Biodiversity data analysis workshop - Day 3

Dominic Henry & Lizanne Roxburgh
Conservation Science Unit

dominich@ewt.org.za lizanner@ewt.org.za









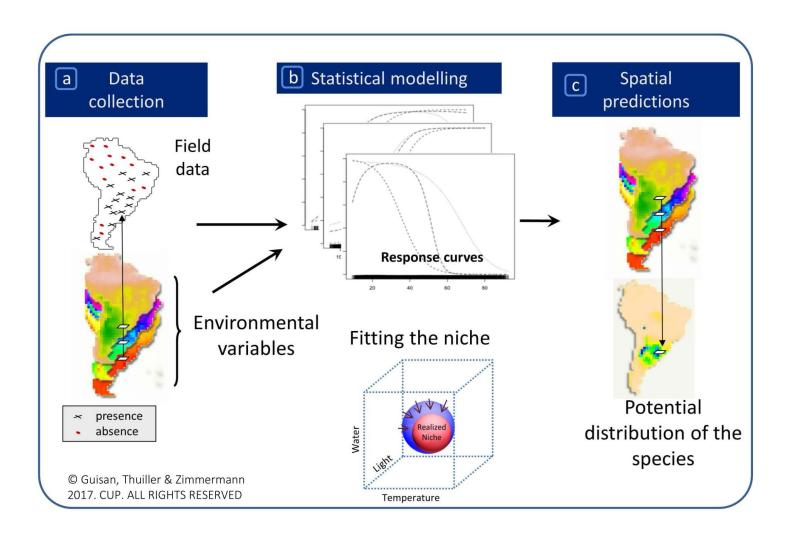
Programme – Day 3



08h30 - 10h00	Session 1 - SDM theory and background
10h00 - 10h30	Tea break
10h30 – 12h15	Session 2 – SDM in R practical demo
12h15 – 13h30	Lunch
13h30 – 15h00	Session 3 – SDM Assignment
15h00 – 15h30	Tea break
15h30 – 17h00	Session 4 – SDM Assignment

Species distribution models (SDMs)





Terminology



- Species distribution models
- Habitat suitability models
- Environmental niche models
- Climate envelope models
- Resource selection functions

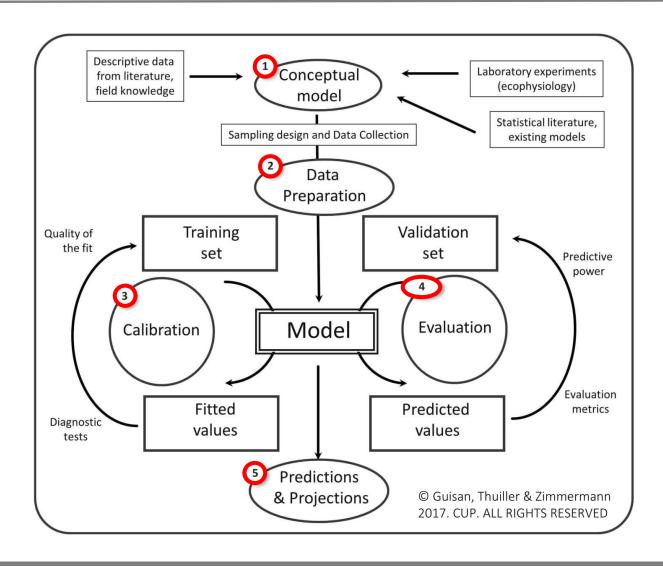
Fields of application



- Guiding field surveys to find new populations
- Predicting invasion of non-native species
- Assessing disease risk
- Supporting conservation prioritisation and reserve selection
- Projecting the impacts of climate change on species' ranges
- Testing evolutionary theory and speciation mechanisms

Methodological steps





What drives species distributions?

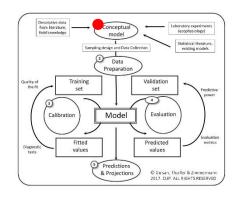


Fundamental biogeographic questions

Where, when and why?

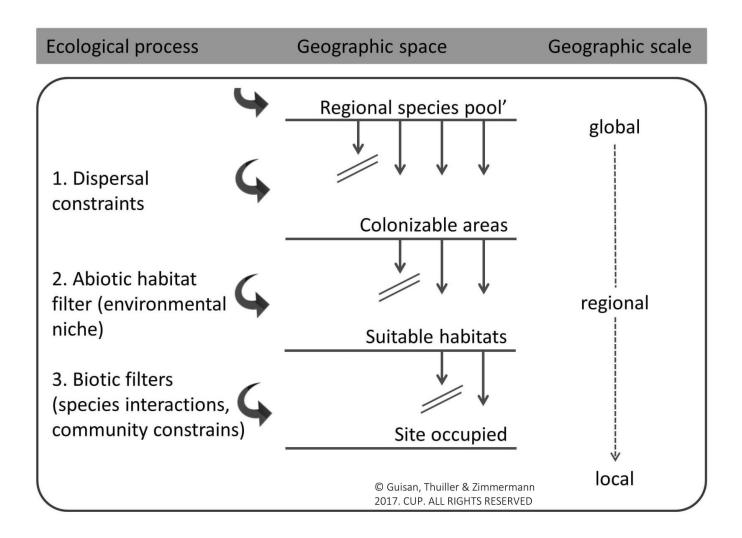
Three primary conditions for site occupancy

