

Ecological niche models and species distributions

Jesús N. Pinto-Ledezma and Jeannine Cavender-Bares

In this lab, we will explore some Correlative models using data at species level (occurrence data points) and environmental data at large spatial scales. We will use occurrence data of the Hyacinth macaw (*Anodorhynchus hyacinthinus*), the largest macaw in the world and that is distributed in the Pantanal of Bolivia and Brazil and the Cerrado in Brazil. The environmental was obtained from from **Ecoclimate** (<http://ecoclimate.org>). Notice that you can obtain environmental data from WorldClim.

Part of the explanation of algorithms was extracted from the vignette of the dismo package.

Set up your data and your working directory

Set up a working directory and put the two data files in that directory. Tell R that this is the directory you will be using, and read in your data:

```
setwd("path/for/your/directory")
```

Install and load the following packages.

```
library(vegan) # Community Ecology Package
require(raster) # Geographic Data Analysis and Modeling
library(sp) # Classes and Methods for Spatial Data
library(psych) # Procedures for Psychological, Psychometric, and Personality
# Research
library(rgeos) # Interface to Geometry Engine - Open Source ('GEOS')
library(maps) # Draw Geographical Maps
library(maptools) # Tools for Handling Spatial Objects
library(dismo) # Species Distribution Modeling
library(rworldmap) # Mapping Global Data
library(kernlab) # Kernel-Based Machine Learning Lab
library(corrplot) # Visualization of a Correlation Matrix
```

Prepare species data

Get some occurrence data for our species from GBIF, directly within R. This may take some time, given the number of occurrences for the selected species.

NOTE: we need to have an internet connection.

```
hyacinth <- gbif("anodorhynchus", "hyacinthinus*", geo = TRUE, removeZeros = TRUE)
```

Inspect the results.

```
# how many rows and columns?
dim(hyacinth)
## [1] 4084 137
```

```
# select the records that have longitude and latitude data
colnames(hyacinth)
```

Wow, there are a bunch of columns, lets clean our data.

First, remove all NA data.

```
hyacinth <- subset(hyacinth, !is.na(lon) & !is.na(lat))
```

```
dim(hyacinth)
## [1] 3785 137
# show some values
hyacinth[1:4, c(1:5, 7:10)]
```

Second, select only the columns that are necessary for our modeling purposes. Notice that using `colnames(hyacinth)` you can obtain the position of the columns, and the columns 120, 90 and 83 (in that order) are the columns that correspond to species name, longitude and latitude, respectively and that is all what we need.

```
hyacinth_coords <- hyacinth[, c(120, 90, 83)]
```

Lets inspect the results.

```
head(hyacinth_coords)
```

Now let's save the data needed for further analyses.

```
dir.create("Data/OCC")
write.csv(hyacinth_coords, "Data/OCC/hyacinth_data.csv")
```

We can clean our R environment.

```
rm(list = ls())
```

Prepare environmental data

Read the environmental variables

```
datoG0 <- read.table("Data/Environment/bio_var_CCSM_0k_global.txt", h = T)
```

Lets explore some details of the raw data

```
head(datoG0)
datoG0[1:5, 1:5]
dim(datoG0)
str(datoG0)
summary(datoG0)
class(datoG0)
names(datoG0)
```

Note that the environmental data is a data frame, but we need a raster. To transfer the data frame to raster, we just need to use the function **gridded**. Lets try!

```
gridded(datoG0) <- ~long+lat
```

```
class(datoG0)
```

Now use the function **stack** to merge all environmental variables into a single raster

```
clima0k <- stack(datoG0)
```

Lets plot one of the variable to check if everithing is ok.

```
plot(clima0k$bio.1)
map(add = T)
```

This map correspond to the environmental varibale Annual Mean Temperature or BIO-1 accross the world.

