



density estimation

Non Parametric Deepness Metric of Performance applied
to Modelling of Niches
Niche modeling

Camilo Oberndorfer Mejia¹
Miguel Valencia Ochoa²

Advisors:

Daniel Rojas Diaz³

Henry Laniado⁴

Research practice 1 & 2

Research proposal

Mathematical Engineering

School of Applied Sciences and Engineering

Universidad EAFIT

August 2022

¹Mathematical Engineering Student at EAFIT University coberndorm@eafit.edu.co

²Mathematical Engineering Student at EAFIT University(mvalenciao@eafit.edu.co

³Mathematical Engineering professor at EAFIT University, drojasd@eafit.edu.co

⁴Mathematical Engineering professor at EAFIT University, hlaniado@eafit.edu.co

1 Introduction (3 a 4 párrafos)

Species distribution models (SDMs) are numerical tools that estimate the relationship between observations of species occurrence records or abundance with environmental and/or spatial estimates of those sites (Elith & Leathwick (2009)). Often, niche models are either built with presence-only data (the registered sight of species) because absence data (places where the species was sought yet not observed) is unavailable (Merow & Silander Jr (2014)). Presence-only models include MaxEnt and maxLike, among others. Understanding species' spatial occurrence patterns and their dependence on environmental variables gives us insights to predict niche distributions across landscapes, sometimes requiring extrapolation in space and time. Even though the linkage between SDMs and theory is often weak, they are fundamental goals of ecology and evolution.

This research aims to create a non-parametric function that can potentially solve the problems presented by the niche models used nowadays. The environmental data of each pixel of a map can be represented with functional data, this is multivariate data that varies over a continuum (F. Heinrichs (2021)). Once the data of the spots of the map and the species occurrences is normalized, a specific deepness measure is applied to the data of species occurrences to create a replacement of a probability density function to color the corresponding map accordingly. The deepness metric used approximates a calculation by simplicial (Lopez-Pintado *et al.* (2014)). This method for niche modeling will then be compared to the already existing parametric methods Maxent and Max-like and checked for possible applications.

As of 2022, the most used niche modeling tool is Maxent and Maxlike. Both Maxent and Maxlike use underlying parametric assumptions, to generate the species distribution (Elith *et al.* (2010) & Merow & Silander (2014)). These assumptions create a bias that can cause incorrect predictions. Non-parametric methods offer solutions to these problems by not basing themselves on assumptions. That means species occurrence data will not be based on an assumed distribution (ibm (2021)).

2 Statement of the problem

2.1 Statement of the problem (4 a 5 párrafos)

Throughout the centuries humans have observed and recorded consistent relationships between species and the physical environment (M (2013)). For most regions, systematic biological survey data tend to be sparse and/or limited in coverage. Formal biological surveys in which a set of sites are surveyed and the presence/absence or abundance of species at each site are recorded tend to be sparse and costly. Yet, species records are available in the form of presence-only records in herbariums and museum databases. Many of these databases represent well over a century of public and private investment in biological science (Elith *et al.* (2010)) and are a hugely important source. Thus the more readily-available presence-only data and the desire to maximize their utility compels interest in Niche Modelling (Stockwell & Peterson (2002)).

A formal definition for a niche was given by G. Evelyn Hutchinson in 1950 who defined

it as a hypervolume. Hutchinson defined that for a series of environmental variables or covariates each variable had a range of values in which a species can survive. for example, let temperature equal x , let humidity equal y , and suppose a graph of x against y is plotted. Let x be a covariate, such as temperature, then there is a value of $x(x')$ above which a particular species could not survive. There is also another value $x(x'')$ below which the species could not survive. In the same manner, there are two identically defined points for $y(y', y'')$ (Vandermeer (1972)). The region defined by these points in 2-dimensional space is called a Niche. This definition can be expanded into as many dimensions as desired, with every point in this multidimensional space describing states of the environmental variables which are suitable for the survival of the species in question. It is this hypervolume that Hutchinson suggested would adequately conceptualize the fundamental niche.

Generally called Species Distribution Modeling (SDM) or Environmental Niche Modeling (ENM). We can now formally describe Niche Modeling as the estimation of the hypervolume given by a set of covariates suitable for a species within a landscape of interest L . The objective of niche modeling is to create a colored map showing the probability intensity of the presence of a species on every pixel of the map. The estimations done can give way to a better understanding of a species' correlation with its environment and allow policymakers to make informed decisions (Elith *et al.* (2010)).