- Species must be able to reach the site
- Ecophysiology must match the abiotic conditions
- Biotic environment must be suitable



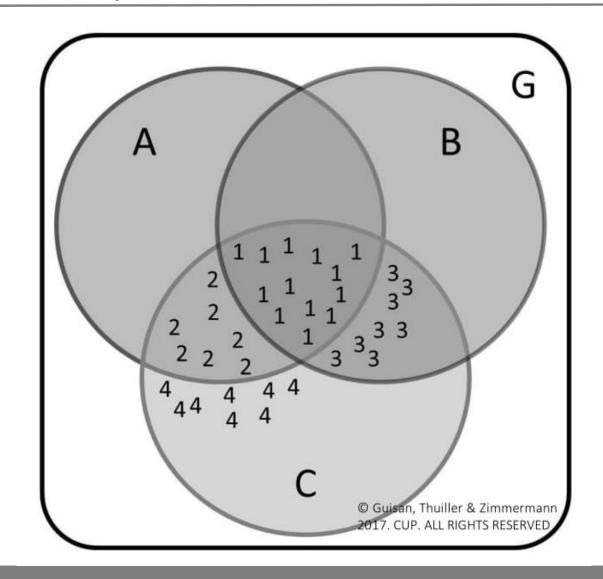
What drives species distributions?





What drives species distributions?





Occurrence outside of realised niche



- Populations might persist in suboptimal habitats
- Source-sink dynamics
- No net population growth but high colonisation pressure
- Constant immigration from neighbouring sites

Other sources of species absence

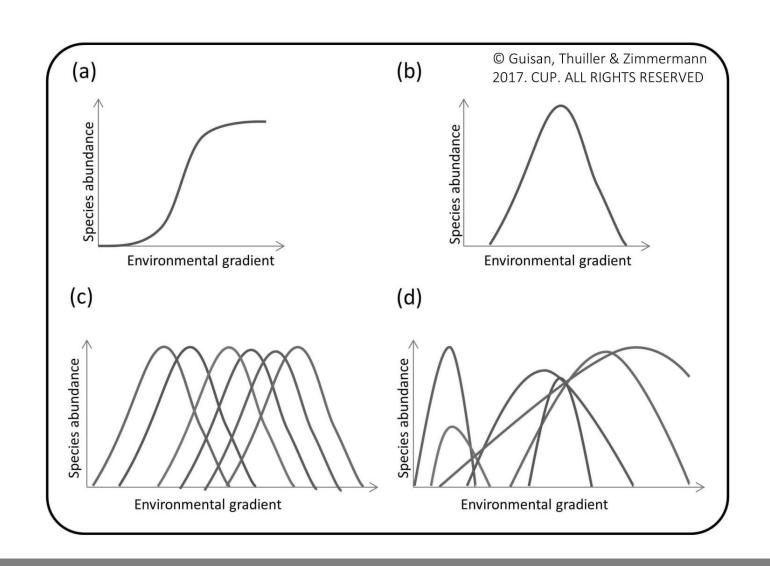


- Natural or human disturbance
- Intrinsic stochasticity
- Temporal extinctions

 Complications therefore arise when suitable habitats are unoccupied and unsuitable habitats are occupied (unexplained variation)