However, we do not need all environmental information across the world, as our focal species is ditributed only in South America (if the species is in a Zoo in US or Europe doesn't mean that the species occur in those places, dah!).

```
# Set a geographical extension (South America long = -90, -30, lat = -60, 15)
e <- extent(c(-90, -30, -60, 15))
```

Using the extent cut the information only for South America

```
clima0k.SA <- crop(clima0k, e)
```

Check if everything is ok.

```
plot(clima0k.SA$bio.9)
map(add = T)
```

```
ncell(clima0k.SA)
```

The next step is to extract the information from the raster of South America.

```
SA0k <- values(clima0k.SA)
```

```
class(SA0k)
```

Obtain the the geographical coordinates for each cell in South America.

```
coord.SA <- xyFromCell(clima0k.SA, 1:ncell(clima0k.SA))
```

Now, merge the environmental values with the geographical coordinates.

```
SA0k <- data.frame(cbind(coord.SA, SA0k))
```

Check the values on specific rows and columns.

```
# Observ specific rows and columns
SA0k[1:5,]
SA0k[1:5, 1:5]
```

```
# Save the selected variables.
write.table(SA0k, row.names = FALSE, "Data/Environment/climaSA0k_all.txt", sep = "\t")
```

We can also save each environmental variable as raster files in our computer, but we do not need to do that.

```
writeRaster(clima0k.SA, "Data/Environment/climaSA0k.asc", format = "ascii", bylayer = TRUE, overwrite =
```

Now we have the environmental information needed for modeling the fundamental niche of the Hyacinth macaw. So, we can clean our R environment.

```
rm(list = ls())
```

Merge both information

First load the occurrence data for our species.

```
hyacinth <- read.csv("Data/OCC/hyacinth_data.csv")
```

As we are working only with one species, we just need the lon/lat columns.

```
head(hyacinth)
```

The columns 3 and 4 correspond to the geographical coordinates.

```
hyacinth_coords <- hyacinth[, c(3, 4)]
```

```
head(hyacinth_coords)
```

Now, load the environmental information saved in the first part of the this lab.

```
climSA0k <- read.table("Data/Environment/climaSA0k_all.txt", h = T)
class(climSA0k)
```

```
gridded(climSA0k) <- ~x + y
climSA0k <- stack(climSA0k)
climSA0k
```

Inspect the both data.

```
plot(climSA0k$bio.1)
points(hyacinth_coords[, "lon"], hyacinth_coords[, "lat"])
```

Great, it seems that everything is good! However, we know that the species can't survive in the Andes, so, we need to decide the areas that the species can invade (the M component in the BAM diagram). This step is not trivial and maybe it can be the most difficult, but let's keep simple. We will define the **M** as a buffer of across all occurrence points.

```
hyacinth_spatial <- hyacinth_coords
coordinates(hyacinth_spatial) <- ~lon+lat
proj4string(hyacinth_spatial) <- CRS("+proj=longlat +datum=WGS84")
hyacinth_spatial
```

This part will take some time due the quantity of points...

```
hyacinth_buffer <- gBuffer(hyacinth_spatial, byid = FALSE, id = NULL, width = 3.0,
                           quadsegs = 5, capStyle = "ROUND",
                           joinStyle = "ROUND", mitreLimit = 1.0)
```

How it looks...

```
plot(climSA0k$bio.1)
plot(hyacinth_buffer, add = TRUE)
plot(hyacinth_spatial, add = TRUE)
```

Using the buffer created in the previous step we will crop the bioclimatic data.

```
hyacinth_bios <- crop(climSA0k, hyacinth_buffer)
hyacinth_bios2 <- mask(hyacinth_bios, hyacinth_buffer)
```

How it look...

```
plot(hyacinth_bios2$bio.1)
plot(hyacinth_spatial, add = TRUE)
```

Now, we will extract the environmental information using the hyacinth occurrences.

```
hyacinth_var <- extract(hyacinth_bios2, hyacinth_coords, cellnumbers = T)
head(hyacinth_var)
```

Merge the environmental data with the species occurrences.

```
hyacinth_var <- cbind(hyacinth_coords, hyacinth_var)
hyacinth_var[1:5, ]
```

Remove of rows with NA values.

```
hyacinth_var <- na.omit(hyacinth_var) # remove NA's
dim(hyacinth_var)
```

Now, lets remove all duplicate values based on the cells position.

```
duplicated(hyacinth_var[, "cells"])
```

