

**MN341 Species distribution models**  
**Dr Peter Long, University of Oxford**

**Introduction**

Species distribution modeling is a technique for linking spatially referenced records of species occurrence – for example collected during appropriately designed field-based biodiversity monitoring programmes – with maps of environmental variables such as elevation, climate, vegetation or human disturbance, in order to create a statistical model of the relationship between a species and its environment, ie the species realised ecological niche. GIS can then be used to express the results of models as habitat suitability maps across a desired spatial extent. These output habitat suitability maps are also a powerful tool to communicate conservation messages to non-scientists

The species records may be a set of presences only or a set of presence and absence records, depending on the detectability of the species and sampling method that has been used. Almost any environmental variables can be used in a distribution model, although it is normal to select a restricted set of variables at a particular spatial scale based on a working hypothesis about the aspects of the environment which may be important to the focal species.

A wide range of statistical approaches have been developed for fitting distribution models including various types of regression models, machine learning and classification methods. Regression approaches such as GLM are simple to implement; however, more complex information theoretic approaches, especially maximum entropy, have proved to be very powerful. Once a model has been built and refined, it is critically important to validate the model on a subset of data which was not used in model construction, in order to objectively assess how well the model performs. Please see the slides at the end of this document for an overview of the modeling process.

**Designing your project**

Species distribution models is a popular dissertation topic at the Madagascar site. There is scope for great flexibility in the ways that this technique can be used and there are many possible extensions and additional analyses that are possible.

Should you decide to do a dissertation on SDM, I recommend focusing on just one of the following taxonomic groups:

- flowering plants
- birds
- reptiles and amphibians
- mammals

Some brief notes to help you in this difficult choice...

### *Flowering plants*

A very diverse group with a couple of hundred species in nearly 60 families recorded from Mahamavo, some of which have very restricted ranges and are globally threatened. Almost every single species is endemic to Madagascar. Many of the flowering plants also have medicinal uses. The flora is predominantly woody trees and shrubs although there are some succulents (Euphorbiaceae) and orchids (Orchidaceae) too. The dry forest is dominated by Fabaceae and Combretaceae but there are also Baobabs (Malvaceae), Figs (Moraceae), Palms (Arecaceae), Mangroves (Rhizophoraceae and others). A major advantage of distribution models of plants is that individuals are usually much more detectable, although often harder to identify.

### *Birds*

The forest bird fauna is very rich in Mahamavo: over 100 species have been recorded here. Our scientists have spotted several groups of Van Dam's vanga (*Xenopirostris damii*), a globally endangered bird species which has hitherto only ever been recorded from Ankarafantsika national park and Analamera forest. This is a very significant range extension for this species as it is the rarest and most threatened species in the whole Vanga family. The other really notable forest birds in the site are Madagascar harrier hawk (*Polyboroides radiatus*), an indicator of forest in excellent condition, and Crested Ibis (*Lophotibis cristata*), a vulnerable and unusual forest ibis. The main method is repeated 10min point counts at ~100 sample sites in the forest, plus opportunistic sampling in lots of other places.

There are around 50 forest bird species, but the following abundant 20 species account for 3/4 individuals detected by point counts): Souimanga sunbird (*Nectarinia souimanga*), Madagascar bulbul (*Hypsipetes madagascariensis*), Madagascar bee-eater (*Merops superciliosus*), Palm swift (*Cypsiurus parvus*), Paradise flycatcher (*Terpsiphone mutata*), Madagascar magpie-robin (*Copsychus albospectularis*), Common newtonia (*Newtonia brunneicauda*), Grey-headed lovebird (*Agapornis canus*), Fork-tailed drongo (*Dicrurus forficatus*), Long-billed greenbul (*Phyllastephus madagascariensis*), Pied crow (*Corvus albus*), Madagascar buzzard (*Buteo brachypterus*), Madagascar red fody (*Foudia madagascariensis*), Lesser vasa parrot (*Coracopsis nigra*), Crested coua (*Coua cristata*), Yellow-tailed kite (*Milvus aegypticus*), Chabert's vanga (*Leptopterus chabert*), Greater vasa parrot (*Coracopsis vasa*), Long-billed green sunbird (*Nectarinia notata*), Madagascar coucal (*Centropus toulou*).

