

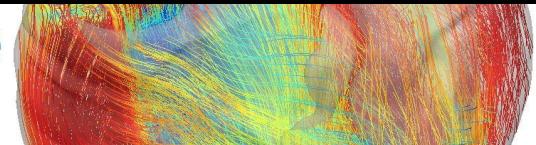


Modeling species spatial distributions in ecology

André Luís Luza (andre.luis-luza@u-bordeaux.fr)

23/01/2025

**Numerics – 2024/2025 – Semaine
« Numérique et transitions »**



Topics

- Introduction
- Sources of Distribution Information
- Models and Data
- Site Occupancy Models and Data
- Interactive Activity
- Wrap-Up and Q&A
- Concluding remarks

Introduction

Fundamental quantities/state variables:

- Species occurrence (z , z_i)
- Species abundance (N , N_i)



Involve field observations and some sort of spatial model (SDMs)

Introduction

Distribution: manner in which a biological taxon is spatially arranged

- individual property

The way we perceive distribution can change in function of taxonomic, temporal, and spatial aspects.



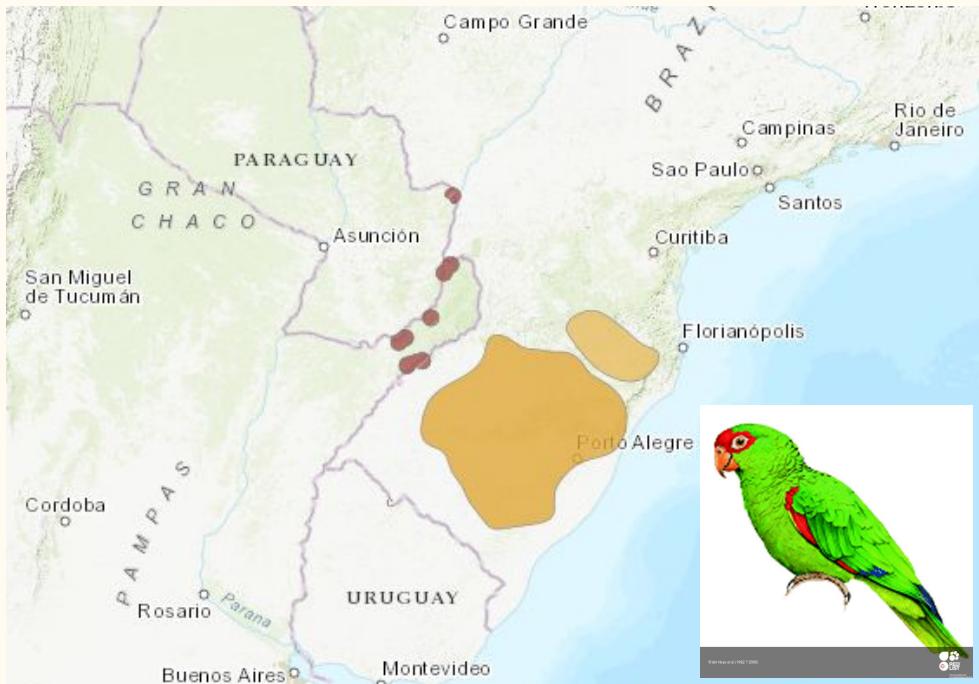
The Greater rhea (*Rhea americana*) range map (Source: IUCN, International Union for Conservation of Nature)

Introduction

Distribution: manner in which a biological taxon is spatially arranged

The way we perceive distribution can change in function of taxonomic, temporal, and spatial aspects.

Distribution can be large or small (endemic species)



The Red-spectacled Amazon (*Amazona pretrei*) range map (Source: IUCN)

Introduction

Distribution: manner in which a biological taxon is spatially arranged

The way we perceive distribution can change in function of taxonomic, temporal, and spatial aspects.

Distribution can be large or small (micro endemic species)



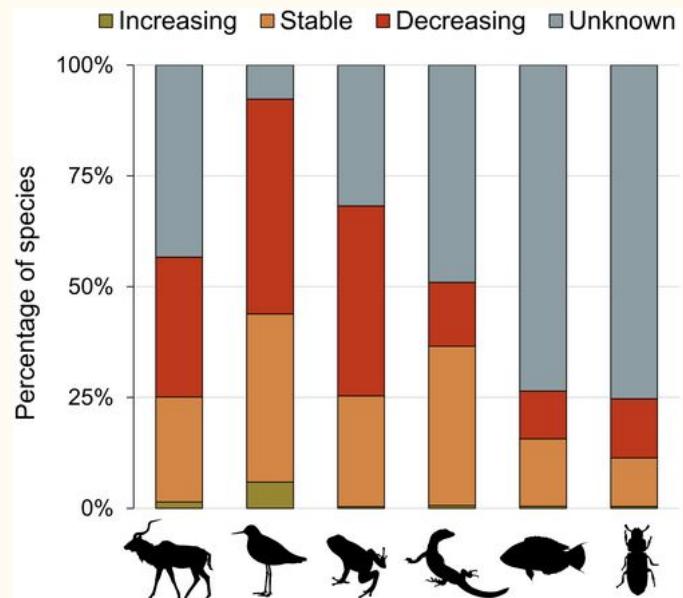
The Admirable-Redbelly-toad (*Melanophryne admirabilis*) range map (Source: IUCN)

Introduction

We're in the midst of a biodiversity crisis

- Fast extinction rates
- Widespread and rapid land use change: *suitable habitats are becoming rarer, smaller and isolated*
- Disease outbreaks and spread of invasive species
- Lost of unknown populations

Knowledge about species distribution
can help to provide solutions
Inform conservation (IUCN)



Source: Finn et al. 2023.
10.1111/brv.12974

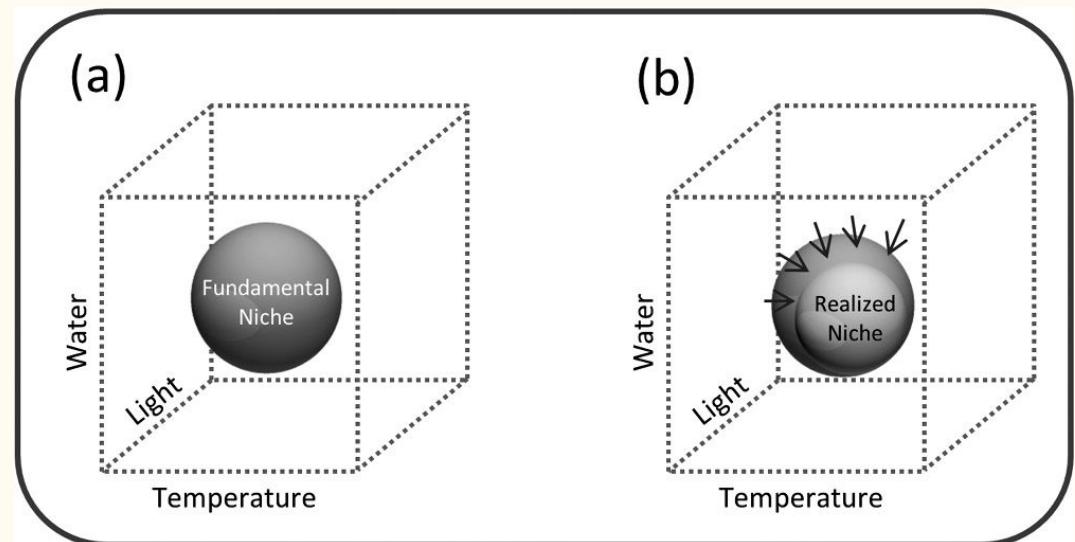
Sources of Distribution Information

Early assessments of species distribution:

- Non-spatial **niche** projections based on empirical (presence) data

Ecological niche: species' responses to a suite of environmental conditions

Non-spatial, non-statistical



(Source: Guisan et al. 2017)

Sources of Distribution Information

Early assessments of species distribution:

- Expert knowledge
 - Hand made range/distribution maps
 - What's inside? And beyond?
 - Do the edges really exist?

Spatial, non-statistical



The Greater rhea (*Rhea americana*) range map (Source: IUCN)

Sources of Distribution Information

Modern assessments (80's) of species distribution:

- Statistical models used to estimate ψ_i , based on the relationships between environmental values and species presence and absence (in some cases).
 - *Species distribution models, bioclimatic models, climate envelopes, ecological niche models, habitat suitability models*
- Abundance maps, occurrence probability maps

Spatial, statistical



Greater rhea abundance map (Source: eBird)

Principles of Species Distribution Models

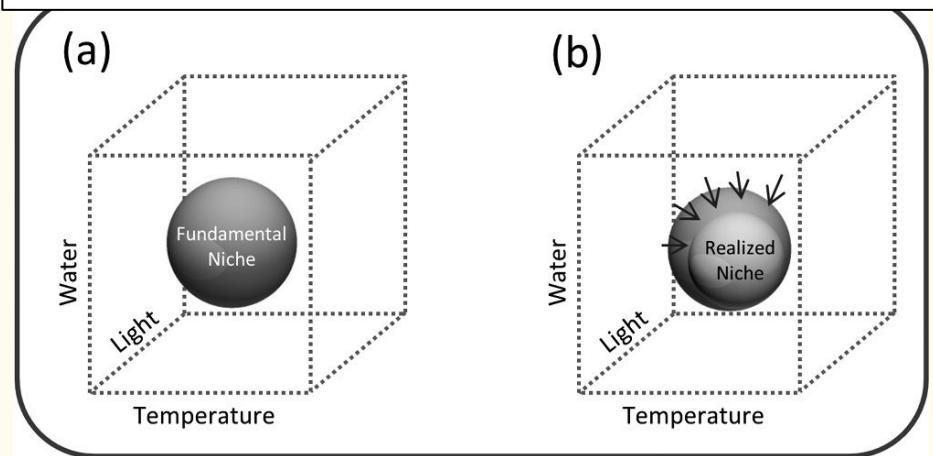
Let $z_i \in \{0,1\}$ be the (latent) realized/true presence of one focal species in $i=1, \dots, i=M$ sites.