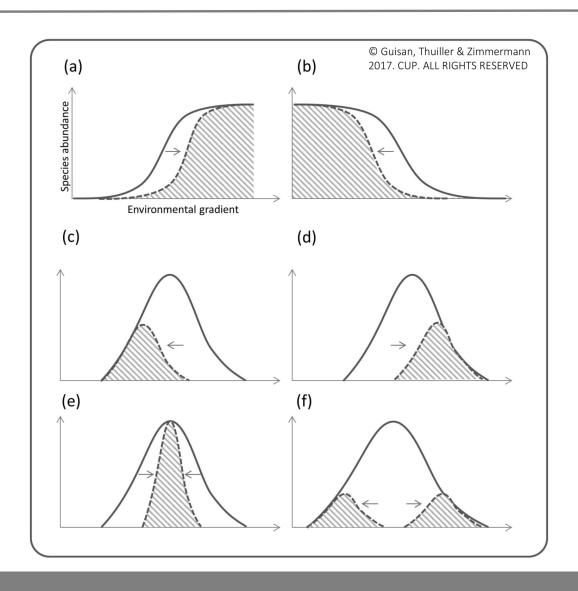
Abiotic environment





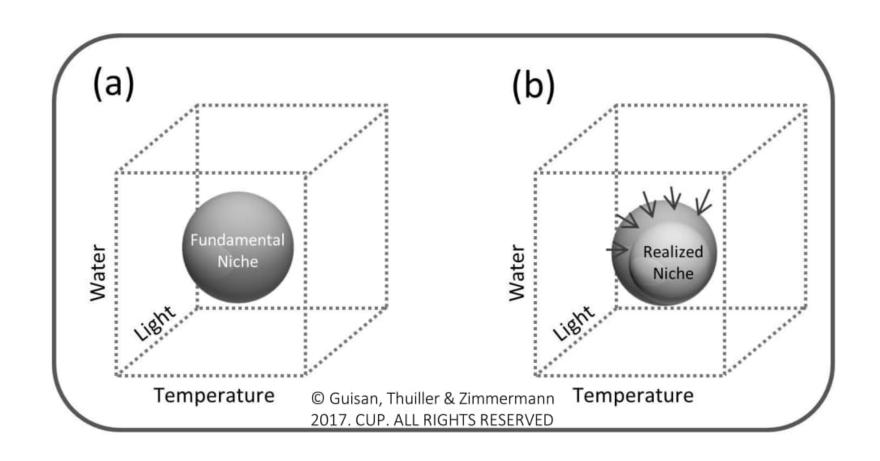
Biotic environment





The realised niche





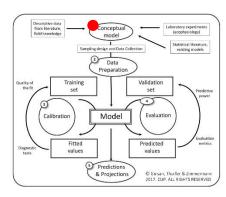
Theoretical considerations



Species-environment equilibrium

Comprehensive set of environmental predictors

Appropriate species observation data



Methodological considerations



Statistical model matches data (categorical, ordinal, continuous)

Measurement error in environmental predictors

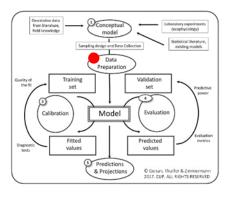
Bias in species data

Independence of species occurrence data points

Data sources - predictors



Category	Example	Source
Climatic	Rainfall, Seasonality	Bioclim - https://www.worldclim.org/bioclim
	Temperature, Aridity	ENVIREM - https://envirem.github.io/
Topographic	Elevation, Aspect, Slope, Terrain Ruggedness	USGS earth explorer/aster global DEM https://asterweb.jpl.nasa.gov/gdem.asp
Land cover	Natural, Agricultural	http://bgis.sanbi.org/landcover/project.asp
Primary productivity	NDVI, NDWI	https://modis.gsfc.nasa.gov/data/dataprod/mod13.php



Bioclim

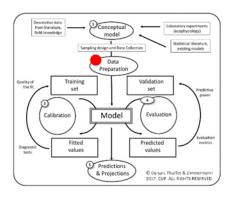


BIO1	Annual Mean Temperature
BIO2	Mean Diurnal Range (Mean of monthly (max temp - min temp))
BIO3	Isothermality (BIO2/BIO7) (* 100)
BIO4	Temperature Seasonality (standard deviation *100)
BI05	Max Temperature of Warmest Month
BI06	Min Temperature of Coldest Month
BI07	Temperature Annual Range (BIO5-BIO6)
BI08	Mean Temperature of Wettest Quarter
BI09	Mean Temperature of Driest Quarter
BIO10	Mean Temperature of Warmest Quarter
BI011	Mean Temperature of Coldest Quarter
BIO12	Annual Precipitation
BIO13	Precipitation of Wettest Month
BIO14	Precipitation of Driest Month
BIO15	Precipitation Seasonality (Coefficient of Variation)
BI016	Precipitation of Wettest Quarter
BIO17	Precipitation of Driest Quarter
BI018	Precipitation of Warmest Quarter
BI019	Precipitation of Coldest Quarter

Data sources - predictors



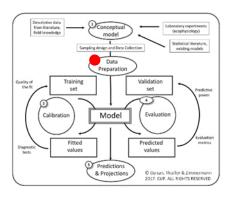
Category	Example	Source
Soils	pH, Structure, Organic content	Prof Mike Cramer (UCT)
Vegetation	Biomes, vegetation classes, Bioregions	http://bgis.sanbi.org/SpatialDataset
Climate change	Various scenarios	https://gisclimatechange.ucar.edu/
Marine	Oceanic properties	NOAA - https://www.ncdc.noaa.gov/data-access/marineocean-data