We also work on wetland birds, conducting extensive surveys using a speedboat boat and visiting lakes on foot. You will see lots of Madagascar malachite kingfishers (*Alcedo vintsioides*), Anhingas (*Anhinga rufa*), Openbills (*Anastomus lammelligerus*) Herons (*Ardea cinerea*, *A. purpurea*, *Ardeola ralloides*, *Butoroides striatus*, *Nycticorax nycticorax*), Egrets (*Bulbulcus ibis*, *Egretta alba*, *E. ardesiaca*, *E. dimorpha*). Raptors (*Buteo brachypterus*, *Milvus aegypticus*, *Falco newtoni*, *Polyboroides radiatus*), Shorebirds (*Callidris ferrunginea*, *C. minuta*, *Charadrius hiaticula*, *C. tricollaris*, *Numenius phaeopus*), Swifts (*Cypsiurus parvus*), Ducks (*Dendrocygna viduata*, *D. bicolor*) Rails (*Dryolimnas cuvieri*, *Ixobrychus minutus*, *Gallinula chloropus*), Stilts (*Himantopus himantopus*), Cormorants (*Phalacrocorax africanus*), Flamingoes (*Phoenicopterus ruber*), Spoonbills (*Platalea alba*), Ibises (*Threskiornis bernieri*, *Plegadis falcinellus*), and Grebes (*Tachybaptus pelzelni*, *T. rufficollis*).

### *Reptiles and amphibians*

We have a very cool herp assemblage in Mahamavo with a species list of over 80, and they're very abundant in the forest. We repeatedly walk routes in the day and at night and sample herps opportunistically. You will encounter chameleons (*Furcifer oustaleti*, *F. angeli*), iguanas (eg *Oplurus cuvieri*), plated lizards (eg *Zonosaurus laticaudatus*), skinks (eg *Trachylepis elegans*), day geckos (eg. *Phelsuma madagascariensis*, *Geckolepis typica*), boas (*Acrantophis madagascariensis* and *Sanzinia madagascariensis*), a lot of different colubrid snakes (eg *Leiheterodon madagascariensis*, *Dromicodryas quadrilineatus*), lots of frogs (*Boophis* spp., *Heterixalus* spp., *Hoplobatrachus* sp. *Ptychadena* spp.), night geckos (eg *Parodura stumpfii*) and the legendary leaf-tailed geckos (*Uroplatus henkeli*, *U. ebenau*, *U. guntheri*). In the wetlands we regularly observe crocodiles (*Crocodylus niloticus*) and turtles (*Pelomedusa subrufa*). There are probably also marine turtles in the seagrass beds.

### *Mammals*

The aim of repeatedly walking lemur routes in the day and at night is to estimate encounter rates, densities and population sizes of diurnal lemurs using distance sampling. The commonest diurnal lemurs is Coquerel's sifaka (*Propithecus coquereli*), but you will also see common brown lemur (*Eulemur fulvus*), mongoose lemur (*Eulemur mongoz*), At night there are hundreds of mouse lemurs; golden-brown mouse lemur (*Microcebus ravelobensis*) and grey mouse lemur (*Microcebus murinus*). You might also see Western Avahi (*Avahi occidentalis*), Milne-Edwards sportive lemur (*Lepilemur edwardsi*), Western fork-marked lemur (*Phaner pallescens*), and Fat-tailed dwarf lemur (*Cheirogaleus medius*).