To $z_i = 1$, three general conditions need to be met

- Dispersal
- **Habitat suitability (ψ)**
- Biotic environment

$$z_i \sim \text{Bernoulli}(\psi)$$

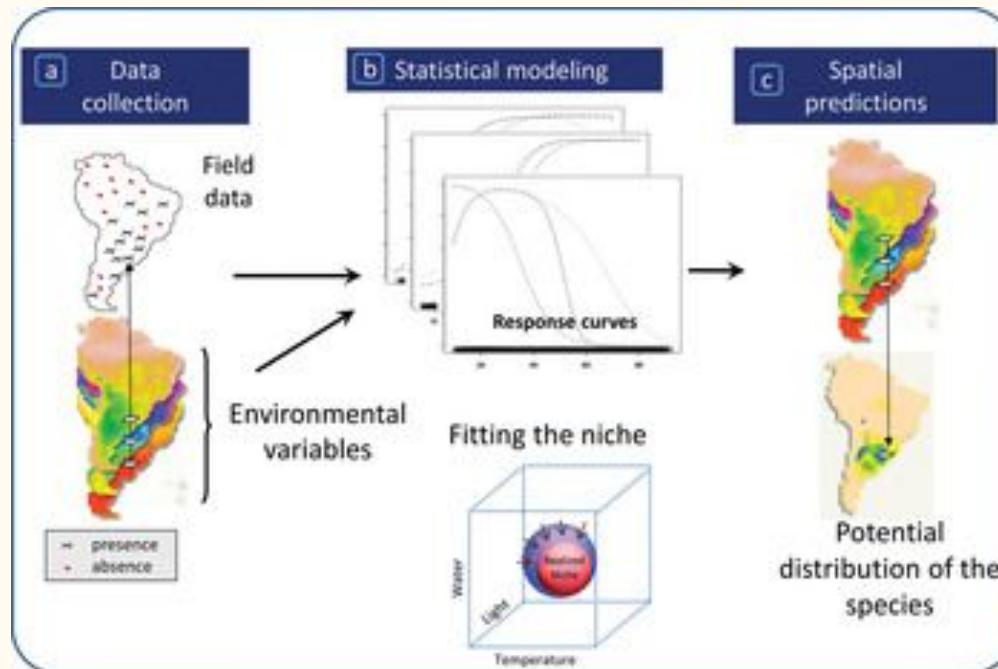
Project the niche into the geographic space



(Source: Guisan et al. 2017: 10.1017/9781139028271)

Principles of Species Distribution Models

Combine knowledge from ecology, natural history,
remote sensing (GIS) and statistics



Once established, the model can be used to make:

- Inference about determinants
- Predictions (interpolation, extrapolation, hindcast, forecast)

(Source: Guisan et al. 2017: 10.1017/9781139028271)

(Hierarchical)Models and Data

Let ψ be the probability the site is occupied, p the detection probability in truly occupied sites (assumed to be $p=1$ for now), and z_i as before

$z_i \sim \text{Bernoulli}(\psi)$ (ecological model)

$y_i | z_i \sim \text{Bernoulli}(z_i p)$ (observation model)

Let $y_i = (y_1, y_2, y_3, \dots, y_M)$ be the set of observations taken at the M sites, which could be

$$y_i = (0, 1, 0, 0, 1, 1, 0, 1) = z_i$$

$$L(y_i) = \begin{cases} \psi, & \text{if } y_i = 1 \\ 1 - \psi, & \text{if } y_i = 0 \end{cases}$$

Pigmy rice rat
(*Oligoryzomys*)



(Hierarchical)Models and Data

Let ψ be the probability the site is occupied, p the detection probability in truly occupied sites (assumed to be $p=1$ for now), and z_i as before

$z_i \sim \text{Bernoulli}(\psi)$ (ecological model)

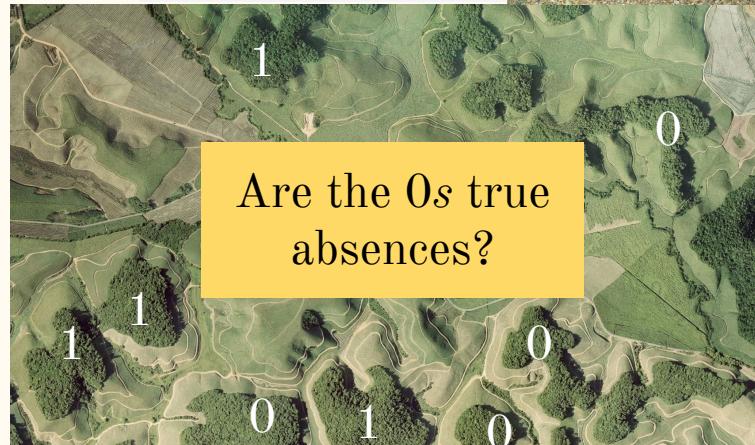
$y_i | z_i \sim \text{Bernoulli}(z_i p)$ (observation model)

Let $y_i = (y_1, y_2, y_3, \dots, y_M)$ be the set of observations taken at the M sites, which could be

$$y_i = (0, 1, 0, 0, 1, 1, 0, 1) = z_i$$

$$L(y_i) = \begin{cases} \psi, & \text{if } y_i = 1 \\ 1 - \psi, & \text{if } y_i = 0 \end{cases}$$

Pigmy rice rat
(*Oligoryzomys*)

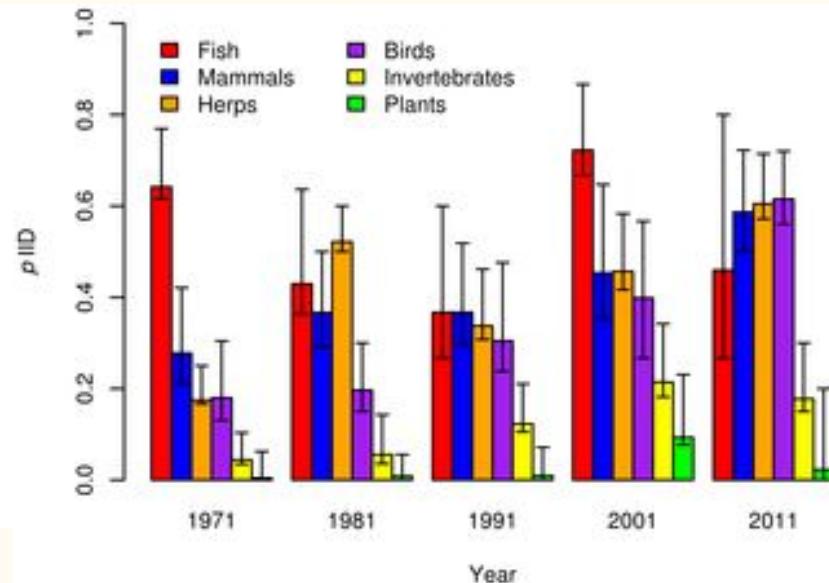


(Hierarchical) Models and Data

The predominant literature on **SDMs** (GLM, Maxent, RF, GAM) assumes **perfect** and **constant detection**:

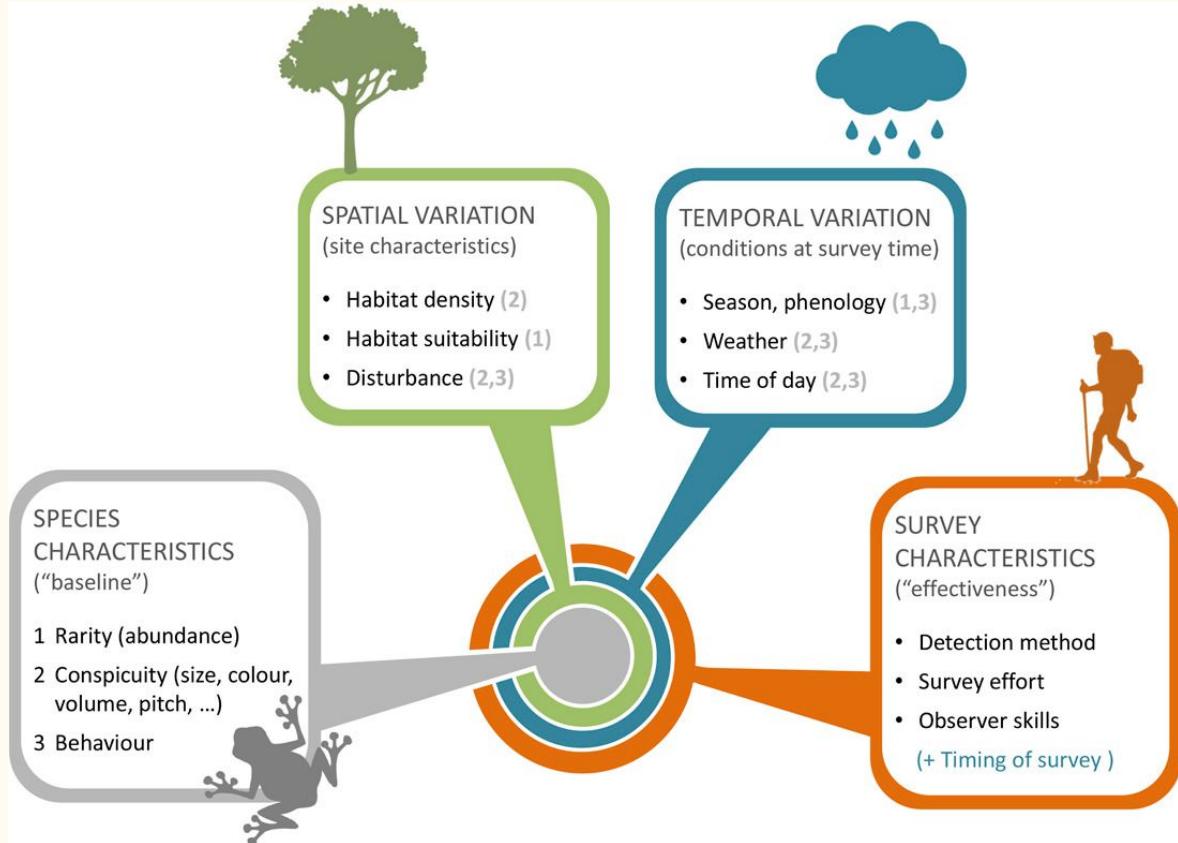
Kellner & Swihart (2014): Review about the consideration of imperfect detection in ecology:

- $23\% \pm 1.8$ (123/537) of articles addressed imperfect detection
- Range of single-survey detectability were 0.29 and 0.71
- 70% of the studies had at least one estimated detection probability per survey < 0.5



(Hierarchical) Models and Data

Which factors affect species detection in truly occupied sites?