Data sources - species

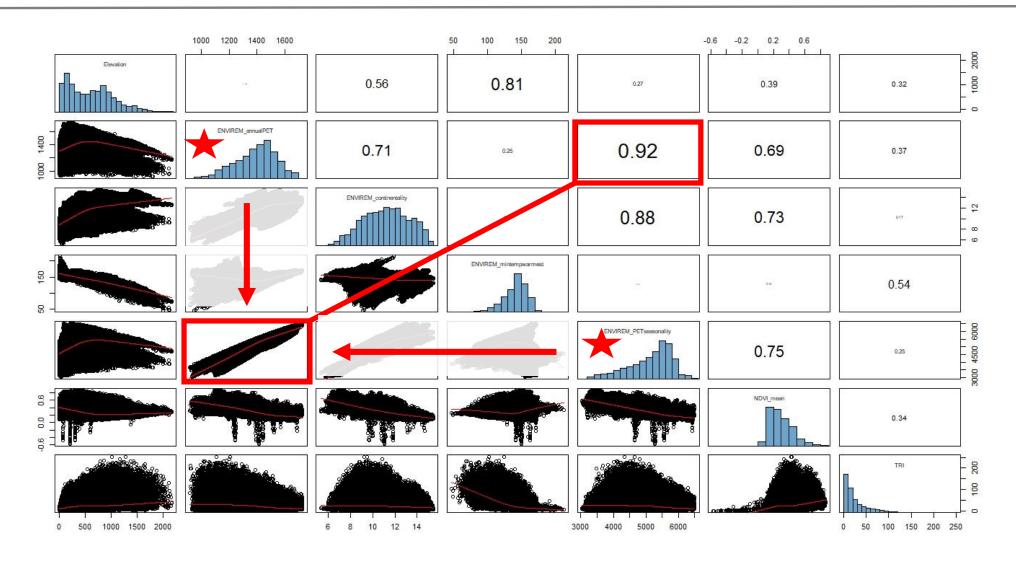


iNaturalist	https://www.inaturalist.org/
GBIF	https://www.gbif.org/
SANBI	http://biodiversityadvisor.sanbi.org/online-biodiversity- data/
Animal Demography Unit	http://www.adu.uct.ac.za/
Map of life	https://mol.org/
LifeMapper	http://lifemapper.org/
AmphibiaWeb	https://amphibiaweb.org/
BirdLife	http://datazone.birdlife.org/home
Plants - BIEN	http://bien.nceas.ucsb.edu/bien/
National collection facilities	
Field surveys	



Predictor correlation





Predictor multi-collinearity



- Multi-collinearity & variance inflation factor (VIF)
- Rule of thumb VIF < 10

Predictor multi-collinearity



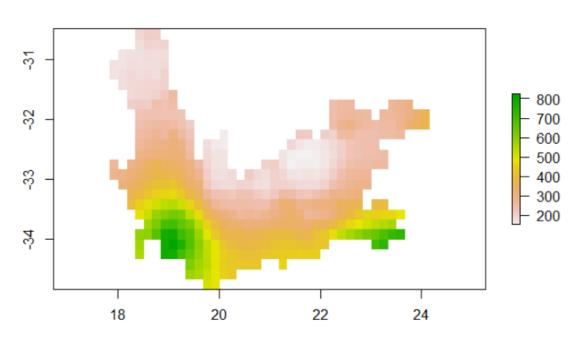


Common spatial resolutions (Bioclim and ENVIREM):

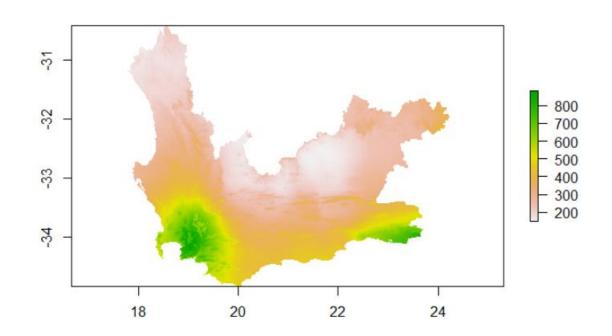
- 30 arc-seconds (~ 1 km)
- 2.5 arc-minutes (~ 5 km)
- 5 arc-minutes (~ 10 km)
- 10 arc-minutes (~ 20 km)