Other terrestrial mammals which are detected opportunistically include carnivores (*Cryptoprocta ferox*, *Eupleres goudotti*, *Galidia elegans*), wild pigs (*Potamochoerus larvatus*) and larger tenrecs (eg *Setifer setosus*). Small mammals trapped by pitfalls include Malagasy mice, tuft-tailed rats (*Eliurus myoxinus*), smaller tenrecs (eg *Microgale brevicaudata*). Fruit bats (*Pteropus rufus*) are observed opportunistically, and a number of species of microbats are caught in mist nets.

### **Extensions**

Additionally you need to decide whether to undertake a straightforward analysis of a set of species or whether you want to use the results for further analysis, in which case you would need to decide where to place the emphasis in your analysis. You may wish to use distribution models to look at species responses to landscape configuration, or use the results to explore patterns of biodiversity and even hence undertake systematic conservation planning. Maybe you might decide to investigate covariate spatial scale effects by comparing models made with, say 30m and 300m resolution data. A related approach is to create focal covariates to explore the scale at which each species responded to landscape features. You might consider looking at the effects of historical environmental changes or future change (eg climate change, land cover change) on species distributions. Perhaps you would be interested to contrast the performance of several different modeling techniques such as Maxent, GLM, ENFA. Alternatively you might want to look at the feasibility of assimilating new data (eg RS) to generate indicators eg of area of occupancy, for use as an aid to monitoring. It could also be possible to generate

detection histories for sample units, model detectability and correct distribution models for imperfect detections. A final, very neat possibility is to use records for Mahamavo from GBIF to make distribution models and compare estimated attributes such as patterns of species richness estimated from field data with that from web databases or to evaluate the accuracy of polygons denoting the global ranges of endangered species and investigate statistical range polygon refinement methods.

These more advanced topics may sound quite complicated, so please don't hesitate to get in touch to discuss further and ask questions

### **Practical aspects: fieldwork, model type, covariates**

In practical terms, I would recommend that when you come to Mahamavo, you join in with one of the biodiversity survey teams and help with general data collection for a few weeks, then query our field database (from 2009, 2010, 2011, 2012) to generate your dataset, which you then join with whichever spatial data you want to use. I maintain all the spatial datasets for Mahamavo, so you will not need to prepare these (covariate preparation is normally the hardest part of distribution modelling), although if there's extra covariates that you want to incorporate, I can show you how to do this.

Full training will be given in the necessary data handling and the complete modeling procedures on site, including validation statistics, cartography etc.

Whichever taxonomic group you choose and regardless of whether you want to pursue one of the many possible extensions, there will be lots of choices to be made about the modelling approach to take in the first place. I recommend using either GLM, a regression method, and/or Maxent, an information-theoretic method. GLM is conceptually simplest, fairly straightforward to do and usually gives very good results. The mathematical basis of Maxent is harder to understand, but in practical terms, it's very easy to make models and they generally perform slightly better than GLM.

You will need to decide on the spatial extent and the scale (spatial grain) of your modelling. (eg the whole Mahamavo watershed at 30m resolution), and also think about the a set of environmental covariates to use which will be relevant to your taxonomic group.

Available covariates at both 30m and 300m scale for Mahamavo include:

*Reflectance derived* (from Landsat 5 and & at 30m, from MERIS and MODIS at 300m)

NDVI –normalised difference vegetation index

EVI – enhanced vegetation index

Tasseled cap greenness, moistness, brightness – indices of healthy green vegetation, moisture availability, bare soil

BRDF volumetric and geometric scattering kernel parameters – proxies for canopy architecture

Texture metrics – proxies for habitat heterogeneity

*Elevation derived (from SRTM)*

Elevation

Slope

Aspect and  $\sin(\text{aspect})$  and  $\cos(\text{aspect})$  as aspect is a circular variable

TWI - Topographic wetness index

Relative insolation

*Distance metrics (from 1:100k topo maps)*

Distance to water

Distance to road

Distance to village

*Configuration metrics (from a classification of*

Patch size

Patch perimeter:area ratio

Distance from patch edge

Patch isolation

*Soils (from 1:100k soils maps)*

Soil class

*Climate (from WORLDCLIM station interpolation and also RS observed from TRMM, DMSP-SSMI, AIRS etc)*

Several climate parameters are available eg mean annual temp, total annual precipitation but please note that the study area is relatively small, with limited elevation range. As such most climate parameters show very limited variation across the study landscape.