The most widely available and generally used methods for the estimation of a species' spatial occurrence are MaxEnt and MaxLike. Both of these methods share an underlying distribution, the Gibbs distribution, and work by minimizing the distance of the estimation to the occurrence (Merow & Silander (2014)). The main difference being the first maximizes the relative entropy while the latter maximizes the likelihood. The main problem with both of these methods is the assumption of the distribution. Such assumptions can lead a model to be biased and to incorrect estimations or extrapolations, thus a non-parametric method is a possible solution to avoid these problems.

2.2 Formalization of the problem

In this paper, we will work with presence-only data, i.e., a set of locations within L where the species has been observed. Let Y be a binary variable, that expresses $Y = 0$ as absence and $Y = 1$ as presence, and let z denote a vector of environmental covariates in every location within L .

Define $f(z)$ to be the probability density of covariates across L , $f_1(z)$ to be the probability density of covariates across locations within L where the species is present. The quantity that we wish to estimate is $f_1(z)$ to then solve the probability of the presence of the species, conditioned on environment covariates: $Pr(Y = 1|z)$. Because of the Bayes rule we have:

$$f_1(z) = Pr(z|Y=1)$$

$$Pr(Y = 1|z) = \frac{f_1(z)}{f(z)} Pr(Y = 1) \quad (1)$$

Due to the lack of capability to assure the real distribution of $f_1(z)$ we will focus on the estimation of this distribution.

Es necesario hacer más énfasis en $f_1(z)$.
Lo que se quiere estimar es $Pr(Y=1|z)$, pero no se puede, se aplica la regla de Bayes y se deduce que $f_1(z)$ es clave para el objetivo del modelado.

3 Objectives

3.1 General objective

Propose a non-parametric niche modeling technique ~~through~~ ^{based on} depth measures ~~that allow~~ ^{to assign} assigning an appropriate ~~probability intensity~~ ^{records} to a map of environmental covariates from a set of ~~registers~~ ^{records} of the presence of a species. ^{suitability}

3.2 Specific objectives

- Develop a non-parametric method for niche modeling based on a simplicial like depth measure.
- Implement the proposed technique ~~in~~ ^{into} a Matlab toolbox ~~that~~ ^{to} create a colored map based on the probability intensity of a species occurring in each pixel of the map.
- Asses the performance of the proposed method with a diverse set of virtual species ~~results~~ including accuracy and computation time.

4 Justification

The traditional process for studying species' spatial occurrence patterns is a slow and expensive process that can be invasive to some ecosystems; this makes niche modeling a better alternative for conducting these studies. The traditional method consists of doing an excursion to the place of interest with biologists to set traps, kill animals present in the ecosystem, and revise them to establish the biodiversity of the area. For more information on traditional methods, Hoel (1943) describes the other traditional methods and their disadvantages. Niche modeling, in contrast, lets biologists check for the probability intensity of a species being in a certain area using the sampled data of similar places in terms of environmental characteristics. Niche modeling presents several practical advantages over traditional sampling methods which makes it an important tool for the efficiency of these studies and environment conservation.

Some implemented parametric models have some issues with a high sampling bias and computation time; for example, Max-like model results depicted in Merow & Silander (2014) show poor accuracy with a high computation time. The most explored flaw in parametric models is the sampling bias, this means the model is usually wrong when analyzing uncommon species which have a big importance in biodiversity studies. Non-parametric models can reduce this bias by not depending on parameter estimation but on information about all samples.

A non-parametric method also adds a different way to tackle the problem by centering more on robustness and simplicity, possibly giving a better solution to the species' spatial occurrence patterns recognition (Frost (2017)).

A niche modeling method has some possible benefits on price, facility, and time to entities that need or perform studies of species' spatial occurrence. The theoretic qualities of a

non-parametric method include a simpler algorithm, a low computation time, and a more robust metric for the map coloring function; These qualities improve the practicality of niche modeling methods as a tool for studying species distributions. These improvements together lead to a different way to tackle map coloring problems based on probability intensity and set the ground for better methods, or at least slight improvements for the future.

