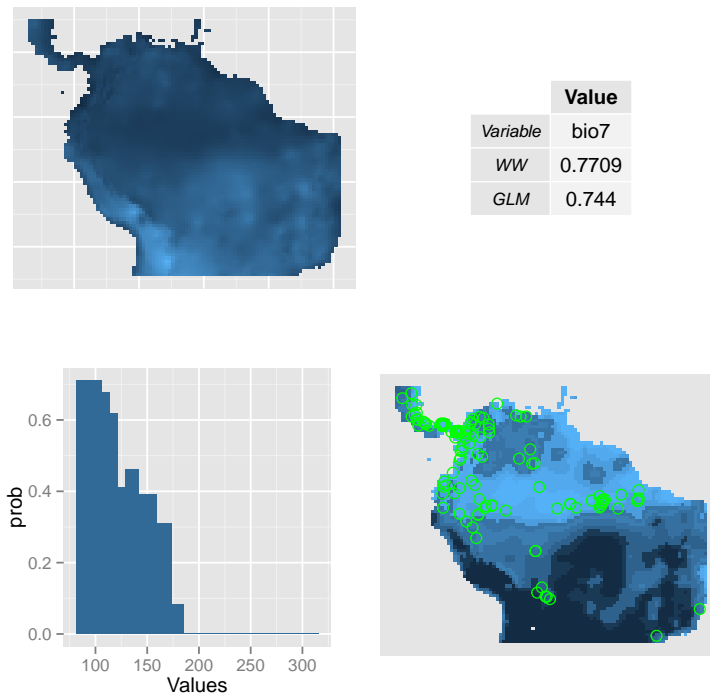


Figure 1: Output of WW2 on the Bradypus data set. (A) the best variable, (B) the AUC of the best model, (C) the membership function, and (D) the predicted distribution with the presence points in green.

Table 2: The lookup table for the membership function on the Bradypus data. Each row is a segment of the range of the variable (factors). The density of points in each segment in the background (g1) and the points where species occur (g2) is used to calculate the odds ration, and then membership in each segment. The width is the size of the segment interval

	factors	levels	g1	g2	odds	prob	width
1	(81,106]	1	0.12	0.30	2.46	0.71	25.00
2	(106,114]	2	0.12	0.26	2.09	0.68	8.00
3	(114,121]	3	0.11	0.17	1.62	0.62	7.00
4	(121,129]	4	0.09	0.06	0.69	0.41	8.00
5	(129,142]	5	0.10	0.09	0.85	0.46	13.00
6	(142,159]	6	0.09	0.06	0.64	0.39	17.00
7	(159,174]	7	0.10	0.04	0.45	0.31	15.00
8	(174,185]	8	0.09	0.01	0.09	0.08	11.00
9	(185,199]	9	0.08	0.00	0.00	0.00	14.00
10	(199,316]	10	0.09	0.00	0.00	0.00	117.00

3 Experiments

3.1 K-fold Validation

K-fold validation provides a robust evaluation of the accuracy of a method on independent data. The data was prepared by generating a data set using presample, and then labeling each row from one to 5. K-fold validation proceeds by sequentially holding back one fifth of the data each time for evaluation, and developing the model using the remaining four-fifths. Table 3 shows the results of the five-fold validation.

Table 3: Results of a five-fold validation of Bradypus model. WW is WW2 and GLM is generalised linear regression. the suffix tr is the AUC on the training set and te is the AUC on the test set.

	name	WWtr	GLMtr	WWte	GLMte
1	bio7	0.77	0.74	0.76	0.78
2	bio7	0.77	0.75	0.75	0.75
3	bio7	0.77	0.74	0.78	0.79
4	bio7	0.77	0.75	0.75	0.75
5	bio7	0.77	0.74	0.78	0.79
6	mean	0.77	0.74	0.76	0.77
7	sd	0.00	0.01	0.02	0.02

The accuracies were similar on test and training sets and between the WW2 and the GLM models, and the best variable *bio7* was chosen consistently. The performance of WW2 not dissimilar to the GLM.

226 3.2 Multi-dimensional option

227 Table 4 lists the results for multidimensional models developed by combining
 228 the prediction of two or more variables with a fuzzy AND operator and then
 229 evaluating the resulting AUC.

Table 4: The AUC of environmental variables on multi-dimensional prediction of the Bradypus data: WW is WhyWhere and GLM is a third order generalised linear model.

	name	WW	GLM
1	bio7.biome	0.79	0.90
2	bio7	0.78	0.75
3	bio12	0.73	0.70
4	bio6	0.72	0.67
5	bio17	0.71	0.68
6	bio12.bio6	0.70	0.69
7	bio12.bio17	0.70	0.69
8	bio16	0.69	0.66
9	bio12.bio16	0.68	0.66
10	biome	0.66	0.64
11	bio5	0.64	0.62
12	bio1	0.61	0.59
13	bio8	0.60	0.57
14	bio12.bio5	0.51	0.71
15	bio7.bio8	0.50	0.66
16	limit	0.00	0.00

230 The best result is a combination of *bio7* and *biome* variables (the AUC of
 231 WW=0.7969 and GLM=0.906). The variable *biome* is a categorical variable
 232 expressing ecosystem type, and so more like a habitat variable than a numeric
 233 climatic variable.