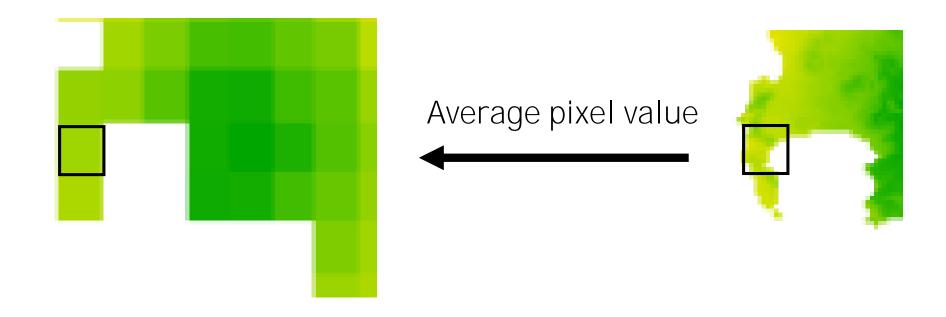
30 arc-seconds



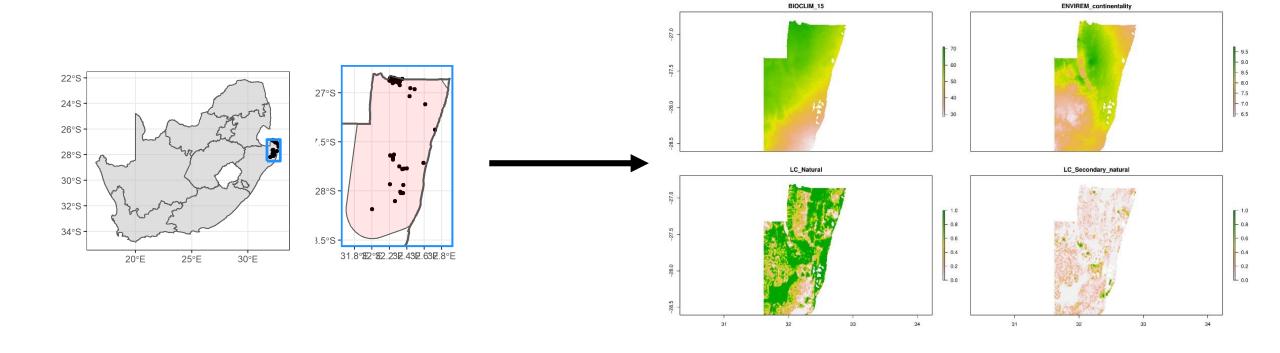


10 arc-minutes

30 arc-seconds





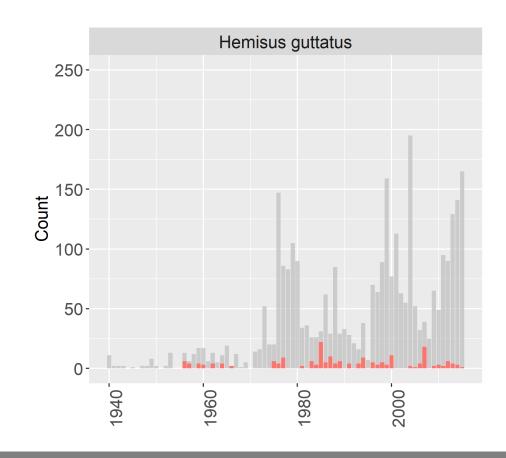


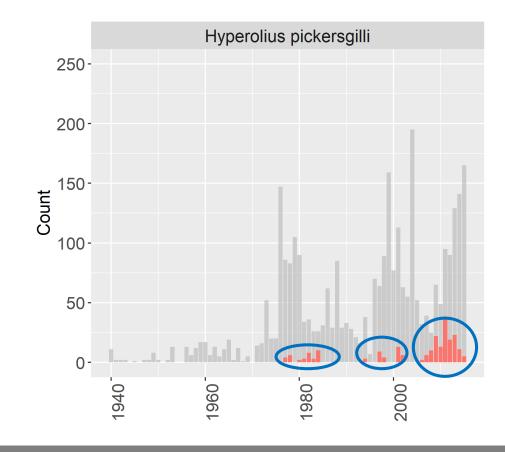


False positives



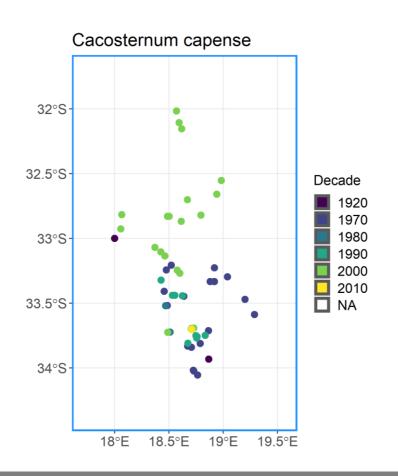
Spatial and/or temporal bias

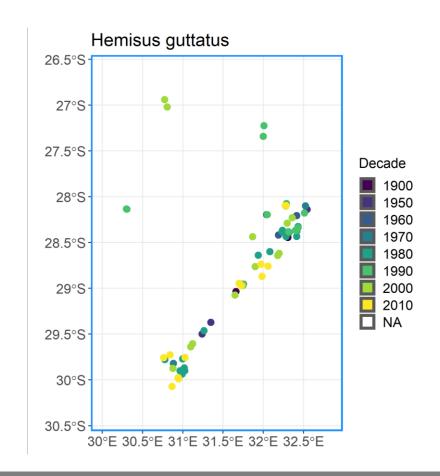






Spatial and/or temporal bias





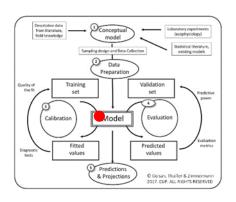


- False positives
- Spatial and/or temporal bias
- Spatial autocorrelation
- Pseudo-replication (duplicate records)
- Positional uncertainty (georeferencing error)

Modelling methods



Method	Software
Maximum Entropy	MAXENT standalone dismo R package
Generalized linear model (GLM)	Base R
Generalized additive model (GAM)	<i>mcgv</i> R package
Boosted regression trees (BRT)	dismo R package
Multivariate adaptive regression splines (MARS)	earth R package
Ensemble Models	BIOMOD standalone biomod2 R package



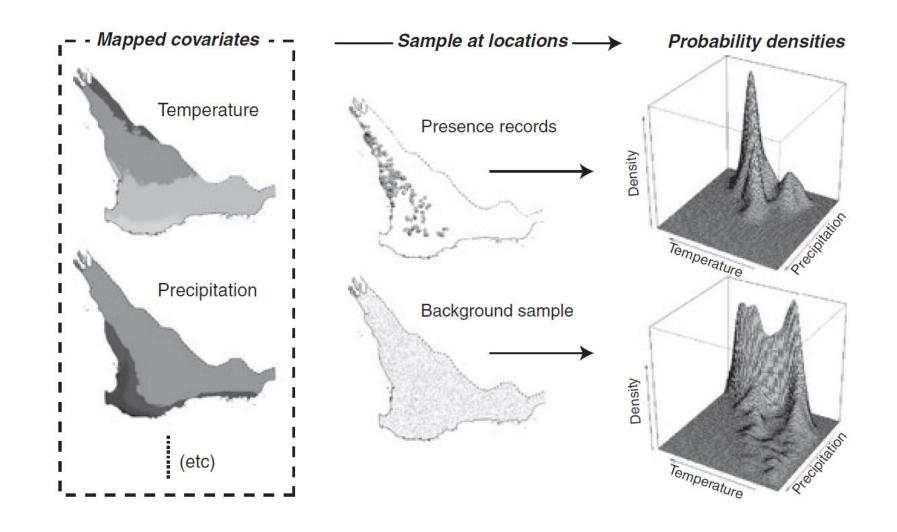
Maxent in R



- MaxEnt takes a list of species presence locations as input (often called presence-only (PO) data)
- Additionally a set of environmental predictors (e.g. precipitation, temperature) across a user-defined landscape that is divided into grid cells.
- From this landscape, MaxEnt extracts a sample of background locations that it contrasts against the presence locations.
- Presence is unknown at background locations.
- Then estimates a relative occurrence rate (or habitat suitability)