This research enriches the literature on niche models as it opens a new investigation path using a non-parametric estimation of probability intensity. Non-parametric estimation is an important step towards replacing old models as new parametric estimation-based models are usually worse than Maxent created in 2004. The used estimator for replacing simplicial calculation of deepness can be applied in similar problems with functional data where the final product is a probability intensity. These added benefits from this research open a whole lot of opportunities to improve the tools for map coloring problems, especially for niche modeling

5 Scope (2 o 3 párrafos)

The tools for the research will include Matlab toolboxes and functions, reference papers about niche modeling, and premade data in the form of virtual species and ASCII maps containing environmental information of different areas. The main Matlab tools for the research will be the Matlab shapes toolbox (mat (2022)) and a Matlab function developed by assessor Daniel Rojas Diaz and student Luisa Toro for virtual species generation. Also, Maxent software will be used to check the accuracy of the results of the map coloring function. These tools give the method an easy-to-access and easy-to-use algorithm for comparisons with future work.

Although the proposed method for the research can solve some of the issues of parametric estimation of $f_1(z)$ it can still share or even create problems of its own. For example, Its accuracy has a high dependence on the accuracy of the sampled data of environmental characteristics and species spatial occurrences as it uses this data to create the colored map. Its dependency on the sampled data also brings an issue on accurately predicting species occurrences in places with not very explored niches. Either way, Using sampled data from niches instead of sampled areas has the advantage of being a better estimator for habitats that share characteristics with the more explored areas. Little information is found about non-parametric niche modeling methods, making a comparison to a more similar model virtually impossible. As by Elith & Leathwick (2009) says, the factors that influence the robustness and accuracy of the model are the extents of the extrapolation, the interplay between environmental variables, and the consideration of scale. This makes the method have room for improvement for future research.

With the mentioned limits and tools, the research's expected results are the following: A map coloring function for probability intensity of species spatial occurrence in a Matlab toolbox for future applications and modifications of the method. The results contrast with the colored maps generated by Maxent software to check for the hypothesized accuracy and computation time. A paper detailing the process and tools ultimately used for the research and recommendations for future improvement. These deliverables aim to set a clear

Recomiendo cambiar el orden de presentación del scope:

- 1) ¿qué queremos hacer?
- 2) ¿qué dificultades prevemos?
- 3) ¿qué no haremos?

ground base for future non-parametric methods that attempt to create colored maps based on probability intensity for niche modeling or other similar subjects.

6 State of the art (5 a 6 párrafos)

→ Hace falta una exploración mas amplia de bibliografía.

The study of species' spatial occurrence is widely used as a decisive factor when establishing agricultural centers, roads, or constructions (Gaston & Fuller (2009)). Government and non-government organizations have adopted these methodologies for large-scale, real-world biodiversity mapping applications, like humming-birds distribution prediction for conservation planning (Tinoco *et al.* (2009)), environmental correlates of the European Wildcat (Monterroso *et al.* (2009)) or ants potential invasive expanding distribution (Ward (2006)).

The idea to explain a species' distribution with environmental/geographical features has long been studied, early work includes Murray (1866), by Murray (1866) and Schimper (1898), by Schimper (1898). Early quantitative approaches used multiple linear regression and linear discriminant function analyses to associate species and habitat (Capen (1981)). The application of generalized linear models (GLMs) allowed for non-normal error distributions, additive terms, and nonlinear fitted functions (Beery *et al.* (2021)).

The distribution of a species is typically represented as a map that indicates the spatial extent of the species. These representations can be classified into three categories (Beery *et al.* (2021)). Firstly, raw species observation data simply shows all the locations where a species is known to be present or absent. Secondly, statistical models, this category encapsulates SDMs and correlates observations with environmental characteristics. Thirdly, expert range maps are often based on a complex combination of heterogeneous information sources, including personal observations, understanding of the species' habitat preferences, local knowledge/reports, etc. The focus of this paper will be on the map coloring of statistical models.

→ lo del map coloring debe quedar muy claro este es el statement of the problem

In equation 1 we can see three terms, these are $f_1(z)$, $f(z)$ and $Pr(Y = 1)$. Some methods take $Pr(Y = 1)$ as a constant since it corresponds to the prevalence of the species, which is normally calculated by the proportion of occupied sites (Fernandez-Manjarrés (2018)). Naturally different ways of estimating $Pr(Y = 1)$ have arisen.