234 3.3 Alternate segmentation

235 We also evaluated some alternative methods of segmenting the response function,
 236 shown on Table 5: an even distribution of cuts over the range of the
 237 variable, distribution by quantile frequency, and an entropy optimizing method.
 238 The even method segments the variable evenly over the range. The quantile
 239 method segments the variable so that as far as possible each segment contains
 240 an equal number of data.

241 The number of segments is determined using the Freedman-Diaconis rule
 242 [5] for optimal binning of histograms. The entropy method uses the R package
 243 **discretization** and the routine `cutPoints()` that perform cuts for the Minimum
 244 Description Length Principle. This analysis is not a comprehensive evaluation
 245 of discretization method, but serves to validate the performance of the quantile
 246 approach in comparison of some readily available alternatives.

Table 5: Results of a five-fold validation of these data with other segmentation approaches: even, quantile and entropy.

	name	WWtr	GLMtr	WWte	GLMte	segment
1	mean	0.77	0.74	0.76	0.77	even
2	mean	0.77	0.74	0.76	0.77	quantile
3	mean	0.75	0.72	0.72	0.74	entropy

247 The segmentation methods that performs best on the test set are the even
 248 and quantile methods, and their performance is similar on the training set.
 249 The entropy method performed less well. This provides support for the quan-
 250 tile approach that was used in WW1, although further testing of segmentation
 251 approaches may lead to improvements.

252 3.4 Prediction on World Data

253 To demonstrate the system on a larger data set we use the Global Ecosystems
 254 Database, a set of 940 global data sets of environmental variables, previously
 255 prepared and used in the WW1. The multi-agency distribution of the original
 256 CD includes many groups of variables listed in [17] including satellite green-
 257 ing, monthly temperature and rainfall and many others in a range of different
 258 resolutions in raster and vector formats.

Table 6: The top ten variables in a single variable WhyWhere model of *Bradypus* using the 940 variable World dataset. AUC is the WW AUC and BAUC is the AUC from GLM.

	name	AUC	BAUC
1	fnocwat	0.85	0.80
2	srzsoil	0.85	0.68
3	wrzsoil	0.85	0.58
4	wrroot	0.84	0.71
5	lcprc08	0.84	0.74
6	lwerr05	0.84	0.76
7	wrcia01	0.84	0.57
8	wrcia03	0.84	0.54
9	SALINITY_ANN_AVG	0.84	0.68
10	wrsil02	0.83	0.65

259 Table 6 lists the top 10 variables identified by the algorithm in predicting
 260 *Bradypus* on the World dataset. Out of 940 variables, the top variables were
 261 *fnocwat*: Navy Terrain Data - Percent Water Cover. The next three variables
 262 are soil classifications – *srzsoil*: Staub and Rosenzweig Zobler Soil Units, and
 263 Webb et al Soil Particle Size Properties Zobler Soil Types. The fifth variable is
 264 climatic: Leemans and Cramer August Precipitation (mm/month) which corre-
 265 sponds to the dry season in Amazonia. Note that the accuracy of the GLM is

less than the WW2 in this case.

It might be inferred from this limited study that habitat features have greater predictive power than climatic variables over the region of distribution for this species. By way of interpretation, *Bradypus variegatus* does leave the trees in search of food and while it crawls along the forest floor poorly, it does swim well [6]. Its distribution may be closely linked to the flooded forest (defined as a seasonal inundation of the forest floor) facilitating access to other trees for food. The identification of soil classification may be indicative of the soils supporting a the flooded forest ecosystem.

Habitat variables are proximal causes of species which necessarily produce higher accuracy than the more distal climatic variables. It is well known that habitat features are crucial in identifying suitable areas of land for the conservation of threatened species, and due to the proximal relationship should be a more important determinant of species decline than distal factors such as climate change. This is not a definitive examination of determinants of *Bradypus*, but serves to illustrate the potential expositions available from this approach.

4 Discussion

This study evaluated the accuracy of a segmented model and algorithm in an updated version of **WhyWhere** against a generalized linear model, and also modelling species response to climate variables and accuracy on a large data set containing mixtures of continuous and categorical environmental data.

The accuracies were similar on the *Bradypus* data set and **WhyWhere** was superior on the large World dataset, attributed to selection of the best WW2 variables and handling of categorical as well as continuous variables. The benefit of WW2 are more accurate species prediction, and potential insight into proximal cause of the species occurrences. The results verify the findings of the previous version of the **WhyWhere**, showing progress that could be made in modelling species response to the environment by using segmented models of few variables mined from large databases of environmental variables.