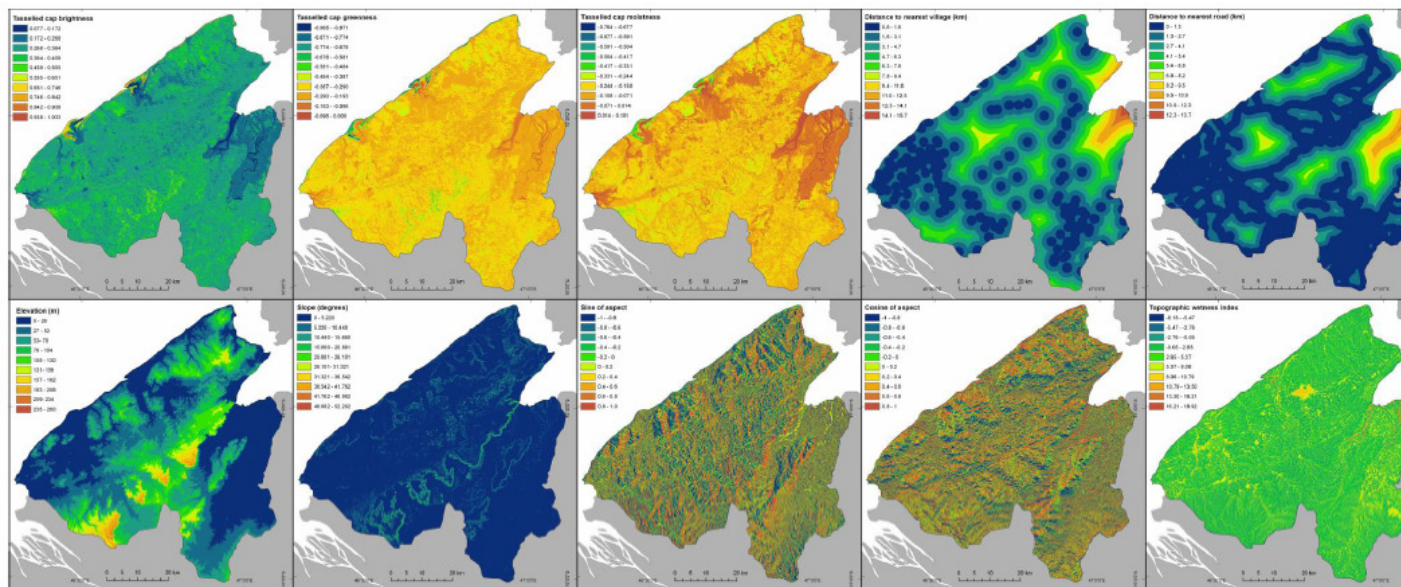


Figure. Maps of some covariates for Mahamavo at 30m resolution: TC1,2,3, village distance, road distance, elevation, slope cos(aspect), sin(aspect), TWI.

You will have to think carefully about the organisms in your chosen group to select an appropriate set of covariates. For example, for reptiles, relative insolation is important for thermoregulation, whereas it is unlikely to matter to birds and mammals.

Some of these variables are static, such as elevation, whereas others like reflectance and configuration change with phenology and land cover change. It is possible to create covariates to capture these effects, such as the first and second principal components of a monthly time series of NDVI, EVI or TC1 to capture mean greenness and seasonality in greenness.

At this stage, please remember that the final way you do your project will be shaped by your reading and your interests and you don't actually need to decide the answers to all of these issues immediately.