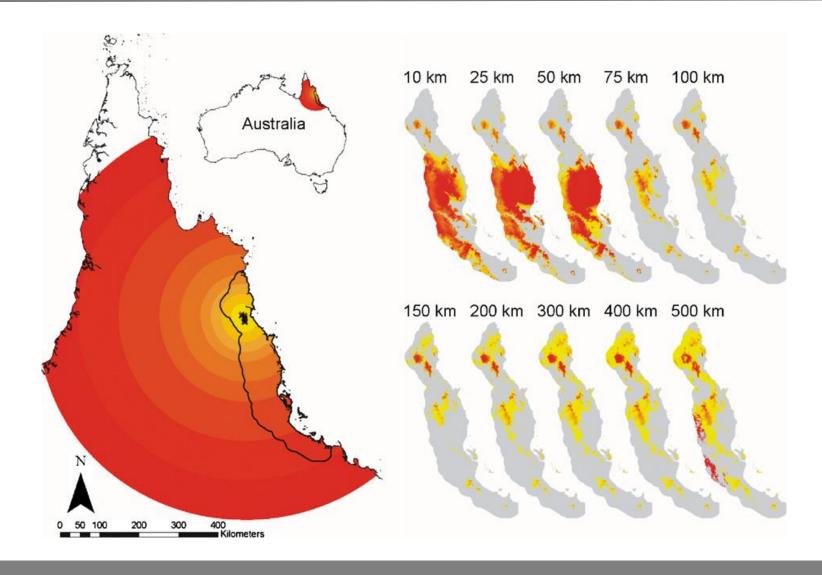
Maxent in R





Background points





Background points



Methods in Ecology and Evolution



Methods in Ecology and Evolution 2012, 3, 327-338

doi: 10.1111/j.2041-210X.2011.00172.x

Selecting pseudo-absences for species distribution models: how, where and how many?

Morgane Barbet-Massin¹*, Frédéric Jiguet¹, Cécile Hélène Albert^{2,3} and Wilfried Thuiller³

¹Muséum National d'Histoire Naturelle, UMR 7204 MNHN-CNRS-UPMC, Centre de Recherches sur la Biologie des Populations d'Oiseaux, CP 51, 55 Rue Buffon, 75005 Paris, France; ²Department of Biology, McGill University, 1205 Docteur Penfield, Montréal, QC, Canada; and ³Laboratoire d'Ecologie Alpine, UMR-CNRS 5553, Université Joseph Fourier, Grenoble I, BP 53, 38041 Grenoble Cedex 9, France

Recommend ~ 10 000

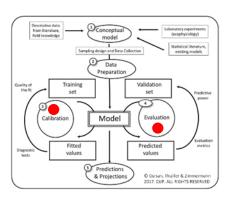
Train and test data



 Common to test how well predictions match observations (model validation)

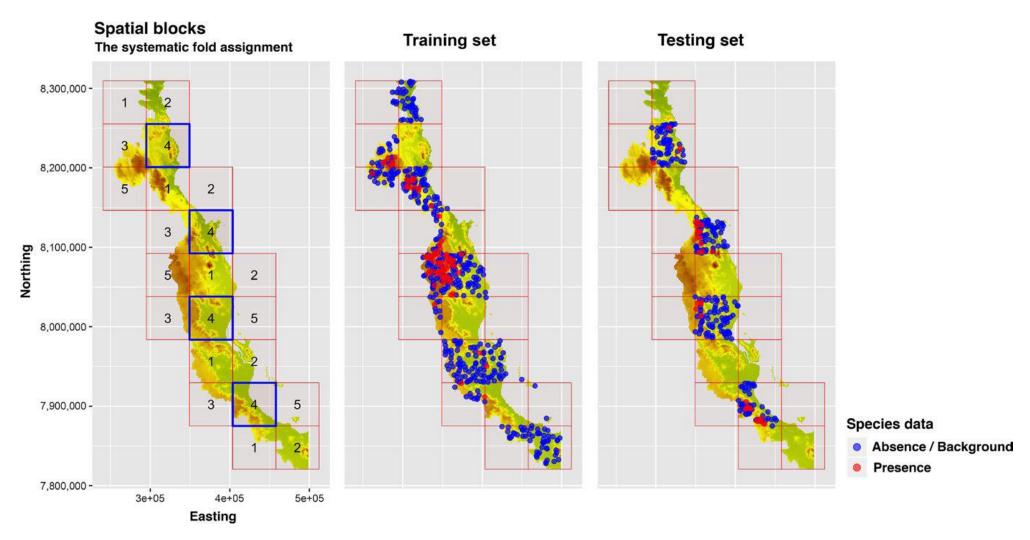
- Train (calibration) = data used to create model
- Test (validation) = data used to evaluate predictive performance