Wow, there are a lot of duplicate cells.

```
a <- which(duplicated(hyacinth_var[, "cells"]) == T)
a
```

Using the object a, we can delete the duplicate cells.

```
hyacinth_var_noDUPS <- hyacinth_var[-a, ]
dim(hyacinth_var_noDUPS)
```

Interestingly only 130 occurrences do not have duplicates. We can save this information in our OCC subfolder.

```
write.table(hyacinth_var_noDUPS, row.names = FALSE, "Data/OCC/hyacinth_var.txt", sep = "\t")
```

Inspect how the 130 looks in the map.

```
plot(hyacinth_bios2$bio.1)
points(hyacinth_var_noDUPS[, "lon"], hyacinth_var_noDUPS[, "lat"])
```

The last step to build the Hyacinth macaw environmental niche is to establish a background (pseudoabsences). To do that we will use the environmental data.

```
clima0k <- read.table("Data/Environment/climaSA0k_all.txt", h = T)
dim(clima0k)
```

Set the background.

```
ext <- extent(hyacinth_bios2)
id.back <- data.frame(randomPoints(hyacinth_bios2, n = 2*nrow(hyacinth_var_noDUPS),
                                ext = ext))
```

```
# This is other way to stablish random points
```

```
#id.back <- sample(1:nrow(hyacinth_bios2), nrow(hyacinth_var_noDUPS), replace = TRUE) # nrow(hyacinth_v  
#length(id.back)
```

How it looks. These random points represent **pseudoabsences** for our species.

```
plot(hyacinth_bios2$bio.1)
points(id.back[, "x"], id.back[, "y"])
```

Now we can extract the bioclimatic data using the random points.

```
background_var <- extract(hyacinth_bios2, id.back, cellnumbers = T)
dim(background_var)
names(background_var)

background_var <- cbind(id.back, background_var)
background_var[1:5, ]

write.table(background_var[, c(1:2, 5:23)], "Data/Environment/background.txt", row.names = F, sep = "\t")
```

Well, we are almost there.

Variable selection

Now, we need to decide which variables are needed to model the niche (Grinnellian niche or Fundamental niche) for the Hyacinth macaw. Here we can rely the selection to the specialist or use statistical tools or both.

Let's use statistical tools.

First we can create a temporal object that will help us in this process.

```
hyacinth_var_tmp <- data.frame(hyacinth_var_noDUPS)
names(hyacinth_var_tmp) <- c("x", "y", "cells", "id", names(hyacinth_var_tmp)[5:23])
hyacinth_var_tmp$id <- rep(1, nrow(hyacinth_var_noDUPS))

background_tmp <- data.frame(background_var)
background_tmp$id <- rep(0, , nrow(background_var))

var_tmp <- rbind(hyacinth_var_tmp, background_tmp) # Here 1 means presence and 0 absence
head(var_tmp)
```

One way to select variables is to explore the correlation among the predictors. We can use a threshold of $|r| < 0.7$ which means that correlated variables below this threshold can be considered as unproblematic (Dormann et al. 2013), however you can use a more conservative threshold, such as < 0.5 . Let's see.

```
# We first estimate a correlation matrix from the predictors. We use Spearman rank correlation coefficients
cor_mat <- cor(na.omit(var_tmp[, c(5:23)]), method = 'spearman')

# We can visualise this correlation matrix. For better visibility, we plot the correlation coefficients
corrplot.mixed(cor_mat, tl.pos = "lt", tl.cex = 0.5, number.cex = 0.5,
               addCoefasPercent = TRUE)
```

Wow, looks like a lot of the predictors are highly correlated and inspecting predictor by predictor can take for ever. However we can also ask to the computer make the work for us, in this case we will build univariate models and select the best fitting model. To do that, we will use the function **selectBelow07** which is an adaptation from a function written by Damariz Zurell.

```
source("R-Functions/selectVariables.R")

head(var_tmp)
tail(var_tmp)
bios <- names(var_tmp)[5:23]
var_sel <- selectBelow07(pred_names = bios, response_name = "id",
                        data = na.omit(var_tmp), cor_mat = cor_mat, threshold = 0.7)
```