(Hierarchical)Models and Data

Let ψ , the probability the site is occupied, p the detection probability in truly occupied sites (now acknowledge that $p < 1$), and z_i as before

$z_i \sim \text{Bernoulli}(\psi)$ (ecological model)

$y_i | z_i \sim \text{Bernoulli}(z_i p)$ (observation model)

Let $y_i = (y_1, y_2, y_3, \dots, y_M)$ be the set of observations taken at the M sites, which could be

$$y_i = (0, 1, 0, 0, 1, 1, 0, 1)$$

ψ and p not identifiable
(only ψp)

Apparent species distribution

$$L(y_i) = \begin{cases} \psi \cdot p, & \text{if } y_i = 1 \\ \psi(1 - p) + (1 - \psi), & \text{if } y_i = 0 \end{cases}$$



(Hierarchical) Models and Data

To make ψ and p identifiable we must conduct $j=1, \dots, j=J$ surveys to all / subset of the sites.
Consider now a **detection/non-detection occupancy data**, with $J=4$ replicated surveys (days)

$z_i \sim \text{Bernoulli}(\psi)$ (ecological model)

$y_{ij} | z_i \sim \text{Bernoulli}(z_i p)$ (observation model)

Let $y_{ij} = (y_{1,1}, y_{1,2}, y_{1,3}, \dots, y_{M,J})$ be the set of observations taken at the M sites at J surveys. For a single site it could be

$$y_{i=1,J} = (1, 0, 0, 1),$$

$$\text{L}(y_{i=1,J} | z_{i=1}=1) = \psi \cdot p^J \cdot (1-p)^J$$

$$y_{i=2,J} = (0, 0, 0, 0),$$

$$\text{L}(y_{i=2,J} | z_{i=2}=1) = \psi \cdot (1-p)^J$$

$$\text{L}(y_{i=2,J} | z_{i=2}=0) = 1-\psi$$

For any given survey budget, there is a trade-off between sampling more sites, and applying more effort per site.



(Hierarchical)Models and Data

Heterogeneous occupancy data

Let $y_{ij} = (y_{1,1}, y_{1,2}, y_{1,3}, \dots, y_{M,J})$ be the set of observations taken at the M sites at J surveys. For a single site it could be

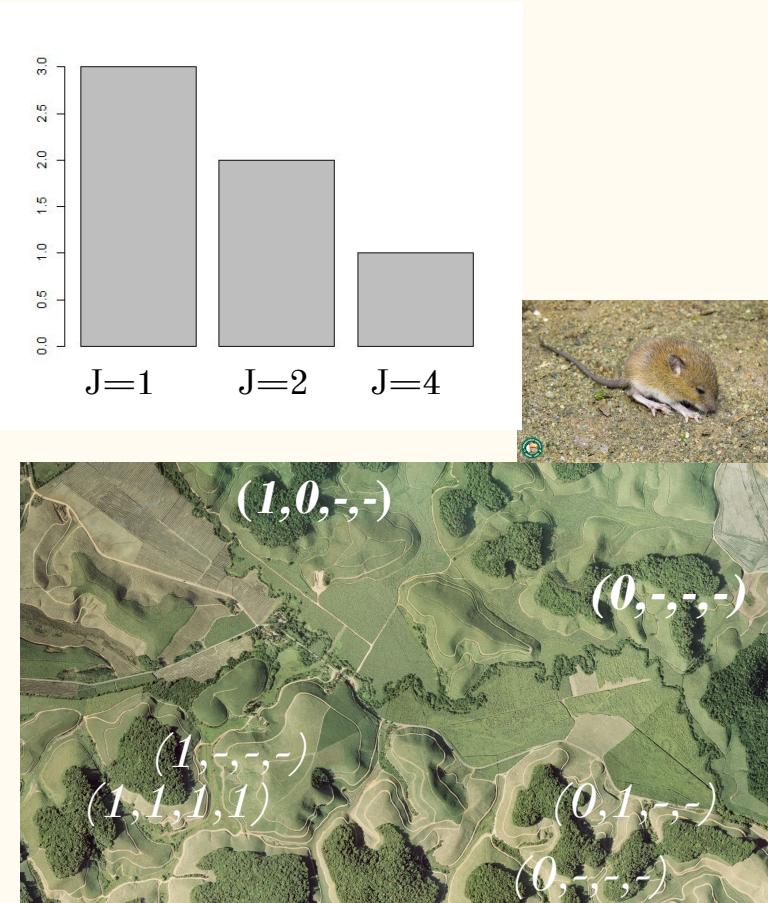
$$y_{i=1,J} = (1, 0, -, -),$$

$$\text{L}(y_{i=1,J} | z_{i=1} = 1) = \psi \cdot p \cdot (1-p)$$

$$y_{i=2,J} = (0, -, -, -),$$

$$\text{L}(y_{i=2,J} | z_{i=2} = 1) = \psi \cdot (1-p)$$

$$\text{L}(y_{i=2,J} | z_{i=2} = 0) = 1 - \psi$$

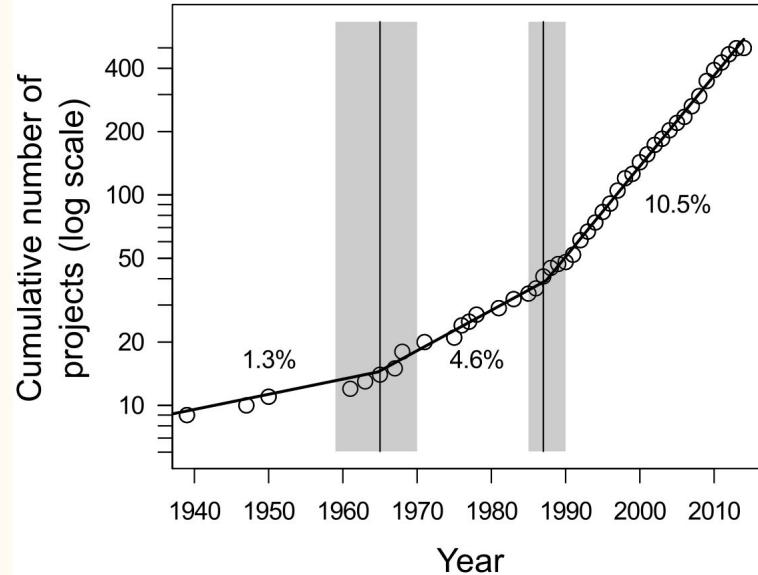


(Hierarchical)Models and Data

To model spatial distribution we need more data, right?

Pooling opportunistic data: Planned data + Museum data
+ Citizen-science projects (broad niche coverage)

- GBIF, iNaturalist, eBird
- Prevalence of non-sampled sites $y_i = \text{NA}$
- Presence-only data (lacking $y_i = 0$)
- When $y_i = 1$:
 - Limited temporal replication (J visits)
 - Biased (non-random) environmental sampling
 - Limited data about the observation process (covariates, V_{ij})
 - Big and heterogeneous data (varied detection histories)



Cumulative number of ecological and environmental citizen science projects as revealed by a systematic global search (Pocock et al. 2017,
10.1371/journal.pone.0172579)

(Hierarchical)Models and Data

There are three conditions for $y=1$ (ie. record coordinates in a database)

- The species must occur ($z=1$)
- **The site must be sampled/visited**
- **The species must be detected**

TABLE 1. List of challenges for the use of citizen science data for monitoring biodiversity. There are four broad categories and 10 individual challenges

Category	Challenge
Observer behaviour	Spatial bias
	Observer differences
	Reporting preferences
	False positive errors



(Hierarchical)Models and Data

How to obtain absences from presence-only data?

- Background absences (proportional)
- Absences in the neighborhood



(Hierarchical)Models and Data

How to obtain absences from presence-only data?



