SEEC Stats Toolbox

Species distribution modelling in R







SEEC - Statistics in Ecology, Environment and Conservation

Our imperfect knowledge of the natural world







Our STILL imperfect knowledge of the natural world



Rahlao et al. (2010). Weed Research 50: 537





What we actually want



Geerts et al. (2016) Biological Invasions. DOI 10.1007/s10530-016-1226-y





What is a species distribution model (SDM)?



Franklin (2009). Mapping species distributions. Cambridge University Press





Uses of SDMs

- Conservation prioritization
- Climate change predictions
- Rare species detection
- Invasive species screening
- Generating hypotheses correlates of distribution
- Habitat suitability





Eight steps to your own SDM

- 1. Occurrence data
- 2. Environmental data
- 3. Background samples
- 4. Study extent
- 5. Data cleaning
- 6. Modelling
- 7. Checking your model
 8. Projecting your model



Southern Bald Ibis Geronticus calvus (Boddaert, 1783)



www.hbw.com











Sources:

- Your own data
- gbif.org/species Clobal Bindh
- newposa.sanbi.org



Useful to keep ID, locality, date... for data cleaning purposes





🗷 R!	Studio Source Editor	—		\times
State	ts_Toolbox_SDMs.R* ×			
	🐑 💼 🔄 Source on Save – 🔍 Ž 🗸 📳	📑 Run 🔄	Source	• =
13				~
14	#			
15	# det presence and background data:			
17	#We are going to get our data using the package robif. This allows us to download directly from GBTE.			
18	library(rgbif)			
19	#Search for the GBIF key for the species in question. This is a unique identifier that tells GBIF what dat	ta we want:		
20	<pre>(key = name_suggest(q='Geronticus calvus', rank='species'))</pre>			
21	key = key\$key			
22	#See how many records in GBIF are available. We only want records that have coordinates so we say "georef"	erenced=T":		
23	occ_count(taxonkey = key, georeterenced=T)			
24	#Download data from GBF, keep only records with coordinates:			
26	presenter - occ_search(taxonkey-key, filet=10000, flascoordinate=1)			
27	#Keep only dataframe part of GBIF data:			
28	pres = presGBIF\$data			
29	names(pres)			
30	#Keep only some data fields:			
31	<pre>pres = pres[,c('key','scientificName','decimalLongitude','decimalLatitude','basisOfRecord','year','country</pre>	ycode','loca	ality')]	
32	#To save time, save these data to a CSV file so that you can use this instead if you run this script at a	later stage	2:	
33	write.csv(pres, pres.csv, row.names=F)			
2/1				

















- Choosing an extent:
 - Size of the study area
 - Area from which
 background samples
 are selected









See Webber et al. (2011). Diversity and Distributions 17: 978





• Smaller extent:

- Produces a smaller predicted area and perhaps emphasis on incorrect predictor variables.

• Larger extent:

- Produces a larger predicted area, which may be unrealistic, and climate variables often dominate.



Barbet-Massin et al. 2010. Ecography 33: 878





- Depends on your objective: occupied or suitable habitats (Merow et al., 2013)
 - Occupied habitat accessible area.
 - Identify where a species is currently found.
 - Identify environmentally limiting factors.
 - Suitable habitat area that you wish to contrast with against the presences.
 - Climate change.
 - Invasive species.





🗷 RSt	udio Source Editor	_		\times
C Stats_	Toolbox_SDMs.R* *			
00	🗊 🔒 🖸 Source on Save 🔍 Ž 🗸 🗐	•	Bource	• =
185	#we are going to select only pseudo-absences within 2 degrees (~240 km) of presences. You could also select only			~
186	#records within Koppen-Geiger climate zones that presences occur in, or something similar.			
18/	Inprary(rgeos)			
188	pressp = prescient #create new dataframe from presences			
109	condition (d) = condition (d			
191	press(f) = gruffer(press), 2002 #Give the spartal points object a projection (in this case har tebeestideks+)			
192	#Plot the buffer area on a map:			
193	plot(sa)			
194	plot(presBuff, add=T, border='red')			
195	points(pressp) #Add presence records			
196	#This next part will create a spatial points object from our pseudo-absences, overlay these on our buffer object			
197	#and then we select only pseudo-absences falling within the buffer object:			
198	backsp = backClean #Create new dataframe from pseudo-absences			
199	coordinates(backSp) =~ x+y #Transform dataframe into a spatial points object			
200	crs(backSp) = '+init=epsg:2052' #Give the spatial points object a projection (in this case Hartebeesthoek94)			
201	overButterback = over(backSp, presButt)			
202	backClean = backClean[:1s.na(overButterback),]			
203	#Add the pseudo-absences to our map:			
204	points(backs), core blue) #Add pseudo-absences, including close beyond burle zone			
205	points (press) #Add presence records			
207	France (France) and by an			