The function return a list of three objects: 1. AIC values for each model 2. The correlation matrix 3. A vector of the selected variables.

```
var_sel$AIC
var_sel$cor_mat
var_sel$pred_sel
```

According to the **selectBelow07** seven predictors best fit the data and also have low correlation. Assuming this is correct, we will use these variables for further analyses.

```
pred_sel <- var_sel$pred_sel
pred_sel
```

Building ecological niche models

```
head(hyacinth_var_noDUPS)
hyacinth_var_sel <- hyacinth_var_noDUPS[, c("lon", "lat", "bio.1", "bio.6", "bio.8",
                                             "bio.17", "bio.18", "bio.3", "bio.16")]
head(hyacinth_var_sel)

background_sel <- background_var[, c("x", "y", "bio.1", "bio.6", "bio.8",
                                       "bio.17", "bio.18", "bio.3", "bio.16")]
head(background_sel)
```

The first step is to set data for training and test. To do that, we will use the TRUE occurrences and the background (pseudoabsences).

```
id.ocur <- sample(1:nrow(hyacinth_var_noDUPS), round(0.75*nrow(hyacinth_var_sel)))
length(id.ocur)
```

```
id.back <- sample(1:nrow(background), round(0.75*nrow(background_sel)))
length(id.back)
```

```
hyacinth_bios_sel <- stack(hyacinth_bios2$bio.1, hyacinth_bios2$bio.6,
                           hyacinth_bios2$bio.8, hyacinth_bios2$bio.17,
                           hyacinth_bios2$bio.18, hyacinth_bios2$bio.3,
                           hyacinth_bios2$bio.16)
```

```
plot(hyacinth_bios_sel)
```

Set data for training.

```
training <- prepareData(x = hyacinth_bios_sel,
                        p = hyacinth_var_sel[id.ocur, 1:2],
                        b = background_sel[id.back, 1:2], xy = T)
head(training)
```

Set data for testing.

```
test <- prepareData(x = hyacinth_bios_sel,
                    p = hyacinth_var_sel[-id.ocur, 1:2],
                    b = background_sel[-id.back, 1:2], xy = T)
head(test)
```

Now, we will fit some models using different algorithms. Bioclim, DomianGower, SVM, GLMz.

Fit models

Bioclim model

The BIOCLIM algorithm computes the similarity of a location by comparing the values of environmental variables at any location to a percentile distribution of the values at known locations of occurrence ('training sites').

```
Bioclim.model <- bioclim(x = training[training[, "pb"] == 1, -c(1:3)])
```

```
Bioclim.model
```

```
plot(Bioclim.model)
```

Lets plot the response variables, or single variable response curves for a model.

We can use this model to explore how the species responds to these particular variables (based on the particular algorithm that we used and the relationship between the geographic and environmental space), using a response function that creates response plots for each variable, with the other variables at their median value.

```
response(Bioclim.model)
```

Gower model

The Domain algorithm computes the Gower distance between environmental variables at any location and those at any of the known locations of occurrence ('training sites'). For each variable the minimum distance between a site and any of the training points is taken.

```
Gower.model <- domain(x = training[training[, "pb"] == 1, -c(1:3)])
```

```
Gower.model
```

We can also explore the response for each variable.

```
response(Gower.model)
```

Support Vector Machines or SVM model

Support Vector Machines are an excellent tool for classification, novelty detection, and regression.

Support Vector Machines (SVMs; Vapnik, 1998) apply a simple linear method to the data but in a high-dimensional feature space non-linearly related to the input space, but in practice, it does not involve any computations in that high-dimensional space. This simplicity combined with state of the art performance on many learning problems (classification, regression, and novelty detection) has contributed to the popularity of the SVM (Karatzoglou et al., 2006).

```
head(training)
```

```
svm.model <- ksvm(pb ~ bio.1 + bio.6 + bio.8 + bio.17 + bio.18 + bio.3 + bio.16,  
                 data = training)
```

Generalized Linear Models

A generalized linear model (GLM) is a generalization of ordinary least squares regression. Models are fit using maximum likelihood and by allowing the linear model to be related to the response variable via a link function and by allowing the magnitude of the variance of each measurement to be a function of its predicted

value. Depending on how a GLM is specified it can be equivalent to (multiple) linear regression, logistic regression or Poisson regression.

```
glm.model <- glm(pb ~ bio.1 + bio.6 + bio.8 + bio.17 + bio.18 + bio.3 + bio.16,  
                data = training, family = binomial(link = "logit"))
```

Until now, we just fitted different models for the Hyacinth macaw using four different algorithms. The next step is use these fitted models to make spatial predictions.