### Reading:

Fielding, Bell (1997) A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environmental Conservation* 24 (1): 38–49

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

Guisan, Zimmerman (2004) Predictive habitat distribution methods in ecology. *Ecological Modeling* 135: 147-186

Hirzel, Guisan (2002) Which is the optimal sampling strategy for habitat suitability modeling? *Ecological Modelling* 137: 331-341

Liu et al (2005) Selecting thresholds of occurrence in the prediction of species distributions. *Ecography* 28: 385-393

Pearson (2004) Modeling species distributions in Britain: a hierarchical integration of climate and landcover data. *Ecography* 27: 285-298

Franklin (2009) Mapping species distributions. Cambridge

Phillips (2008) Maxent handbook

Moilanen, Wilson, Possingham (2009) Spatial conservation prioritisation. Oxford

Margules & Sarkar (2007) Systematic conservation planning. Cambridge

Moilanen & Kujala (2008) Zonation manual

--

Dr Peter Long  
Department of Zoology  
University of Oxford  
South Parks Road  
Oxford  
OX1 3PS

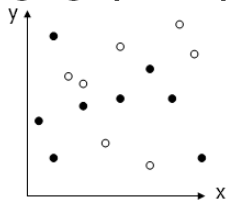
[peter.long@zoo.ox.ac.uk](mailto:peter.long@zoo.ox.ac.uk)  
01865 281878

--

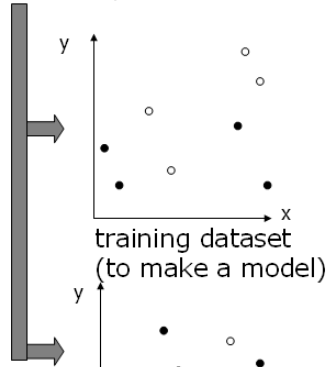
Last updated 8th September 2011

## SDM: biodiversity data requirements

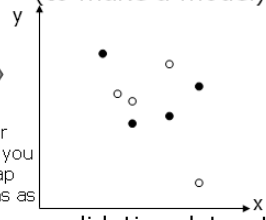
1. Collect species presence (and absence) data in **geographical space**.
2. Randomly divide the data into partitions:



N.B. Make sure you sample across major environmental gradients in the study area eg. By stratifying sampling by elevation.



training dataset  
(to make a model)

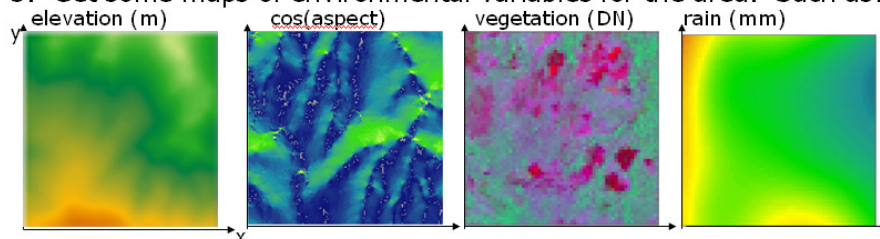


validation dataset  
(to test the model)

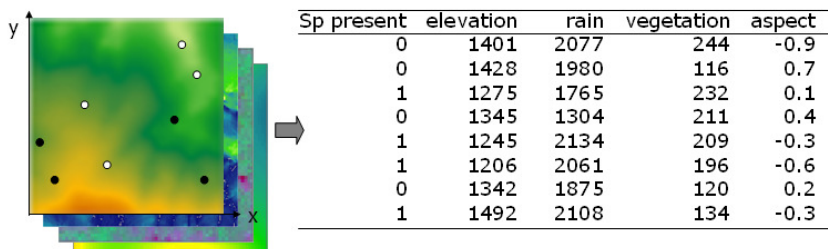
N.B. 2 partitions in this example for simplicity, but in k-fold partitioning you might use 10 or 20. If you bootstrap then there will be as many partitions as records.