- Useful references:
 - Anderson & Raza (2010). Journal of Biogeography 37: 1378
 - Barve et al. (2011) Ecological Modelling 222: 1810
 - Merow et al. (2013). Ecography 36: 1058











Sources:

- Bioclim www.worldclim.org/bioclim
 - Can download directly in R using the "getData" function in the raster package
- BioOracle www.oracle.ugent.be
- Google is your friend





Considerations:

- Spatial resolution



Pearson & Dawson (2003). Global Ecology & Biogeography 12: 361





- Considerations:
 - Temporal resolution







Considerations:

- Selecting variables
 - Include variables that directly limit a species (e.g. min. temperature)
 - Include variables that are resources (e.g. nutrients)
 - Include variables that link to physiology (e.g. water availability)





📧 RSt	udio Source Editor		_		×
Stats	Toolbox_SDMs.R* ×				
	🗊 🔚 🗌 Source on Save 🛛 💁 🎽 🚛	📑 Run	••	🕞 Source	
65 66 67	# #Get environmental layers: #	_			^
68 69	library(raster)				
70 71	#One way to get WORLDCLIM data is to use the getData function in raster. We want the Bioclim variables a #resolution of 10 minutes	at a			
72 73	bio = getData(name='worldclim', var='bio', res=10, download=T) #This creates a raster stack with all the	e Biocli	m var	iables.	
74	#The alternative way of reading in rasters is shown below. You should also use the approach below once	you have			
75 76	#downloaded the Bioclim data using the getData function above. In other words, you can comment out the #above with the getData function for future use of this script.	line			
77	(bioFiles = list.files('wc10/')) #List the files in the wc10 folder that was created using the getData	function	1		
78	(bioFiles = bioFiles[grep('.bil', bioFiles, fixed=T)]) #Select only the .bil files (these are the actual (bioFiles = bioFiles[s(1,12,10,2,11)]) #Beender file pares to get into pymerical order (BIO1, BIO2,))	l raster	s)		
79 80	(DIOFTIES = DIOFTIES[C(1,12;19,2;11)]) #Reorder The names to get into numerical order (BIO1, BIO2)				
81	<pre>bio = stack(paste0(getwd(), '/wc10/', bioFiles), RAT=F)</pre>				
82	<pre>#Rescale Bioclim variables (BIOS 1 to 11) that have been multiplied by 10 (or 1000 for BIO4) (For examp) """"""""""""""""""""""""""""""""""""</pre>	le, BIO1	, whi	ich is	
83	#mean annual temperature, has values from -209 to 314, whereas these should be -20.9 to 31.4 degrees C. #The *10 format is the default format for Bioclim):				
85 -	for $(r i c(1:3,5:1))$				
86	bio[[r]] <- bio[[r]]/10				
87	}				
88	blo[[4]] <- blo[[4]]/1000				
90	#	-			











- When true absences unavailable
- A priori considered equally likely to contain individuals of a species

Merow et al. (2013). Ecography 36: 1058





What background samples are used for



Merow et al. (2013). Ecography 36: 1058





Sampling bias

Target Group Sampling

- Use coordinates of related species or in same functional group to select background
- Accounts for sampling bias
- Create a bias grid for Maxent

• Useful refs:

Hijmans et al. (2000) Conservation Biology 14: 1755

- Elith et al. (2010) Methods in Ecology and Evolution 1: 330
- Merow et al. (2013)
- Phillips et al. (2009) Ecological Applications 19: 181











🔞 RSt	udio Source Editor	_	· 🗆	×
C Stats	Toolbox_SDMs.R* *			
$\langle \phi \phi \rangle$	🗊 🔒 🖸 Source on Save 🛛 Save 🗐 🤮 🖉 📲	un 📘 🖻	🔶 📑 Source	• Ē
208	#Create a sampling bias grid for Maxent			~
209	#First we need to create a dataframe of all our observations (presences and pseudo-absences). This will give	an ind	dication	
210	#of sampling effort:			
211	allDat = rbind(presUncleaned, backUncleaned)			
212	#Next we rasterize this dataframe:			
213	allRas = rasterize(cbind(allDat\$x, allDat\$y), y=bio[[1]], field=1)			
214	plot(allRas, xlim=c(15,35), ylim=c(-35,-15)) #Plot			
215	#Get ID (cell number) of cells in which there are presences (of all species):			
216	<pre>presRasID = which(values(allRas)==1)</pre>			
217	#Get the coordinates of the above cells:			
218	<pre>presRasCoords = coordinates(allRas)[presRasID,]</pre>			
219	#Use a Gaussian kernel to smooth the number of records. Provides a density estimate that will be our bias gri	d :		
220	library(MASS)			
221	<pre>dens = kde2d(x=presRasCoords[,1], y=presRasCoords[,2], h=c(2,2), n=c(nrow(allRas),ncol(allRas)))</pre>			
222	#Rasterize the density estimate:			
223	biasGrid = raster(dens)			
224	plot(biasGrid, main="Bias grid") #Plot			
225	#write this to a raster file:			
226	writeRaster(biasGrid, 'biasGrid.tif', overwrite=T)			
227				
220	и.	_		











