

# Optimising Species Distribution Models: Sample size, positional error, and sampling bias matter

Vítězslav Moudrý<sup>1</sup>, Manuele Bazzichetto<sup>1</sup>, Ruben Remelgado<sup>2</sup>, Rodolphe Devillers<sup>3</sup>, Jonathan Lenoir<sup>4</sup>, Rubén G. Mateo<sup>5</sup>, Jonas J. Lembrechts<sup>6</sup>, Neftalí Sillero<sup>7</sup>, Vincent Lecours<sup>8</sup>, Anna F. Cord<sup>2</sup>, Vojtěch Barták<sup>1</sup>, Petr Balej<sup>1</sup>, Duccio Rochini<sup>1,9</sup>, Michele Torresani<sup>10</sup>, Salvador Arenas-Castro<sup>11</sup>, Matěj Man<sup>12</sup>, Dominika Prajzlerová<sup>1</sup>, Kateřina Gdulová<sup>1</sup>, Jiří Prošek<sup>1,12</sup>, Elisa Marchetto<sup>9</sup>, Alejandra Zarzo-Arias<sup>13,14</sup>, Lukáš Gábor<sup>1</sup>, François Leroy<sup>1</sup>, Matilde Martini<sup>9</sup>, Marco Malavasi<sup>15</sup>, Roberto Cazzolla Gatti<sup>9</sup>, Jan Wild<sup>1,12</sup>, Petra Šímová<sup>1</sup>

<sup>1</sup>Department of Spatial Sciences, Faculty of Environmental Sciences, Czech University of Life Sciences Prague, Kamýcká 129, 16500 Praha-Suchdol, Czech Republic.

<sup>2</sup>Chair of Computational Landscape Ecology, Dresden University of Technology, Helmholtzstr. 10, 10169 Dresden, Germany.

<sup>3</sup>UMR Espace-Dev, Institut de Recherche pour le Développement, Univ Réunion, La Réunion, France.

<sup>4</sup>UMR CNRS 7058 “Ecologie et Dynamique des Systèmes Anthropisés” (EDYSAN), Université de Picardie Jules Verne, 80000 Amiens, France.

<sup>5</sup>Departamento de Biología and Centro de Investigación en Biodiversidad y Cambio Global (CIBC-UAM), Universidad Autónoma de Madrid, Madrid, Spain.

<sup>6</sup>Research Group of Plants and Ecosystems (PLECO), Department of Biology, University of Antwerp, Antwerp, Belgium.

<sup>7</sup>Centro de Investigação em Ciências Geo-Espaciais (CICGE), Faculdade de Ciências da Universidade do Porto, Alameda do Monte da Virgem, 4430-146 Vila Nova de Gaia, Portugal.

<sup>8</sup>Université du Québec à Chicoutimi, 555 Boul. de l'Université, Saguenay, Canada.

<sup>9</sup>BIOME Lab, Department of Biological, Geological and Environmental Sciences, Alma Mater Studiorum University of Bologna, via Irnerio 42, 40126, Bologna, Italy.

<sup>10</sup>Free University of Bolzano/Bozen, Faculty of Agricultural, Environmental and Food Sciences, Piazza Università / Universitätsplatz 1, 39100, Bolzano/Bozen, Italy.

<sup>11</sup>Dept. of Botany, Ecology and Plant Physiology, Faculty of Sciences, University of Cordoba, Spain.

<sup>12</sup>Institute of Botany of the Czech Academy of Sciences, Zámek 1, 252 43 Průhonice, Czech Republic.

<sup>13</sup>Universidad de Oviedo, Oviedo, Spain.

<sup>14</sup>Department of Biogeography and Global Change, Museo Nacional de Ciencias Naturales (MNCN-CSIC), Madrid, Spain.

<sup>15</sup>Department of Chemistry, Physics, Mathematics and Natural Sciences, University of Sassari, Via Vienna 2, Sassari 07100, Italy

November 30, 2023

## Abstract

Species distribution models (SDMs) have proven valuable in filling gaps in our knowledge of species occurrences. However, despite their broad applicability, SDMs exhibit critical shortcomings due to limitations in species occurrence data. These limitations include, in particular, issues related to sample size, positional error, and sampling bias. In addition, it is widely recognized that the quality of SDMs as well as the approaches used to mitigate the impact of the aforementioned data limitations are dependent on species ecology. While numerous studies have experimentally evaluated the effects of these data limitations on SDM performance, a synthesis of their results is lacking. However, without a comprehensive understanding of their individual and combined effects, our ability to predict the influence of these issues on the quality of modelled species-environment associations remains largely uncertain, limiting the value of model outputs. In this paper, we review studies that have evaluated the effects of sample size, positional error, sampling bias, and species ecology on SDMs outputs. We integrate their findings into a step-by-step guide for critical assessment of spatial data intended for use in SDMs.