- Ideal case is to use independent data
- MaxEnt uses internal data split



Train and test data





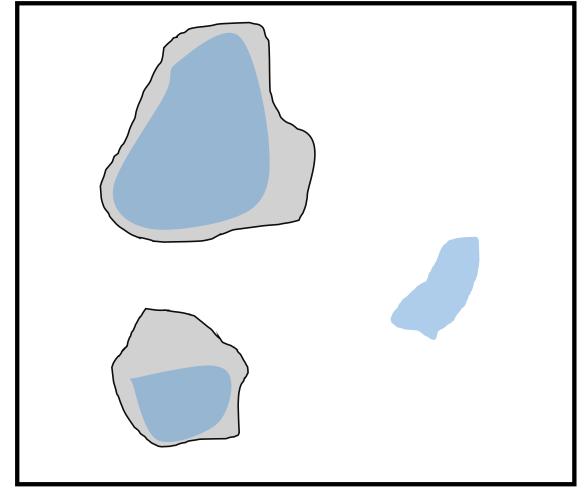
Model evaluation – confusion matrix



	Observed present	Observed absent
Predicted present	True presence (TP)	False presence (FP)
Predicted absent	False absence (FA)	True absence (TA)

Latitude





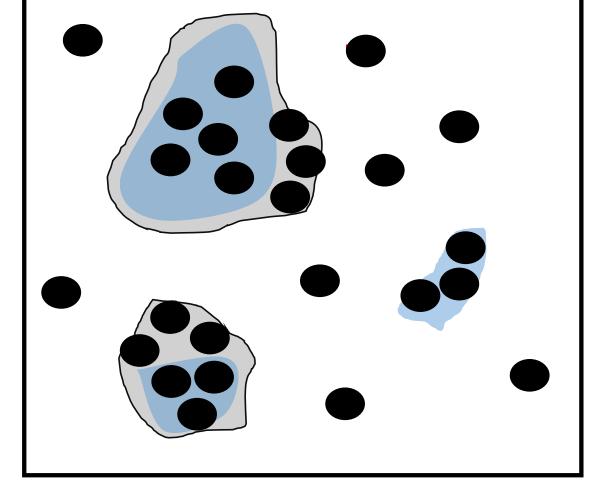
Actual distribution

Species distribution model

Longitude



Latitude





Actual distribution



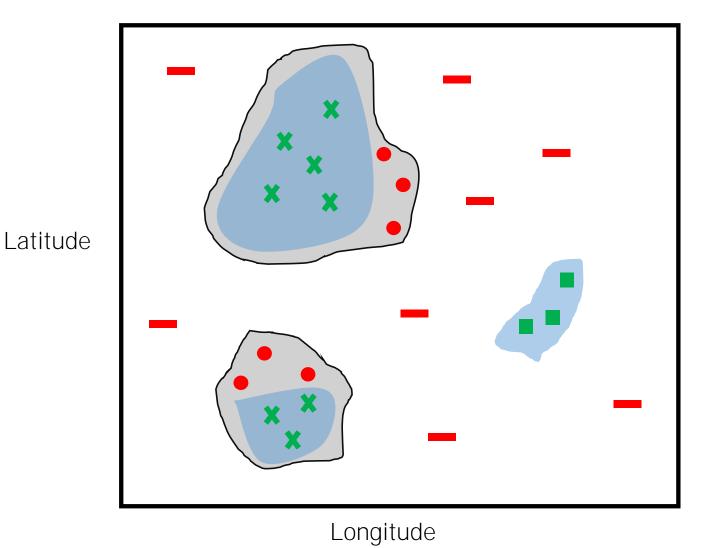
Species distribution model



Presence/absence

Longitude





Actual distribution

Species distribution model

- True absence
- False absence
- **X** True presence
- False presence



Sensitivity = TP /
$$(TP + FA) = 9/(9+3) = 0.75$$

Specificity =
$$TA / (TA + FP) = 13/(13 + 2) = 0.86$$

Ohsarvad prasant

Predicted present

Predicted absent

obscrved present	
9	2
3	13

Ohsarvad ahsant

Evaluation metrics - AUC



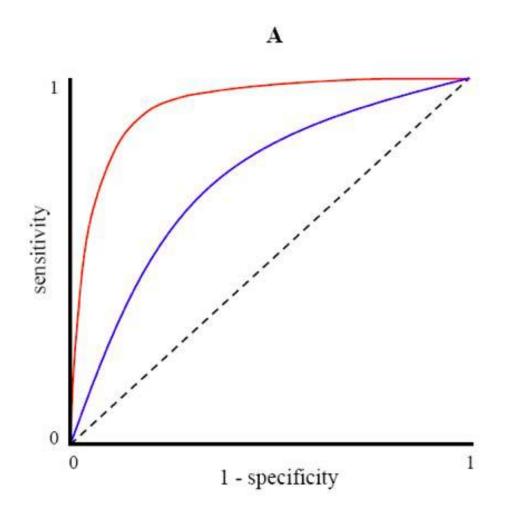
 Most common is the AUC test statistic (area under the receiver operating characteristic curve)