Predictions

Notice that in here we will use the fitted models to predict the fundamental niche (Grinnellian) for our species in South America.

Bioclim

```
Bioclim0k <- predict(object = Bioclim.model, x = hyacinth_bios_sel)
```

Lets see how the bioclim model looks like! and after that repeat for the other three algorithms.

```
plot(Bioclim0k)
```

What about if we plot the prediction and the TRUE occurrences used for fit the model.

```
plot(Bioclim0k)  
points(training[training[, "pb"] == 1, "x"], training[training[, "pb"] == 1, "y"])
```

Gower

```
Gower0k <- predict(hyacinth_bios_sel, Gower.model)
```

Inspect the Gower model.

```
plot(Gower0k)
```

SVM

```
svm0k <- predict(hyacinth_bios_sel, svm.model)
```

```
plot(svm0k)
```

GLMz

```
GLM0k <- predict(hyacinth_bios_sel, glm.model)
```

```
plot(GLM0k)
```

Let's standardize the GLM model as we want values between 0 and 1

```
source("https://raw.githubusercontent.com/jesusNPL/DetectionTT/master/normalize.R")  
GLM0k <- normalize(ras = GLM0k)
```

And plot it again.

```
plot(GLM0k)
```

Now lets plot all four models # Ver los resultados

```
par(mfrow = c(2, 2))
plot(Bioclim0k, main = "bioclim")
plot(Gower0k, main = "Gower")
plot(svm0k, main = "svm")
plot(GLM0k, main = "glm")
```

Model evaluation

Now we will evaluate the performance of each model. To evaluate the performance of our models, we will use the testing data. Repeat for all four models.

Cross-validation of models with presence/absence data. Given a vector of presence and a vector of absence values (or a model and presence and absence points and predictors), confusion matrices are computed (for varying thresholds), and model evaluation statistics are computed for each confusion matrix / threshold.

Bioclim

```
Bioclim.eval <- evaluate(p = test[test[, "pb"] == 1, 1:2],
                        a = test[test[, "pb"] == 0, 1:2],
                        model = Bioclim.model,
                        x = hyacinth_bios_sel)
```

```
Bioclim.eval
str(Bioclim.eval)
```

Gower

```
Gower.eval <- evaluate(p = test[test[, "pb"] == 1, 1:2],
                      a = test[test[, "pb"] == 0, 1:2],
                      model = Gower.model,
                      x = hyacinth_bios_sel)
```

SVM

```
svm.eval <- evaluate(p = test[test[, "pb"] == 1, 1:2],
                   a = test[test[, "pb"] == 0, 1:2],
                   model = svm.model,
                   x = hyacinth_bios_sel)
```

GLMz

```
glm.eval <- evaluate(p = test[test[, "pb"] == 1, 1:2],
                   a = test[test[, "pb"] == 0, 1:2],
                   model = glm.model,
                   x = hyacinth_bios_sel)
```

Now, lets plot the performance of the models. Looking to the AUC values, the **simplest** model (bioclim) have a better performance.

```
# validate results
par(mfrow = c(2, 2))
plot(Bioclim.eval, "ROC")
plot(Gower.eval, "ROC")
plot(svm.eval, "ROC")
plot(glm.eval, "ROC")
```

Finally, lets produce “shadows” from “ghosts”.

Species geographical distributions

First we need to set threshold for each model, the threshold (cut-off) is used to transform model predictions (probabilities, distances, or similar values) to a binary score (presence or absence).