As previously mentioned, both MaxEnt and MaxLike are the state of the art regarding niche modeling. MaxEnt produces a suitability index. because of the way $f(z)$ and $f_1(z)$ are defined, MaxEnt assumes $f_1(z)$ can be found through the Gibbs distribution. $f_1(z) = f(z)e^{\eta(z)}$ where $\eta(z) = \alpha + \beta * h(z)$ and α is a normalizing constant that ensures $f_1(z)$ sums to 1 and β is a vector of coefficients applied to the different terms of the model. Finally $h(z)$ is the constrained covariates. This represents an ill-defined 'suitability index' because MaxEnt does not correspond to an explicit model of species occurrence.

→ Esto es más planteamiento del problema que estado del arte

7 Proposed methodology (5 o 6 párrafos)

Species spatial distribution pattern recognition problems are usually solved via three alternative processes. The traditional methods are very common in underdeveloped areas and are the

Este párrafo es más estado del arte

most invasive and expensive of the options. Neural networks trained via individual samplings of the species have a high computational cost and need a lot of samples to train properly. Lastly, there are probability intensity models which don't need as many samples nor do they invade ecosystems as traditional methods.

Probability intensity models share a data preprocessing procedure that will be considered for the research. This process starts by organizing data in a matrix comparing the values of environmental variables for all samples of the presence of a species. Then the data is normalized to prevent a scale bias for certain variables, the normalization selected is range normalization as it is simple to apply and has little computation time. It also transforms values from a big scale to values between zero and one that represents their position in the interval from its minimum sampled value to the maximum sampled value. Alam (2020) gives a complete explanation on range normalization. After this process, the model in itself begins.

The non-parametric method and the parametric methods differ mainly in the probability intensity calculations. Parametric methods try to find parameters from the normalized data for an underlying distribution. To replace probability intensity, the best alternative for non-parametric models is to use deepness measures. Simplicial is the most common method for calculating deepness but it has a high computation time with the expected amount of data. To replace simplicial calculations, we propose a different approach described ahead.

A Simplicial based calculation for depth can be greatly beneficial for a non-parametric model as it changes the computation time for better or worse and can make the model more simple or complex to apply. The alternative proposed is to first reduce dimensionality by applying Principal component analysis (PCA) to handle data more easily and then constructing a frontier that encloses all samples trying to reduce the space lacking samples. Abdi & Williams (2010) provides a good explanation of the PCA technique. After defining a frontier, the euclidean distance R_i between every sample I and the frontier is calculated and a sphere for that sample is constructed as a ball with a center on the sample and a R_i radius. After all the spheres are constructed. This concludes the process with the sample data leaving us with a sphere for every sample representing its influence.

For defining an equivalent measure for a probability intensity of a species being present in a certain area. The environmental values of every pixel of a map are located on the same plane as the spheres generated by the last process. The depth of this pixel is defined as the number of spheres it appears inside and how close it is to the middle of the spheres. This process makes areas with similar environmental values to those where the analyzed species have been located have a high chance of containing the species as well without using an underlying distribution or doing extensive excursions to the area to revise. This process is a good first step in making non-parametric models useful for niche modeling but it can be experimented with by using different distances of the samples or even using a different metric instead of depth for future research to make it a more practical tool in the future.

No es necesario este nivel de detalle para la técnica propuesta. Es mejor dejarlo como una introducción a la estimación de densidad mediante profundidad

Table 1: Schedule

Activity	Weeks																	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Define workflow and methodology																		
Literature check																		
Toolbox developing																		
Validation																		
Discussion																		

8 Schedule, commitments, and deliverables

Meetings will be held every week Monday 2:00 pm with both students and the Advisor Daniel Rojas Diaz and occasional meetings with advisor Henry Laniado. The expected deliverable for the research are a Matlab toolbox improvement and a paper detailing the results of it tested with multiple data.

9 Intellectual property

According to the internal regulation on intellectual property within Universidad EAFIT, the results of this research practice are product of *Miguel Valencia Ochoa, Camilo Oberndorfer Mejia* and *Daniel Rojas Diaz*.

In case further products, beside academic articles, that could be generated from this work, the intellectual property distribution related to them will be directed under the current regulation of this matter determined by Universidad EAFIT (2017).