Month1-ObsB



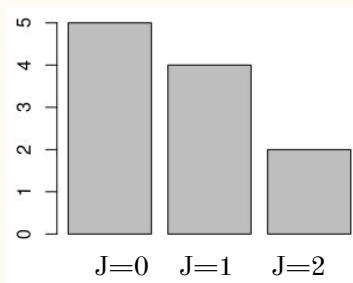
Johnson et al. 2023, 10.1111/2041-210X.13834

- Background absences
- Absences in the neighborhood
- Community data (target-group)



Month1-ObsA



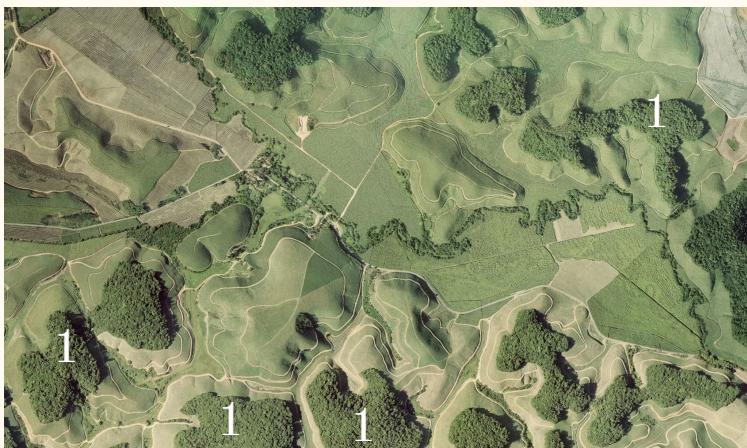


Month1-ObsB

Focal species



Month1-ObsA



Wrap-up of data

Available data sets can be:

Spatial coverage
Availability
Bias

(Gold) Standard
(planned repeats)

Heterogeneous (replicates in
some sites)

Single-visit (no repeat)

Presence-only

BINARY RECORDS

site 1	1	0	0	...	1
site 2	0	-	1	...	-
site 3	0	0	0	...	0
site 4	1	0	1	...	-
...
site S	0	0	1		1

BINARY RECORDS

site 1	1	0			
site 2	0				
site 3	0				
site 4	1	0	1	...	-
...
site S	0				

BINARY RECORDS

site 1	1				
site 2	0				
site 3	0				
site 4	1				
...	...				
site S	0				

BINARY RECORDS

site 1	1				
...	...				
site 4	1				
...	...				
site S	0				

(Hierarchical) Models and Data

Instead of just estimating ψ and p , we often want to estimate the effect of site/visit covariates on these parameters.

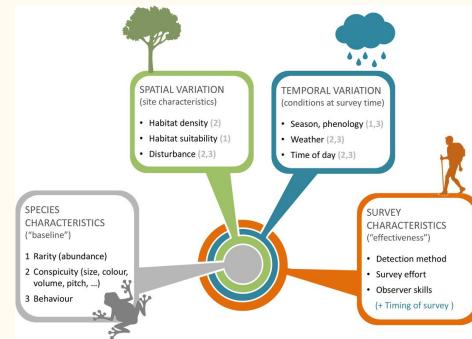
The logistic regression is a natural choice for variables whose values are [0,1],

$$\log (\psi_i / (1 - \psi_i)) = X_i' \beta$$

$$\log (p_{ij} / (1 - p_{ij})) = V_{ij}' \alpha$$

β and α : vectors of coefficients

X_i and V_{ij} : site-level and site-/survey-level covariates



(Hierarchical) Models and Data

Instead of just estimating ψ and p , we often want to estimate the effect of covariates on occupancy and detection per site and/or visit.

$z_{it} \sim \text{Bernoulli}(\psi_{it})$ (ecological model)

$y_{itj} | z_{it} \sim \text{Bernoulli}(z_{it} p_{itj})$ (observation model)

$$\log(\psi_{it} / (1 - \psi_{it})) = X_{it}'\beta + \omega_{(s)} + \eta_t$$

$$\log(p_{itj} / (1 - p_{itj})) = V_{itj}'\alpha$$

β and α : vectors of coefficients

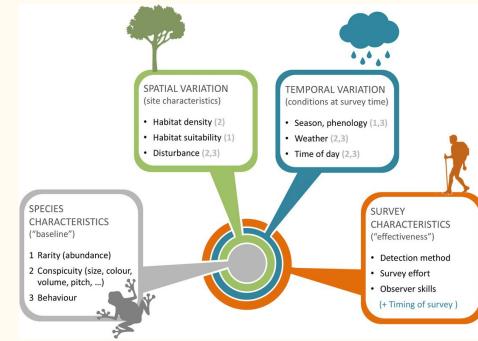
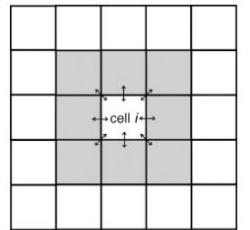
$\omega_{(s)}$: spatial random effect

η_t : temporal random effect

X_{it} and V_{itj} : site/year/survey-level covariates



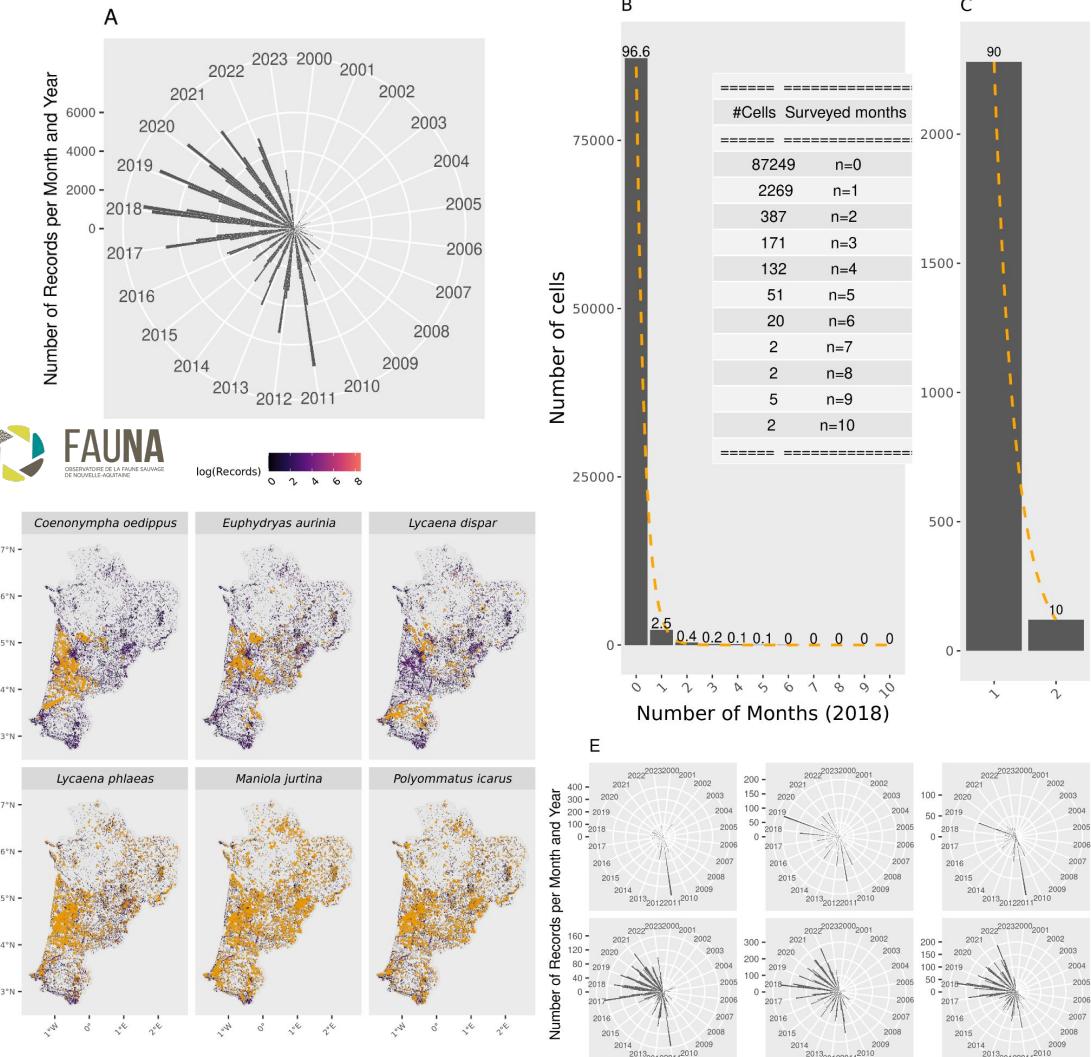
Spatial neighbors
of site i



(Hierarchical) Models and Data

Example: Nouvelle Aquitaine Butterflies

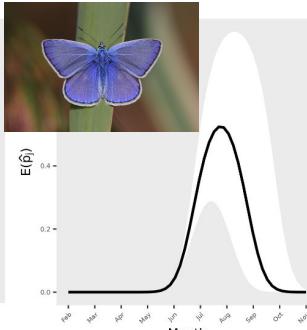
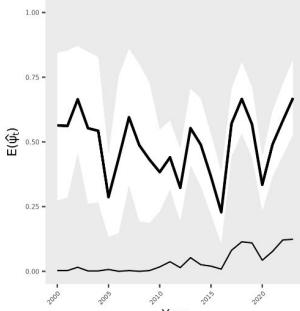
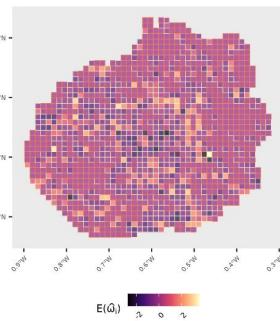
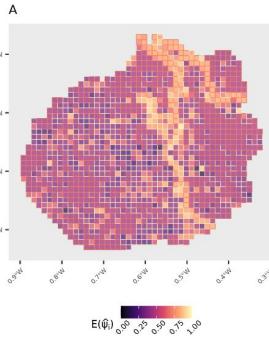
- A) Distribution of records over time
- B) Empirical distribution of visits to sites
- C) Distribution of visits to sites often used in simulations
- D) Spatial distribution of effort and records
- E) Distribution of records across years and months, for six species