- Worthy of a lecture on its own
- Species names
 - Synonyms
 - Misapplication
 - Useful resources:
 - theplantlist.org
 - www.eol.org
 - www.ncbi.nlm.nih.gov
 - www.itis.gov
 - https://ropensci.org/tutorials/taxize_tutorial.html





- Coordinate errors
 - Spatial resolution
 - Zero lat/lon
 - Swapped lat and lon
 - How to detect:
 - Points in the sea (terrestrial organisms) or on land (marine organisms)
 - Country name mismatch
 - Elevational mismatch
 - Outliers with respect to environmental data





Pseudoreplication







Collinearity

- Predictors that are highly correlated with one another (r > 0.8)
- Can be a problem if one wants to understand environmental factors that limit species' distributions (Merow et al., 2013)



 Not a problem if predicting distribution is the sole aim (Elith et al., 2011. Diversity and Distributions 17: 43)





R package biogeo

📵 RS	tudio Source Editor	-		×
Stats	_Toolbox_SDMs.R* ×			
$\langle \phi \phi \rangle$	🐒 🔒 🖸 Source on Save 🛛 🔍 🖉 📲	1 🕪	🕞 Source	• =
109	#Check for fossil records (if using GBIF data):			~
110	levels(as.factor(presCleanSbasisOfRecord))			
111	#To remove fossil records:			
112	prescleanSExclude[prescleanSpasisOTRecord== FOSSIL_SPECIMEN] = 1			
113	prescreansReason[prescreansbasisonRecord== FOSSIL_SPECIMEN] = FOSSIL #Give a reason for exclusion			
114	#check species names.			
116	levels (presclean Species) #Everything is fine All records have the same and correct species name			
117	(certicity) series and correct species name			
118	#Check temporal resolution of presence records:			
119	hist(presClean\$year, main='Histogram of year in which records were made', xlab='Year')			
120				
121	#Identify duplicates in grid cells (using the resolution of the environmental rasters you use. In this	case	= 10):	
122	presClean = duplicatesexclude(presClean, res=10)			
123	summary(as.factor(presClean\$Exclude)) #0 = records that we will keep, 1 = records that will be remove			
124	backClean = duplicatesexclude(backClean, res=10)			
125	summary(as.Tactor(backClean)Exclude)) #0 = records that we will keep, 1 = records that will be remove	1		
120	#Evelude zone lat & land			
127				
120	prescreamsexcrude[prescreamsx:=0 \approx prescreamsy:=0] = 1 prescreamsexcrude[prescreamsx:=0 \approx prescreamsy:=0] = 'zero lat or lop' #cive a reason for evolution			
130	prescreatiskeason[prescreatisk:=0 @ prescreatisy:=0] = 2ero fac or for #dive a reason for exclusion			
131	#Move points in sea to nearest cell on land:			
132	landbat = nearestcell(presclean, bio[[1]))			
133	presClean <- landpat\$dat			
134	landDatB = nearestcell(backClean, bio[[1]])			
135	backClean = landDatB\$dat			
136	#See points that were changed:			
137	corrected = landDat\$dat\$Correction			
138	iCorrected = grep("7", corrected)			
139	landDat\$dat[iCorrected,] #View records that were corrected			
140				
141	#Remove other points in the sea and any records for which we can't get environmental data (from "bio"	our r	aster sta	ack)
142	prescient = missingvalsexclude(Dio[Li]), prescient) cummany(cs_factor(prescient)) #0 = prescient, that we will keep _1 = prescient that will be promove			
143	Summary as ractor (prescream) Exclude() #0 = records that we will keep, I = records that will be remove backcloap	1		
144	Dackerean <- missingvarsexerude(Dio[[1]], Dackerean)			







http://www.earthskysea.org/!ecology/sdmShortCourseKState2012/sdmShortCourse_kState.pdf





- Maxent (R package dismo)
- Check out the vignette help document on SDM from dismo (is in the script)
- Useful refs:
 - Merow et al. (2013) Ecography 36: 1058
 - Yakulic et al. (2012) Methods in Ecology and Evolution 4: 236





• Some basic things to consider:

Maxent under the hood

	constrain	1.5			
	"Feature"	Constraint	Name	min num of presences	
Pr(env presence)	/	mean	linear	>0	Presentation Presentation
	\wedge	variance	quadratic	≥10	 ✓ Linear features ✓ Quadratic features
		proportion above/below threshold	step	≥15	Product features Threshold features
		(as above) & mean	hinge	≥80	✓ Hinge features
	\times	covariance	product (2 variables)	≥80	Run
		proportion in each category	categorical	>0 Use only f	ninge features to avoid overly- its. Elith et al. 2010. The art of
Ada	ENV m B. Smith Missouri B	otanical Garden		modeling Methods	range-shifting species. Ecology & Evol 1:330-342. 53



http://www.earthskysea.org/lecology/sdmShortCourseKState2012/sdmShortCourse_kState.pdf



- Some basic things to consider:
 - Regularization coefficient (β) penalizes overfitting, but it is user-specified.
 - Try a range of β values and evaluate model fit (e.g. using AUC) (Merow et al., 2013)





🔞 RSt	udio Source Editor		_		\times
🕘 Stats	_Toolbox_SDMs.R* ×				
$\langle \phi \phi \rangle$	🖅 🔒 🖸 Source on Save 🛛 🔍 Ž 🗸 📳	📑 Run	50	🕞 Source	• =
Q, ed	Next Prev All Replace Replace All				×
In sele	ction Match case 🗸 Whole word Regex 🖌 Wrap				
307	#Full model with all occurrences:				~
308	betaRCs <- seq(0.05,0.95,0.05)				
309	<pre>bestAUC = read.csv('bestAUC.csv')\$bestAUC</pre>				
310	mxMod = maxent(x = modLayers)				
311	p = p,				
312	$\mathbf{a} = \mathbf{a}$,				
313	<pre>args = c('-p', 'nothreshold', paste0('beta_lqp=', betaRCs[bestAUC]),</pre>				
314	<pre>paste0('beta_hinge=',betaRCs[bestAUC]), 'outputformat=raw'),</pre>				
315	biasfile=biasGrid,				
316	removeDuplicates=T)				
317	•				











Model accuracy

- Area under the receiver operating characteristic curve (AUC)
- Other measures too and suggested that you use some of these (Sensitivity, specificity, Boyce Index...)
- Usually check by bootstrapping:
 - e.g. 100 model runs
 - 70% of data used to build a model (training data)
 - 30% used to evaluate (test data)





🗷 RSt	tudio Source Editor				\times
Stats	_Toolbox_SDMs.R* ×				
	🗊 🔒 🖸 Source on Save 🛛 🔍 🎽 🚛	📑 Run	9	🕀 Source	• =
333 334 335	#- #Model evaluation (AUC) #-				^
335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350	<pre>#</pre>				
351 352 353 354 355	<pre>p = ptrain, a = atrain, args = c('-p','nothreshold',paste0('beta_lqp=',betaRCs[bestAUC]),paste0('beta_hinge=',b biasfile=biasGrid, #Use a bias grid to reduce sampling bias removeDuplicates=T)</pre>	etaRCs[be	stAUC	:])),	
356 357 358 359 360	<pre>#Look at model evaluation (AUC): e = evaluate(mxMod, p=ptest, a=atrain, x=modLayers) #Store AUC values: AUCs[i] = e@auc }</pre>				
362 363 364 365 366	<pre>#Look at model evaluation (AUC): mean(AUCs) #Get mean AUC values 1.96 * sd(AUCs)/sqrt(100) #95% CI #</pre>				



http://www.earthskysea.org/!ecology/sdmShortCourseKState2012/sdmShortCourse_kState.pdf



- Response curves
 - Are they biologically realistic?
 - Is there enough sampling across the full range of the predictor?







Step 8: Projecting your model







Step 8: Projecting your model

- Very easy to make a map, but is it a good map?
- Model extrapolation
 - novel environmental space



Zurell et al., 2012. Diversity and Distributions 18: 628





Step 8: Projecting your model

- Novel environmental space refs:
 - Multivariate environmental similarity surface (MESS) (Elith et al., 2010. Methods in Ecology & Evolution 1: 330)
 - Environmental overlap masks (Zurell et al., 2012. Diversity and Distributions 18: 628)





THE END













Other useful references

- Elith & Leathwick (2009). Annual Rev. Ecol. Evol. Syst. 40: 677
- Elith et al. (2011). Diversity and Distributions 17: 43
- Liu et al. (2013) Journal of Biogeography 40:778
- Renner et al. (2015). Methods in Ecology and Evolution 6: 366
- And many more... (so carry on reading!)





SEEC Stats Toolbox

Want to broaden your stats knowledge? Unsure of what you can do with your data? Still developing your proposal?

Join us for the SEEC Stats Toolbox seminars

Our next seminar:



Topic: Occupancy Models

Who: Prof Res Altwegg When: **Thursday 25 August 2016 (1-2pm)** Where: UCT campus (venue TBC)

More details: www.seec.uct.ca.za





SEEC - Statistics in Ecology, Environment and Conservation