References

- 2021 (Aug). *Statistics - parametric and nonparametric*.
2022. *Draw Shapes and Lines - MATLAB Simulink - MathWorks América Latina*.
- Abdi, Hervé, & Williams, Lynne J. 2010. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, **2**(4), 433–459.
- Alam, Rizwan. 2020. Normalization vs Standardization Explained. *Medium*, May.
- Beery, Sara, Cole, Elijah, Parker, Joseph, Perona, Pietro, & Winner, Kevin. 2021. Species Distribution Modeling for Machine Learning Practitioners: A Review. *Page 329–348 of: ACM SIGCAS Conference on Computing and Sustainable Societies*. COMPASS '21. New York, NY, USA: Association for Computing Machinery.

- Capen, David E. 1981. *The use of multivariate statistics in studies of Wildlife Habitat*. Rocky Mountain Forest and Range Experiment Station, Forest Service, U.S. Dept. of Agriculture.
- Elith, Jane, & Leathwick, John R. 2009. Species Distribution Models: Ecological Explanation and Prediction Across Space and Time. *Annual Review of Ecology, Evolution, and Systematics*, **40**(1), 677–697.
- Elith, Jane, Phillips, Steven J., Hastie, Trevor, Dudík, Miroslav, Chee, Yung En, & Yates, Colin J. 2010. A statistical explanation of Maxent for ecologists. *Diversity and Distributions*, **17**(1), 43–57.
- Fernandez-Manjarrés, Juan. 2018. Using Ecological Modelling Tools to Inform Policy Makers of Potential Changes in Crop Distribution: An Example with Cacao Crops in Latin America. *Economic Tools and Methods for the Analysis of Global Change Impacts on Agriculture and Food Security*, 11–23.
- F.Heinrichs. 2021. Functional Data Analysis(FDA). *towards science*, **1**(1), 1–3.
- Frost, Jim. 2017. Nonparametric tests vs. parametric tests. *Statistics By Jim*, Apr.
- Gaston, Kevin J, & Fuller, Richard A. 2009. The sizes of species’ geographic ranges. *Journal of applied ecology*, **46**(1), 1–9.
- Hoel, Paul G. 1943. The Accuracy of Sampling Methods in Ecology. *The Annals of Mathematical Statistics*, **14**(3), 289–300.
- Lopez-Pintado, Sara, Sun, Ying, Lin, Juan, & Genton, Marc. 2014. Simplicial band depth for multivariate functional data. *Advances in Data Analysis and Classification*, **8**(09), 321–338.
- M, Basvarajaiah D. 2013. *Biodiversity Modeling and Tribal Livelihood Status in Western Ghats*.
- Merow, Cory, & Silander, John A. 2014. A comparison of Maxlike and Maxent for modelling species distributions. *Methods in Ecology and Evolution*, **5**(3), 215–225.
- Merow, Cory, & Silander Jr, John A. 2014. A comparison of Maxlike and Maxent for modelling species distributions. *Methods in Ecology and Evolution*, **5**(3), 215–225.
- Monterroso, P., Brito, J. C., Ferreras, P., & Alves, P. C. 2009. Spatial ecology of the European wildcat in a Mediterranean ecosystem: dealing with small radio-tracking datasets in species conservation. *Journal of Zoology*, **279**(1), 27–35.
- Murray, Andrew. 1866. *The geographical distribution of mammals*. Day and Son.
- Schimper, Andreas Franz Wilhelm. 1898. *Pflanzen-geographie auf physiologischer Grundlage*. Vol. 2. G. Fischer.

- Stockwell, David R.B, & Peterson, A.Townsend. 2002. Effects of sample size on accuracy of species distribution models. *Ecological Modelling*, **148**(1), 1–13.
- Tinoco, Boris A., Astudillo, Pedro X., Latta, Steven C., & Graham, Catherine H. 2009. Distribution, ecology and conservation of an endangered Andean hummingbird: the Violet-throated Metaltail (*Metallura baroni*). *Bird Conservation International*, **19**(1), 63–76.
- Universidad EAFIT. 2017. *Reglamento de propiedad intelectual*.
- Vandermeer, John H. 1972. Niche Theory. *Annual Review of Ecology and Systematics*, **3**, 107–132.
- Ward, Darren F. 2006. Modelling the potential geographic distribution of invasive ant species in New Zealand. *Biological Invasions*, **9**(6), 723–735.