Model application (results)

$$\log(\psi_{it} / (1 - \psi_{it})) = X_{it}'\beta + \omega_{(s)} + \eta_t$$

$$\log(p_{itj} / (1 - p_{itj})) = V_{itj}'\alpha$$

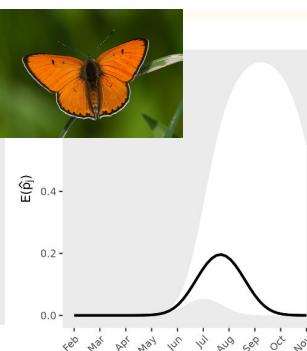
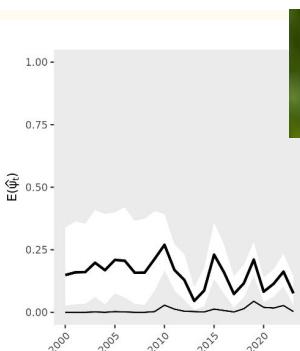
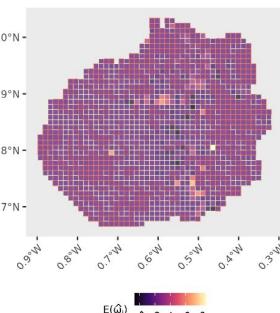
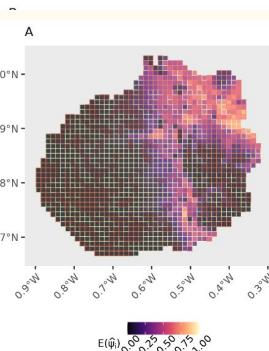


A) Predicted occupancy probability in a buffer of 10km² around Bordeaux

B) Spatial random effect

C) Temporal occupancy trends (since 2000)

D) Detection probability over months



Interactive Activity

We will use R Programming Environment to:

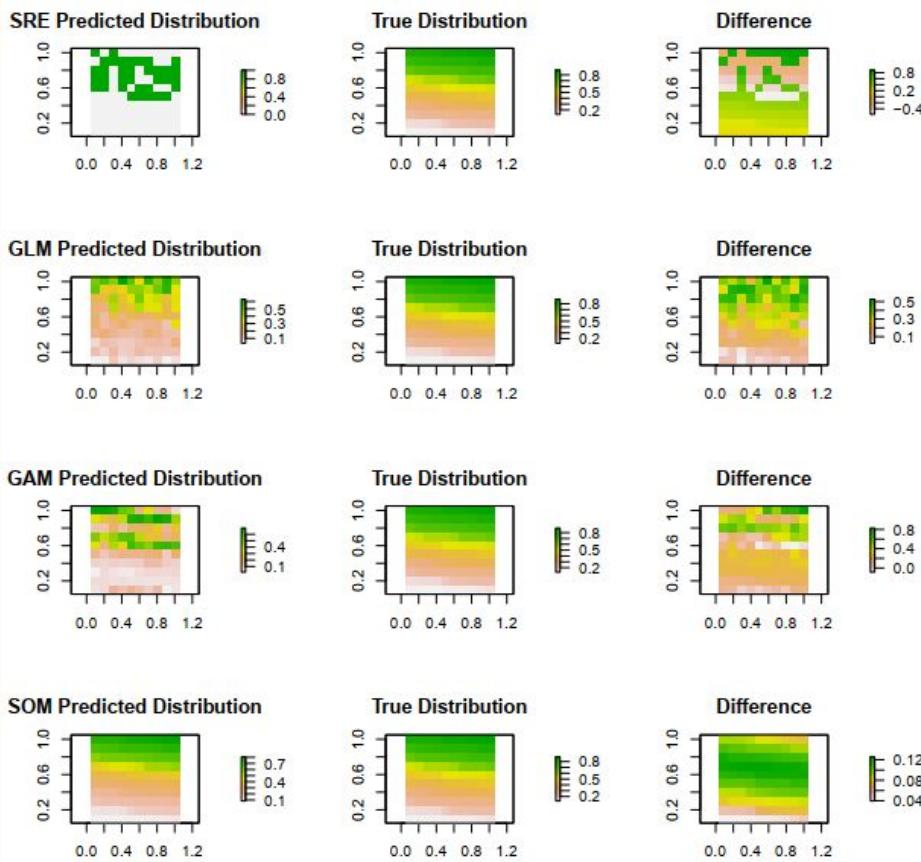
- Simulate some occupancy data (gold-standard data, single-visit/aggregated data)
- Analyze the data using:
 - BIOMOD, GLM, GAM, Hierarchical (Occupancy) Model
 - Create models without and with covariates
- Compare the performance of the different models

Wrap-Up and Q&A

Did you understand that the amount of information in data will determine the most suitable model?

Why is it important to account for imperfect detection?

- The extent of species distribution will be underestimated if $p < 1$;
- Estimates of covariate relationships will be biased to zero/low when $p < 1$;
- Factors that affect the observation may end up being in predictive models of occurrence or may mask the effect of factors that do affect species occurrence

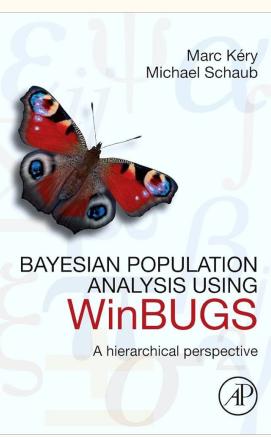
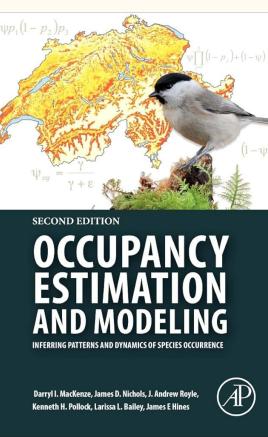
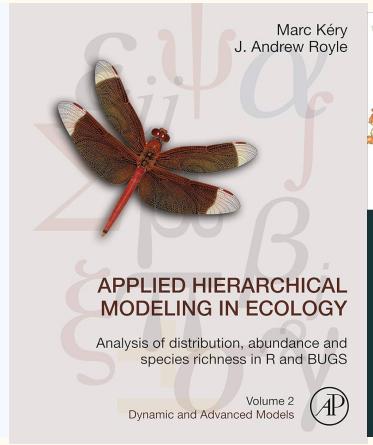
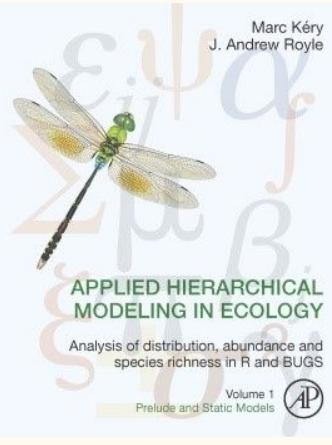
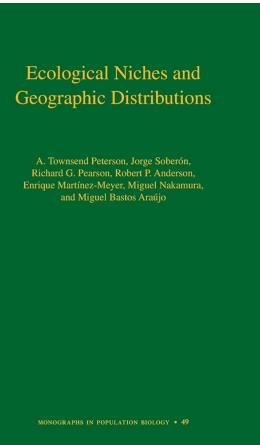
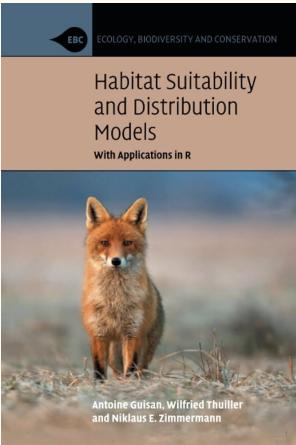


Concluding remarks

- Modern models acknowledge that the data can be generated by ecological and observation processes
 - Hierarchical formulation of SDMs enable to estimate occupancy/habitat suitability and observation probability
-
- If absences are not available, use the most suitable model and acknowledge that you're estimating observation probability/apparent distribution
- SOM enable to deal with several sources of bias in opportunistic data sets→ output useful for conservation proposes (IUCN AOO, extent)



Resources for learning



Tutorials/implementations:

<https://ebird.github.io/ebird-best-practices/intro.html>

<https://ebird.org/explore>

R Code: https://github.com/andreluza/SDM_talk

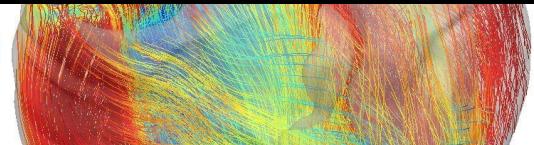


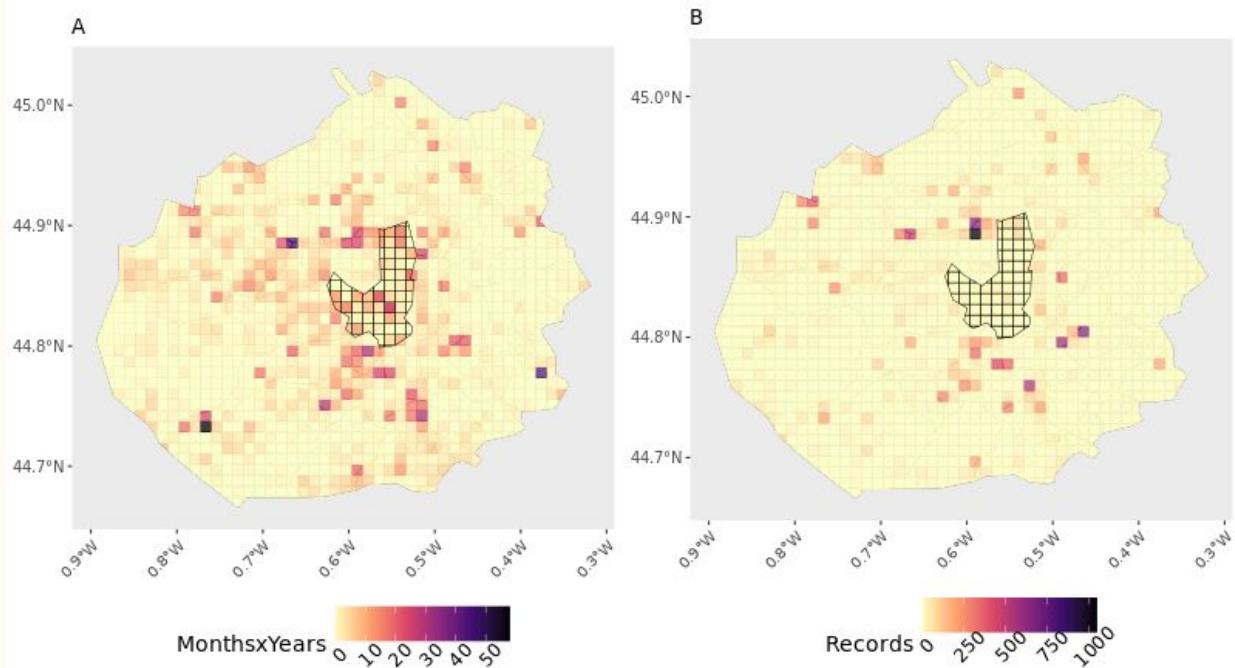
Modeling species spatial distributions in ecology

André Luís Luza (andre.luis-luza@u-bordeaux.fr)

23/01/2025

**Numerics – 2024/2025 – Semaine
« Numérique et transitions »**





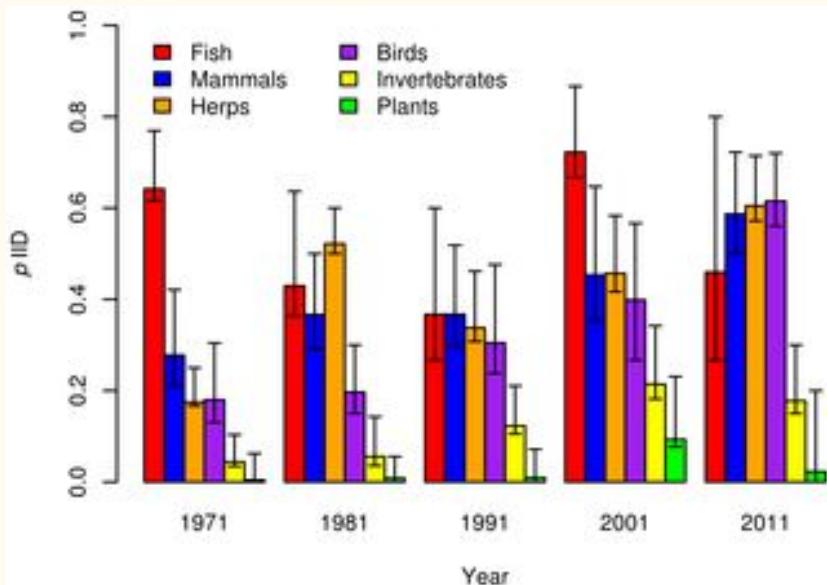
Species distribution models (SDMs)

The predominant literature on **classic SDMs** assumes perfect and constant detection:

$$z_i = y_i$$

Perfect detection is a very strong assumption - observations are not a perfect description of the reality

Classic SDMs model apparent distribution (product of ψ_i and p)

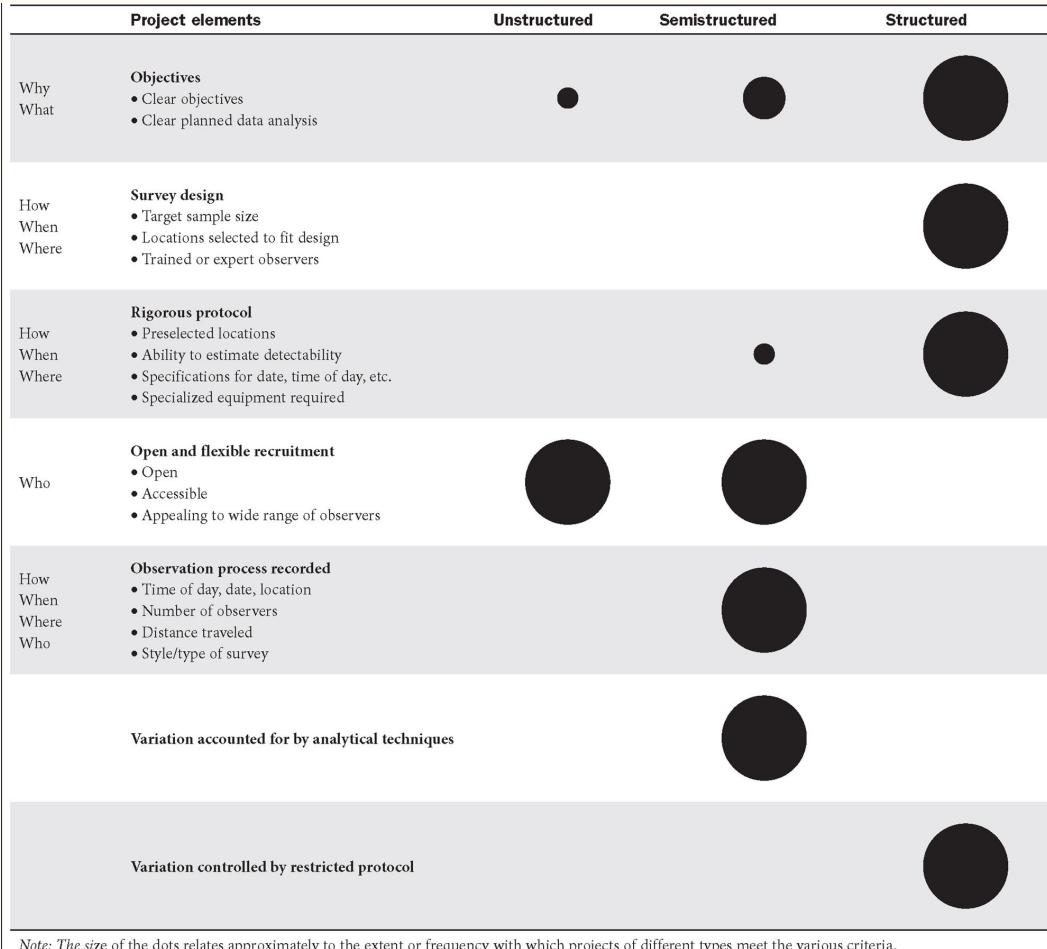


Site occupancy model: An hierarchical SDM

Assumptions

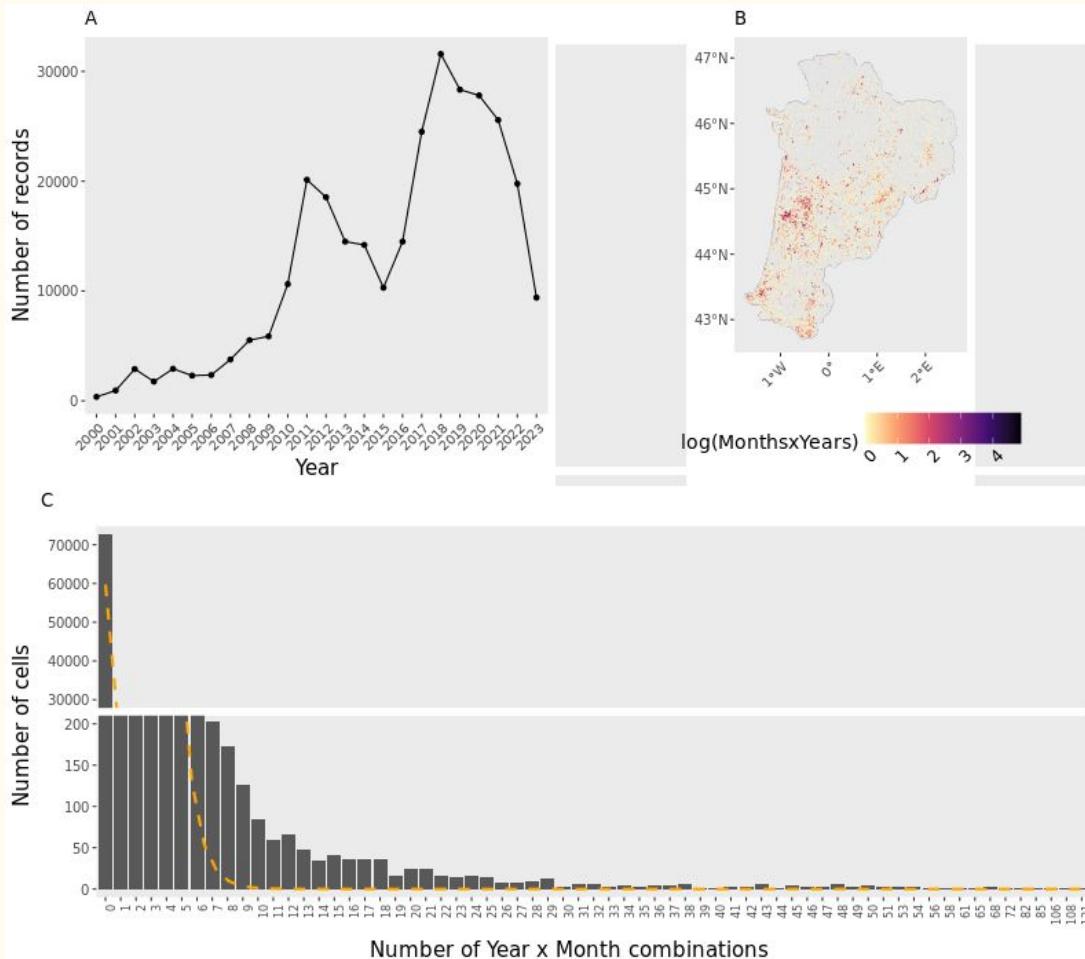
- *Closure*: over the duration of the j visits, z_i must not change
 - If temporary emigration happens, ψ_i = proportion of sites used
 - If colonization and extinction happen, we can “open” the model (multi-season model with closure within seasons)
 - Some advocate the single-visit data guarantee closure (Lele et al. 2012)
- *No false positives* ($z_i=1$ is true)
 - 100% sure the species is there (no misidentification)
- *Homogeneous detection across sites*
 - If heterogeneous, covariates capture variation on p_{ij} (otherwise ψ_i will be underestimated)