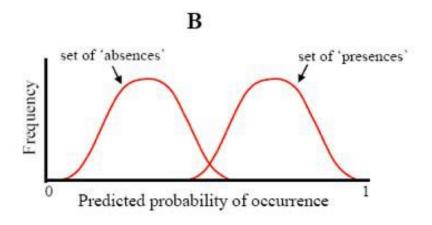
 AUC derived from Receiver Operating Characteristic (ROC) curve

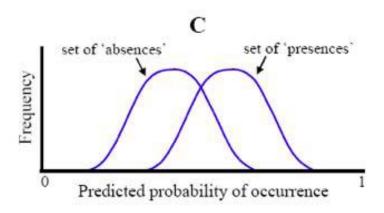
 ROC is defined by plotting sensitivity against '1 – specificity' across the range of possible thresholds

Evaluation metrics









Evaluation metrics



- Calculate the surface below the curve (area under curve)
- Maximum value is obtained when the curve goes through the upper left corner – corresponding to a value of 1 and the lowest values with curve follows 1:1 line
- AUC interpretation
 - AUC > 0.9 = "excellent"
 - 0.80 < AUC < 0.9 = "good"
 - 0.7 < AUC < 0.8 = "fair"
 - 0.6 < AUC < 0.7 = "poor"
 - AUC < 0.6 = "fail"

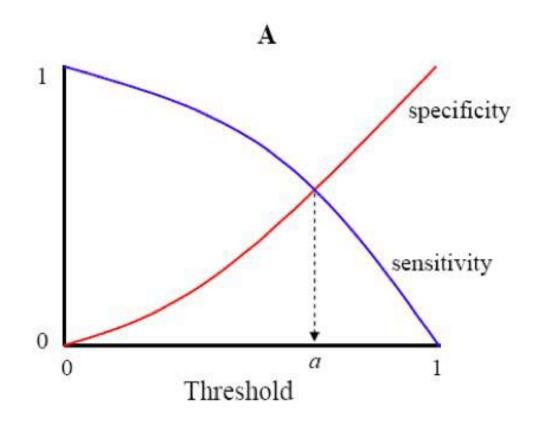
Thresholds of occurrence

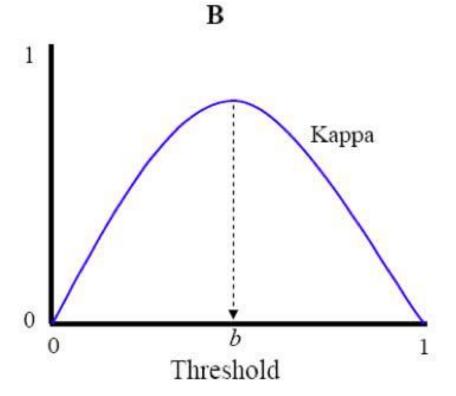


Fixed threshold approach	Taking a fixed value, usually 0.5, as the threshold
Kappa maximization approach	Kappa statistic is maximized
Average probability/suitability approach	Taking the average predicted probability/suitability of the model-building data as the threshold
Sensitivity-specificity sum maximization Approach (max SSS)	The sum of sensitivity and specificity is maximized
Sensitivity-specificity equality approach	The absolute value of the difference between sensitivity and specificity is minimized
ROC plot-based approach	The threshold corresponds to the point on ROC curve (sensitivity against 1- specificity) which has the shortest distance to the top-left corner (0,1) in ROC plot

Thresholds of occurrence

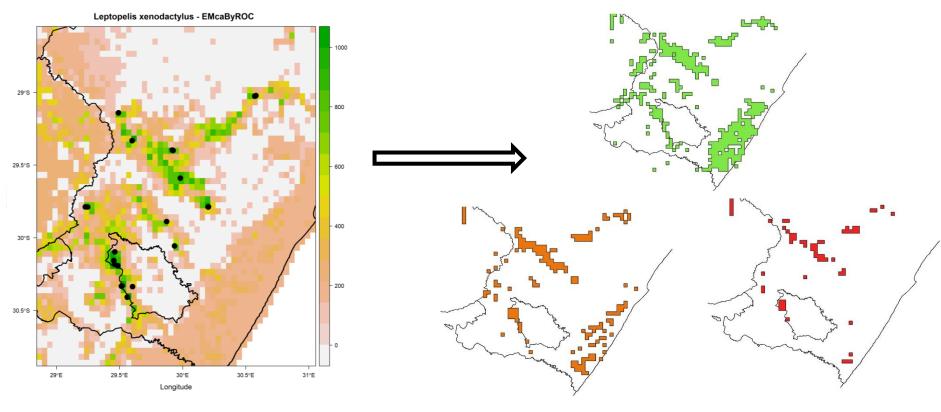






Thresholds of occurrence





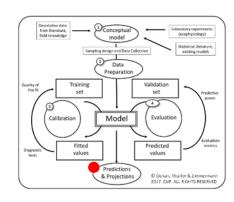
Probability surface

Binary surface

Predictions



- Predicting in novel environments
- Predicting within know range of a species
- Predicting in space and time (climate change scenarios)



Predictions