## SDM: Environmental variables

3. Get some maps of environmental variables for the area. Such as:



4. Overlay the training dataset on each of these maps and extract the environmental information at each point using GIS



## SDM: modelling

5. Make and refine statistical model for the probability of the species occurring in a landscape unit:

$$\text{Pr}(\text{species present in landscape unit}) \sim x_1 + x_2 + x_3 \dots + x_n$$

Lots of ways: GLM, GAM, entropy maximising algorithms, regression trees etc

Eg. Use a GLM to estimate the parameters in an equation of this form:

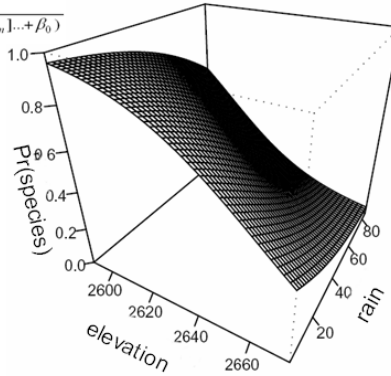
$$\text{Pr}(\text{species}) = \frac{e^{(\beta_1[x_1] + \beta_2[x_2] + \beta_3[x_3] \dots + \beta_n[x_n] + \beta_0)}}{1 + e^{(\beta_1[x_1] + \beta_2[x_2] + \beta_3[x_3] \dots + \beta_n[x_n] + \beta_0)}}$$

The actual model output looks like this:

	Beta	Std. Error	z value	Pr(> z )
(Intercept)	36.275716	17.331823	2.093	0.0363 *
elevation	0.006937	0.002977	2.330	0.0198 *
rain	-0.025242	0.011358	-2.222	0.0262 *

Which means:

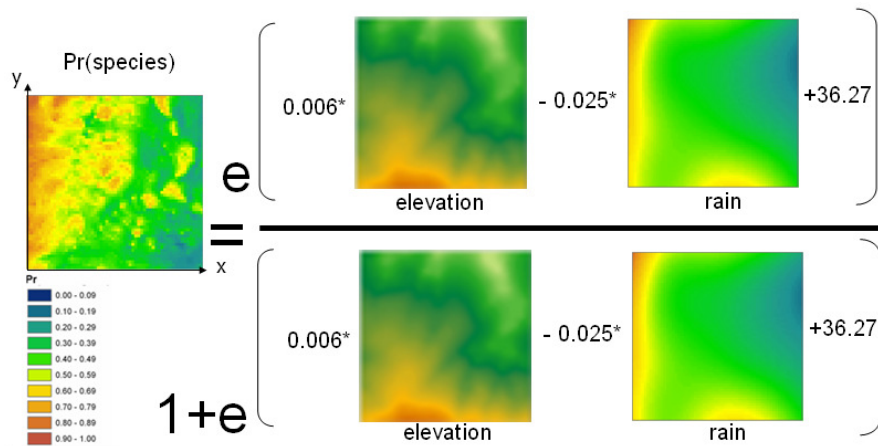
$$\text{Pr}(\text{species}) = \frac{e^{(0.006 * \text{elevation} - 0.025 * \text{rain} + 36.27)}}{1 + e^{(0.006 * \text{elevation} - 0.025 * \text{rain} + 36.27)}}$$



## SDM: Map algebra

6. Use map algebra to evaluate the equation for the species distribution model for every landscape unit to produce a habitat suitability map in **geographical space**

$$\text{Pr}(\text{species}) = \frac{e^{(0.006 * \text{elevation} - 0.025 * \text{rain} + 36.27)}}{1 + e^{(0.006 * \text{elevation} - 0.025 * \text{rain} + 36.27)}}$$



## SDM: validation

7. Overlay the validation data (which was not used to make model) on the habitat suitability map. Interrogate the map to find the modelled probability of occurrence in sampled locations where the true occurrence is known. Plot a ROC curve.