Bioclim

```
Bioclim.thr <- threshold(Bioclim.eval)
```

Inspect the “thr” object and then repeat for the other three models.

```
Bioclim.thr
```

There are different columns in the “thr” object and each one is a measure of threshold, in this example we will use the “spec_sens” column or species sensibility. In the spec_sens, the threshold at which the sum of the sensitivity (true positive rate) and specificity (true negative rate) is highest.

```
bio <- Bioclim.thr$spec_sens
```

Gower

```
Gower.thr <- threshold(Gower.eval)
```

```
gow <- Gower.thr$spec_sens
```

SVM

```
svm.thr <- threshold(svm.eval)
```

```
s <- svm.thr$spec_sens
```

GLMz

```
glm.thr <- threshold(glm.eval)
```

```
g <- glm.thr$spec_sens
```

Lets see what happen...

```

par(mfrow = c(2, 2))
plot(BioclimOk > bio, main = "Bioclim")
plot(GLM0k > g, main = "GLM")
plot(GowerOk > gow, main = "Gower")
plot(svm0k > s, main = "SVM")

```

Looks like the different algorithms produce different shadows... let's see what happens if we combine all four models into a single one "ensemble framework".

```

# Combine all thresholds
thrs <- (bio + gow + s + g)

```

Now combine all predictions into a single prediction.

```

tmp <- stack(BioclimOk, GowerOk, GLM0k, svm0k)
names(tmp) <- c("Bioclim", "Gower", "GLM", "SVM")

```

First, let's sum all predictions and see what happens...

```

map.sum <- sum(tmp) # sum

```

Plot the result.

```

par(mfrow = c(1, 2))
plot(map.sum)
plot(map.sum > 2)

```

Now, what about the mean...

```

map.mean <- mean(tmp) # mean

```

```

map.sd <- calc(tmp, sd) # sd

```

Plot the results...

```

plot(map.mean)

```

```

par(mfrow = c(1, 2))
plot(map.mean)
plot(map.sd)

```

However, this is a problematic approach as the values predicted by the models are not all on the same (between 0 and 1) scale; so you may want to fix that first. Another concern could be weighting. Let's combine the four models weighted by their AUC scores. Here, to create the weights, we subtract 0.5 (the random expectation) and square the result to give further weight to higher AUC values.

```

auc <- sapply(list(Bioclim.eval, Gower.eval, svm.eval, glm.eval), function(x) x@auc)

```

```

w <- (auc-0.5)^2

```

```

map.mean.weight <- weighted.mean(tmp, w)

```

```

plot(map.mean.weight)

```

Finally, let's plot the four predictions

```

par(mfrow = c(2, 2))
plot(map.sum, main = "Sum of all models")
plot(map.mean, main = "Mean of all models")
plot(map.mean.weight, main = "Weighted mean of all models")
plot(map.sd, main = "Standard deviation of all models")

```

That's it, we modeled environmental niches (ghosts) to produce geographical distributions (shadows)...

Exercises

Please respond each question based on the practice.

1. Summarize the data: how many records are there, how many have coordinates, how many records without coordinates have a textual georeference (locality description)?
2. Do you think the observations are a reasonable representation of the distribution (and ecological niche) of the species?
3. There is a best model? Explain your answer.
4. Can we use ENM/SDM to aid the conservation of the Hyacinth macaw? Explain your answer.
5. Based on the lecture and the practice, what can we conclude?

References

- Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., ... Lautenbach, S. (2012). Collinearity: a review of methods to deal with it and a simulation study evaluating their performance. *Ecography*, 36(1), 27–46. doi:10.1111/j.1600-0587.2012.07348.x
- Karatzoglou, A., D. Meyer & K. Hornik, 2006. Support Vector Machines in R. *Journal of statistical software* 15(9).
- Lima-Ribeiro, M.S., S. Varela, J. González-Hernández, G. Oliveira, J.A.F. Diniz-Filho & L.C. Terribile. 2015. ecoClimate: a database of climate data from multiple models for past, present, and the future for macroecologists and biogeographers. *Biodiversity Informatics*, 10: 1-21.
- Vapnik, V., 1998. *Statistical Learning Theory*. Wiley, New York.