# Interpretable ML for biodiversity

An introduction using species distribution models

Timothée Poisot          Université de Montréal

September 29, 2024

1. How do we produce a model?
2. How do we convey that it works?
3. How do we talk about how it makes predictions?
4. How do we use it to guide actions?

## THE STEPS

1. Get data about species occurrences
2. Build a classifier and make it as good as we can
3. Measure its performance
4. Explain some predictions
5. Generate counterfactual explanations
6. Briefly discuss ensemble models

**... think of SDM as a ML problem?** Because they are! We want to learn a predictive
algorithm from data

**... the focus on explainability?** We cannot ask people to *trust* - we must *convince*
and *explain*

§ 1
# Problem statement

We have information about a species

We have a series of observations $y \in \mathbb{B}$, and predictors variables $\mathbf{X} \in \mathbb{R}$

We want to find an algorithm $f(\mathbf{x}) = \hat{y}$ that results in the distance between $\hat{y}$ and $y$ being *small*

The predictor data will come from CHELSA2 - we will start with the 19 BioClim variables

We will use data on observations of *Turdus torquatus* in Switzerland, downloaded from the copy of the eBird dataset on GBIF

We want $\hat{y} \in \mathbb{B}$, and so far we are missing negative values

pseudo-absences

what are the assumptions we make

§ 2

# Training the model

$$P(+|x) = \frac{P(+)}{P(x)} P(x|+)$$

$$\hat{y} = \text{argmax}_j \, P(\mathbf{c}_j) \prod_i P(\mathbf{x}_i|\mathbf{c}_j)$$

$$P(x|+) = \text{pdf}(x, \mathcal{N}(\mu_+, \sigma_+))$$

Can we train the model

assumes parallel universes with slightly less data

is the model good?

coin flip

no skill

constant

| Model | MCC | PPV | NPV | DOR | Accuracy |
|---|---|---|---|---|---|
| noskill | -3.10619e-17 | 0.336873 | 0.663127 | 1.0 | 0.553221 |
| coinflip | -0.326255 | 0.336873 | 0.336873 | 0.25807 | 0.336873 |
| constantpositive | 0.0 | 0.336873 | NaN | NaN | 0.336873 |
| constantnegative | 0.0 | NaN | 0.663127 | NaN | 0.663127 |

k-fold

validation / training / testing

# WHAT TO DO IF THE MODEL IS TRAINABLE?

train it!

re-use the full dataset

# CAN WE IMPROVE ON THIS MODEL?

variable selection

data transformation

hyper-parameters tuning

will focus on the later (same process for the two above)

p plus > p minus means threshold is 0.5

is it?

how do we check this

SPATIALIZED PARTIAL RESPONSE (BINARY OUTCOME)

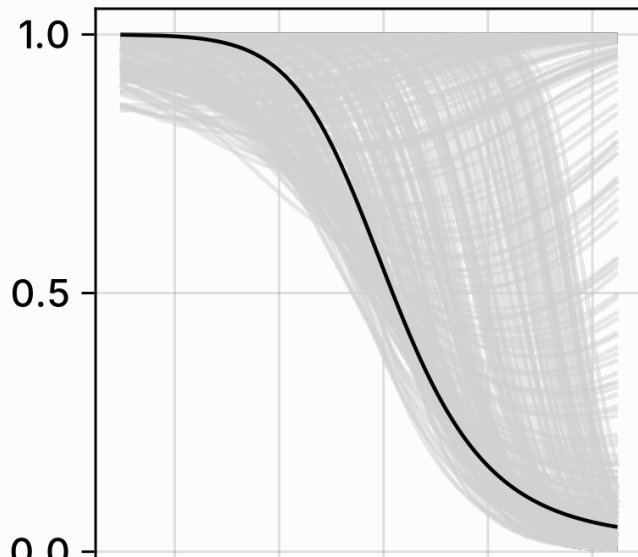Averaging the variables is masking a lot of variability!

Alternative solution:

1. Generate a grid for all the variables
2. For all combinations in this grid, use it as the stand-in for the variables to replace

In practice: Monte-Carlo on a reasonable number of samples.

- partial responses can only generate model-level information
- they break the structure of values for all predictors at the scale of a single observation
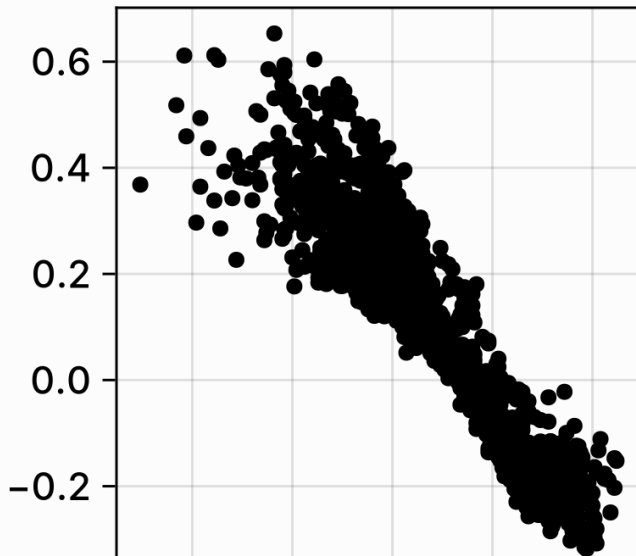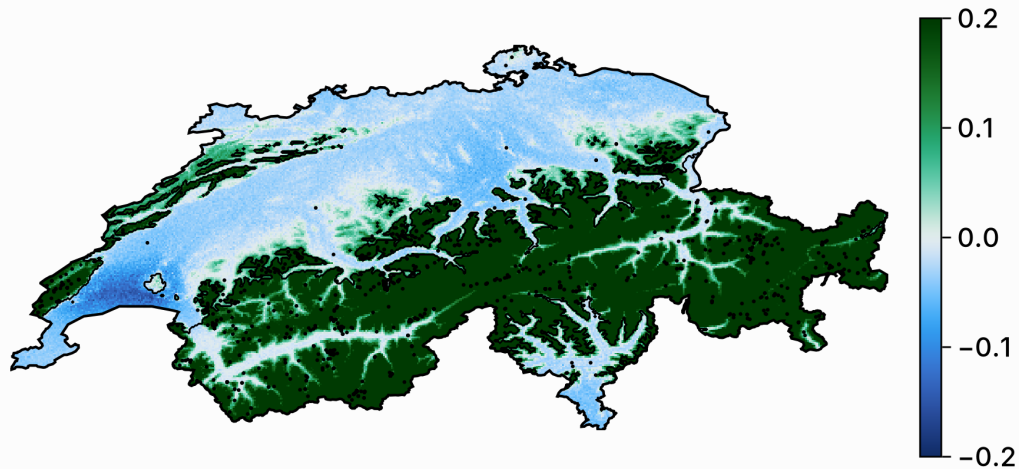- their interpretation is unclear

with shapley

mosaic map

what they are