The recoding of **WhyWhere** into R has greatly improved the programs' utility. Refinements to the algorithm reduce the steps in the species modelling workflow and support more efficient higher dimensional models using novel fuzzy set operators.

It is interesting to note that that the two variables identified in the multi-variable mode are habitat variables (*biome*) and a climate variable (*bio7*). We hypothesize that climate and habitat factors are independent causal factors that together determine species distribution, and that the multi-dimensional WW2 can correctly identify such independent determinants of species response, yielding a parsimonious explanation of the species' response to its environment.

More complex expressions of species distribution models in the R package **WhyWhere** are planned, as are improvements in computational efficiency and testing over a greater range of higher resolution environmental data.

References

- [1] M Austin and J Meyers. Current approaches to modelling the environmental niche of eucalypts: Implication for management of forest biodiversity. *Forest Ecology And Managment*, 85(1-3):95–106, 1996.
- [2] S M Davey and D R B Stockwell. Incorporating Wildlife Habitat Into An Ai Environment - Concepts, Theory, And Practicalities. *AI Applications*, 5(2), 1991.
- [3] M Dowle, T Short, S Lianoglou, and A Srinivasan. data.table: Extension of data.frame. 10 2014.
- [4] U. M. Fayyad and K. B. Irani. Multi-interval discretization of continuous-valued attributes for classification learning. *Artificial intelligence*, 13:1022–1027, 1993.
- [5] David Freedman and Persi Diaconis. On the histogram as a density estimator: a theory. 57(4):453–476–, 1981.
- [6] Virginia Hayssen. *Bradypus variegatus* (pilosa: Bradypodidae). *Mammalian Species*, 42(1):19–32, 2010.
- [7] Paul Heckbert. Color image quantization for frame buffer display. In *SIGGRAPH ’82: Proceedings of the 9th annual conference on Computer graphics and interactive techniques*, pages 297–307, New York, NY, USA, 1982. ACM Press.
- [8] Robert J Hijmans, Susan E Cameron, Juan L Parra, Peter G Jones, and Andy Jarvis. Very high resolution interpolated climate surfaces for global land areas. *International journal of climatology*, 25(15):1965–1978, 2005.
- [9] Leo Joseph and David Stockwell. Climatic modeling of the distribution of some pyrrhura parakeets of northwestern south america with notes on their systematics and special reference to pyrrhura caeruleiceps todd, 1947. *Ornitologia Neotropical*, 13:1–8, 2002.
- [10] H Nix. A biogeographic analysis of Australian Elapid snakes. In R Longmore, editor, *Atlas of Australian Elapid Snakes*, volume 8, pages 4–15. Australian National University, 1986.
- [11] David M Olson, Eric Dinerstein, Eric D Wikramanayake, Neil D Burgess, George VN Powell, Emma C Underwood, Jennifer A D’amico, Illanga Itoua, Holly E Strand, and John C Morrison. Terrestrial ecoregions of the world: A new map of life on earth a new global map of terrestrial ecoregions provides an innovative tool for conserving biodiversity. *BioScience*, 51(11):933–938, 2001.
- [12] Steven J Phillips, Miroslav Dudík, and Robert E Schapire. Maxent software for species habitat modeling, 2007.

- 346 [13] Víctor Sánchez-Cordero, David Stockwell, Sahotra Sarkar, Hawoei Liu,
347 Christopher R Stephens, and Joaquín Giménez. Competitive interactions
348 between felid species may limit the southern distribution of bobcats lynx
349 rufus. *Ecography*, 31(6):757–764, 2008.
- 350 [14] D Stockwell and D Peters. The GARP modelling system: problems and
351 solutions to automated spatial prediction. *International Journal Of Geo-*
352 *graphical Information Science*, 13(2):143–158, 1999.
- 353 [15] D R B Stockwell. Lbs - Bayesian Learning-System For Rapid Expert
354 System-Development. *Expert Systems With Applications*, 6(2), 1993.
- 355 [16] D R B Stockwell. Generic predictive systems: An empirical evaluation
356 using the Learning Base System (LBS). *Expert Systems With Applications*,
357 12(3), 1997.
- 358 [17] David R B Stockwell. *Ecological Niche Modeling: Ecoinformatics in Appli-*
359 *cation to Biodiversity*, chapter 7. CRC Press, 2006.
- 360 [18] David R B Stockwell. Improving ecological niche models by data mining
361 large environmental datasets for surrogate models. *Ecological Modelling*,
362 192:188–196, 2006.
- 363 [19] Robert J. Hijmans & Jacob van Etten. raster: Geographic analysis and
364 modeling with raster data, 2012. R package version 2.0-12.
- 365 [20] L.A. Zadeh. Fuzzy sets. *Information and Control*, 8(3):338 – 353, 1965.