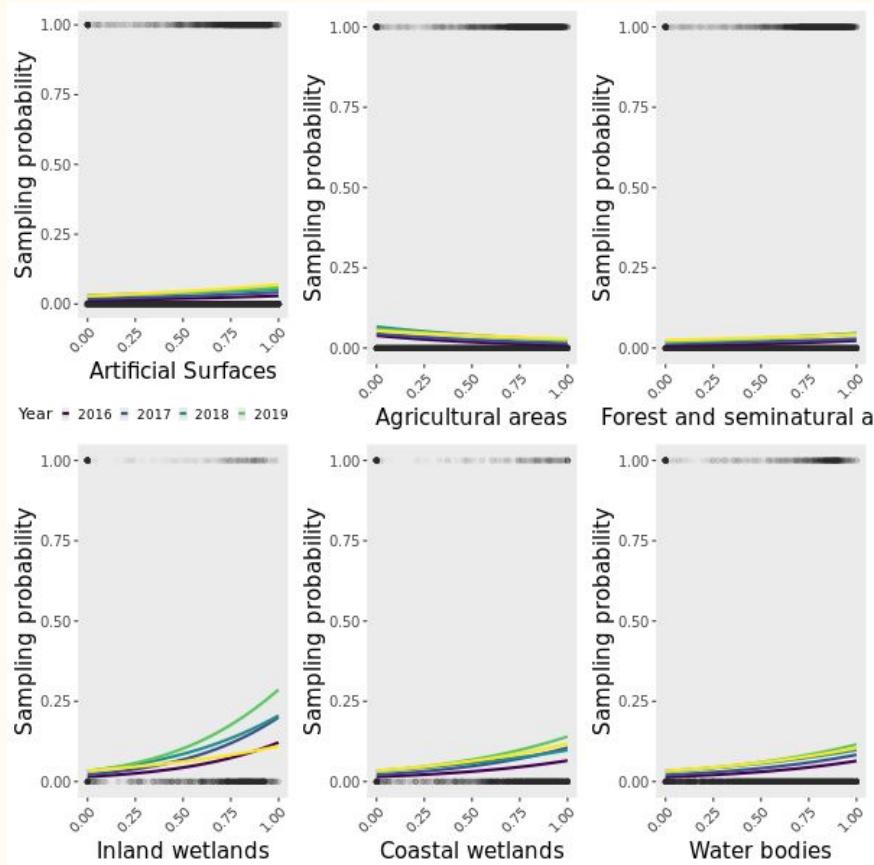
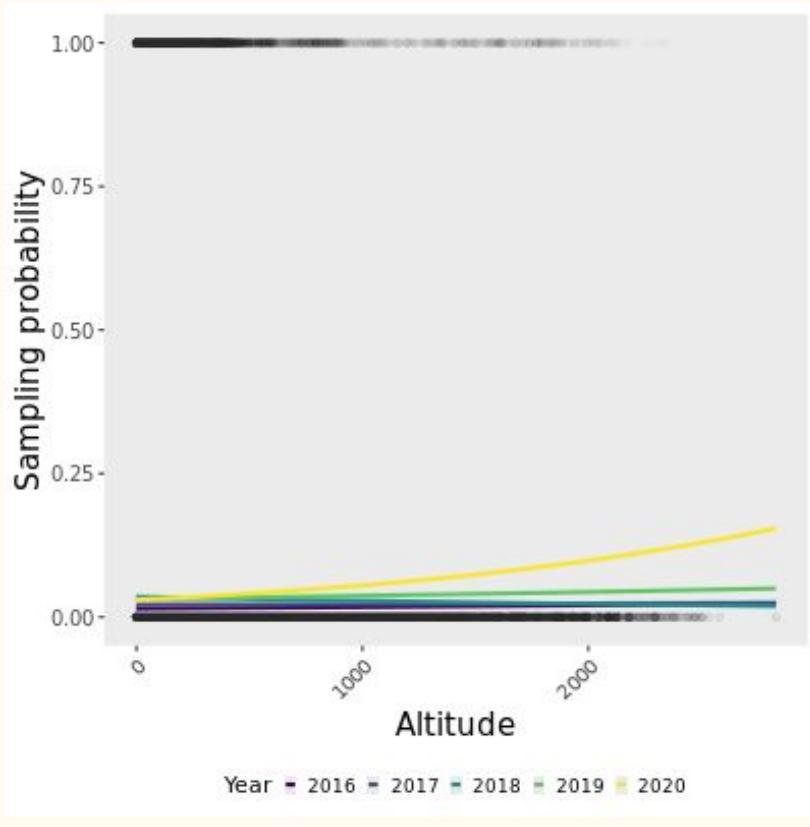
Use of site occupancy models to map geographic distributions



Note: The size of the dots relates approximately to the extent or frequency with which projects of different types meet the various criteria.

Review about the characteristics of unstructured, semistructured and structured citizen science projects (Kelling et al. 2019, BioScience)





Why does this matter?

- Allocating efforts of conservation and population monitoring/management
- Definition of suitable habitats
- Control of invasive species, disease vectors
 - Underestimation of distribution in time (p ignored, ϵ is overestimated)
 - Underestimation of the uncertainty

Temporal dynamics in the occupancy of two *Aedes* species in Rio Grande do Sul

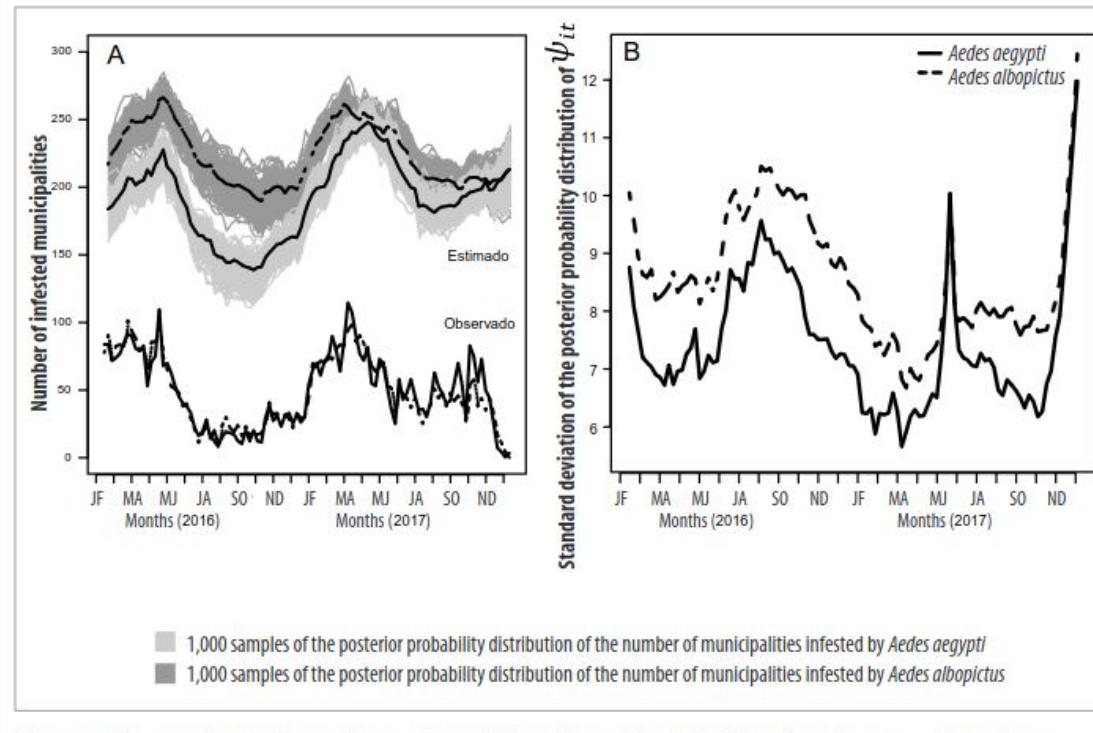


Figure 1 – Temporal variation on the number of infested municipalities (A) and on the uncertainty about infestation (B) by *Aedes aegypti* and *Aedes albopictus*, Rio Grande do Sul, 2016-2017

Why does this matter?