**Keywords:** Complexity, Ecological Niche Modelling, Filtering, Heterogeneity, Sampling, Scale, Training, Quality, Validation

## 1. INTRODUCTION

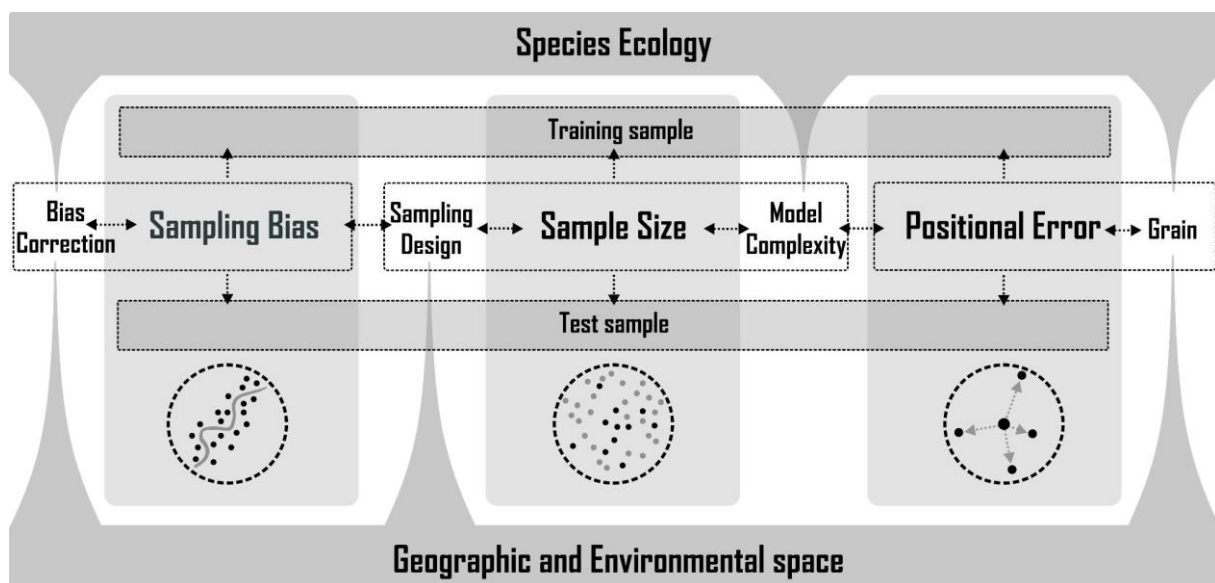
The quantity and quality of biological observations has improved dramatically over the past decades. However, a certain level of uncertainties is inherently present in such data, resulting in uncertainties of scientific inferences based on it (Hortal et al. 2015; Daru and Rodriguez 2023; Hughes et al. 2023). Correlative species distribution models (SDMs; aka habitat suitability models or ecological niche models; Sillero 2011) are useful for tackling the gaps in our knowledge of species occurrence (Elith and Leathwick 2009). These models combine environmental and species occurrence data to build a set of rules describing the environmental space (i.e. species ecological niche) where species were observed and, subsequently, can be used to predict the distribution of the respective species (Ferrier et al. 2017). SDMs support a wide variety of ecological applications, such as the assessment of the spread of invasive species (Guisan et al. 2013; Bazzichetto et al. 2021), the detection of potential impacts of environmental changes on biodiversity (Ehrlén and Morris 2015; Haesen et al. 2023), or the identification of suitable locations for the relocation of endangered species (Guisan et al. 2013; Segal et al. 2021). However, despite their broad applicability, SDMs have critical shortcomings associated in particular with the characteristics of input data, including their quantity and quality (Elith et al. 2002; Barry and Elith 2006; Rocchini et al. 2011; Moudrý and Šímová 2012; Wüest et al. 2020). In this paper, we focus on limitations of species occurrence data (for issues associated with environmental data, see e.g. Fourcade et al. 2018; Araújo et al. 2019; Moudrý et al. 2023).