- Allocating efforts of conservation and population monitoring
- Definition of suitable habitats
- Control of invasive species, disease vectors
 - Underestimation of the spatial distribution

Spatiotemporal variation in distribution of two *Aedes* species in Rio Grande do Sul

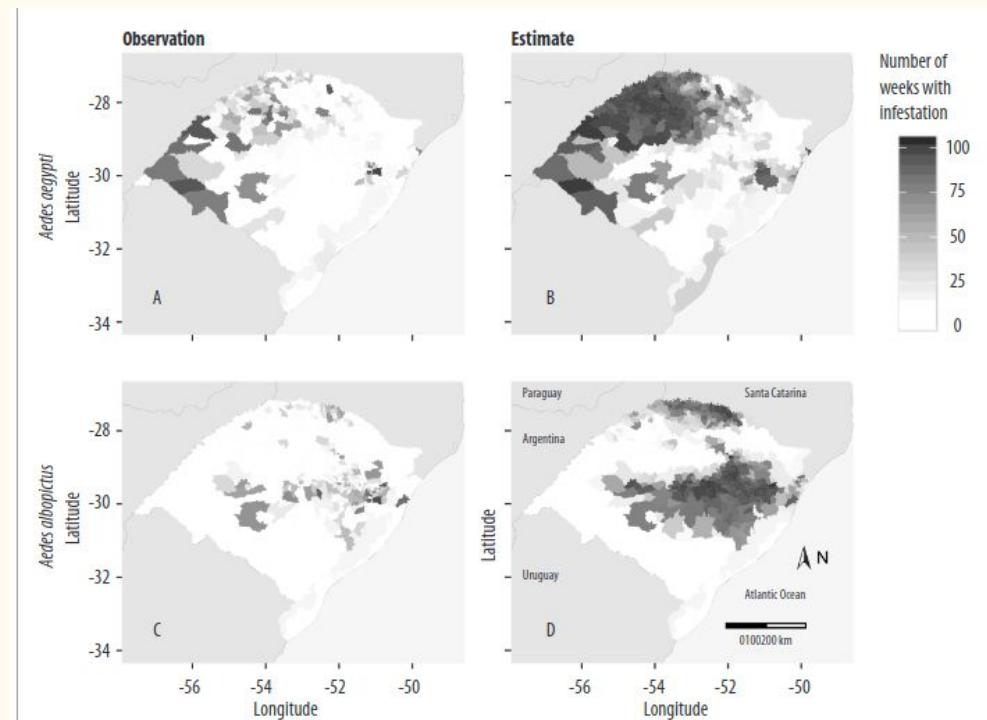
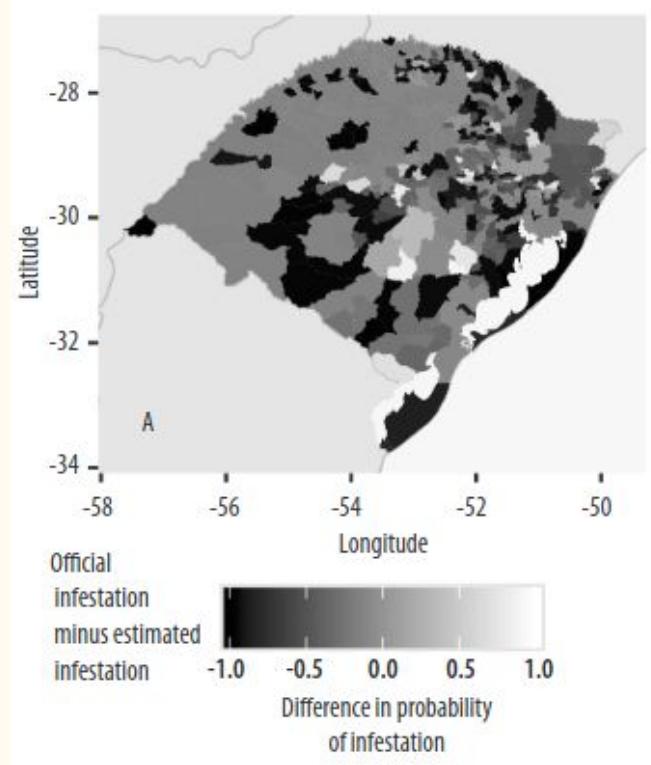


Figure 2 – Geographic distribution of *Aedes aegypti* (A and B) and *Aedes albopictus* (C and D), based on number of weeks with at least one detection of the species (A and C) and with presence of the species according to the model (B and D), Rio Grande do Sul, 2016-2017

Why does this matter?

- Allocating efforts of conservation and population monitoring
- Definition of suitable habitats
- Control of invasive species, disease vectors
 - Attribute low disease risk to municipalities with high ψ_i
 - Autochthonous dengue cases in municipalities not officially infested

Comparison between binary classification and probabilistic classification of infestation by *Aedes aegypti*

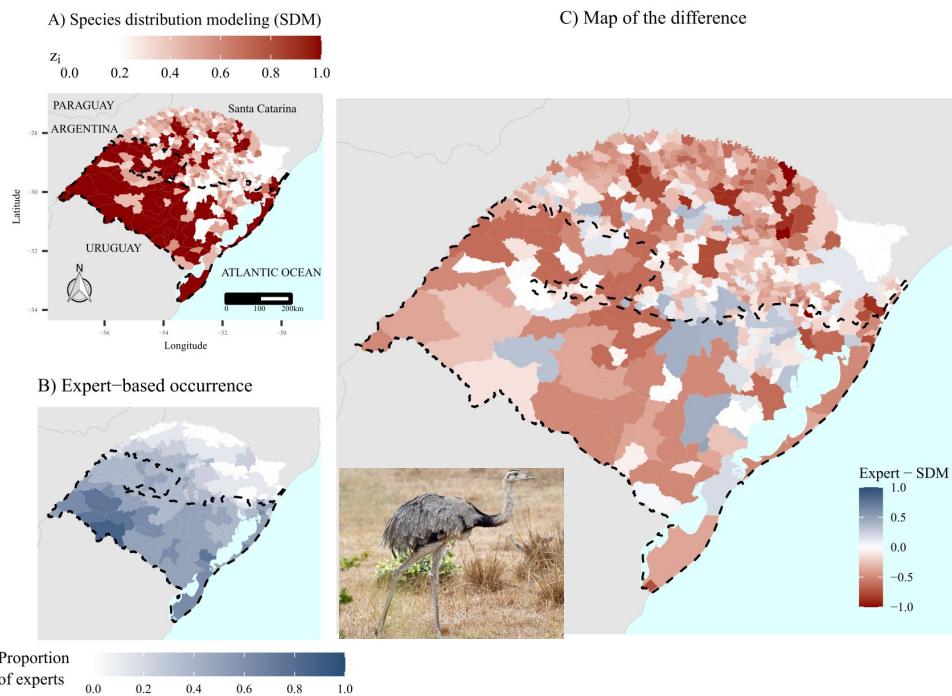


Combination of Sources of Distribution Information

Comparison of sources of distribution data can be informative:

- Prior knowledge
- Where + research is needed
- Suitable vs. unsuitable regions
- Where we need to improve the models

Example: The greater rhea distribution in Rio Grande do Sul, BR, using expert knowledge and SDMs fitted to citizen-science data (Luza et al. 2023. 10.1007/s43388-023-00143-3)

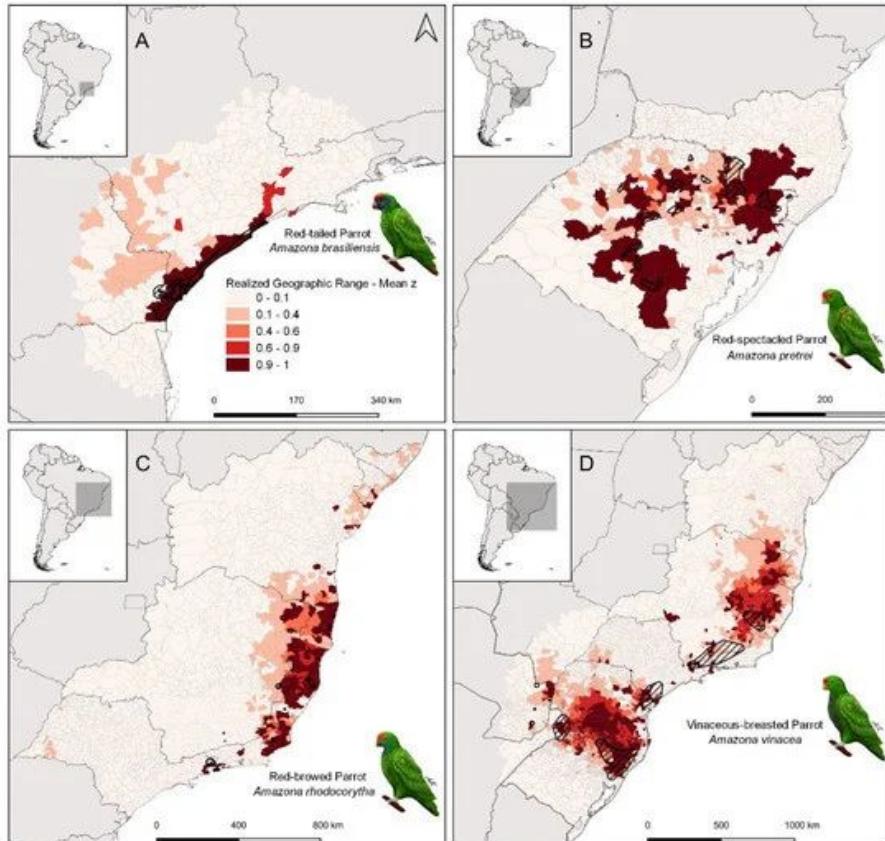


Information might not match (which is interesting)

Statistical tool (regression model) used to fill gaps in knowledge about species distribution and abundance (+demographic and derived parameters)

Often used to make predictions at broad scales (region +).

Parameter of interest for SD modelers: ψ_i or ψ_{it} , the site occupancy probability in site i and/or time t



Distribution of parrot species in Southern Brazil - estimates VS. IUCN (Zulian et al. 2021)