Limitations of species occurrence data can introduce uncertainty and biases in the estimation of species-environment relationships and, consequently, of their modelled distributions (Araújo et al. 2019). In particular, data availability (i.e. sample size) is critical; the smaller is the minimum sample size that can theoretically be used in SDMs, the higher is the number of species that can be modelled (e.g. Stockwell and Peterson 2002). However, measurement errors associated with data acquisition methods (i.e. positional error; Smith et al. 2023) are another major source of uncertainty, which may, in effect, necessitate the use of a larger sample size than had the data been accurate. In addition, the choice of inappropriate sampling strategies (if any) can introduce biases toward certain locations (i.e. sampling bias; Bazzichetto et al. 2023). Moreover, it is well recognized that the quality of SDMs is also

influenced by the species' ecology (Segurado and Araújo 2004; Heikkinen et al. 2006; McPherson and Jetz 2007; Guisan et al. 2007; Collart et al. 2023) and the fact that the effects of different data limitations (e.g. sample size, positional error, and sampling bias) may be species-specific.

As the interest in using SDMs continues to grow, tackling data limitations becomes increasingly critical (Araújo et al. 2019; Wüest et al. 2020; Jansen et al. 2022; Marcer et al. 2022). In this context, data characteristics and limitations are expected to be regularly considered and properly reported during the conceptualization and validation of SDMs (Feng et al. 2019; Zurell et al. 2020; Sillero and Barbosa 2021; Tassarolo et al. 2021; Jansen et al. 2022; Boyd et al. 2023). However, without proper knowledge of the individual or combined effects of sample size, positional error, sampling bias, and species' ecology, our ability to anticipate the effect of these issues on the quality of modelled species-environment associations remains largely uncertain, limiting the value of model outputs (see **Figure 1** for a diagram introducing data characteristics and their relationships considered in this perspective).

A common approach to the evaluation of the effects of data limitations on model performances is to manipulate the input data experimentally or to simulate datasets impacted by various sources of error. Here, we examine studies that manipulated the sample size (**Section 2**), introduced positional errors (**Section 3**) or sampling bias (**Section 4**), to investigate their impact on SDMs' outputs. Accordingly, we also provide guidance for a critical assessment of spatial data to be used in SDMs and identify future directions towards the development of guidelines for optimising the tradeoffs between data limitations and accurate modelling of species-environment relationships (**Section 5**).



**Figure 1.** Sample size, positional error, and sampling bias are the three essential characteristics of species occurrence data addressed in this study. These interconnected characteristics can have a significant impact on the reliability of SDM results. Researchers must thoughtfully address these factors during the collection of species occurrence data (sampling design) and the construction of models (model complexity). Maximising sample size, using sampling bias correction methods, and minimising positional error relative to the study's spatial resolution during model training and testing are all essential steps. Additionally, species ecology and the distribution of species observations in the geographic and environmental space can exacerbate/attenuate the negative effects of small sample size, high sampling bias, and high positional error on the reliability of SDMs results. See the glossary box for the explanation of terms.

## 2. SAMPLE SIZE

Among all possible factors, sample size (see Glossary) has the most profound effect on the performance of an SDM (Thibaud et al. 2014; Santini et al. 2021). Sample size poses an important constraint to the model complexity, i.e. to the number of predictors to be estimated as well as of potential algorithms and their parameters used for modelling. Sample sizes in SDMs can range from just a few (Papeş and Gaubert 2007; Pearson et al. 2007) to millions (e.g., Botella et al. 2023; Gábor et al. 2023a) of records. In the vast literature measuring the effect of sample size on model performance (see **Table 1**), the primary concern has been to determine the minimum adequate sample size required to produce reliable and fit-for-purpose models (e.g., Stockwell and Peterson 2002; Hanberry et al. 2012; Proosdij et al. 2016). In parallel, ecological research investigates to what extent additional time and economic resources should be spent to improve models by increasing the sample size (e.g., Liu et al. 2019). Knowing the minimum (and maximum) sample size required for accurate predictions would theoretically allow optimisation of the resources spent on labour-intensive fieldwork and, therefore, help reduce associated costs. Nonetheless, the extent to which modelling could replace fieldwork remains questionable.

### 2.1. Importance of sample size in model training and testing

Studies focusing on a better understanding of how the sample size impacts the models' accuracy revealed that it is in principle possible to train SDMs with a relatively small sample. Values typically range from 50 to 150 presences (or presences-absences), although values as low as 10 presences or as high as a few hundred have also been reported (**Table 1**). However, it is important to note that studies typically reported minimum values when the model was still relatively useful, not values when the model gave optimal results. Besides, it has been reported that models relying on fewer than approximately 70 presences do not reliably identify the variables affecting distributional patterns (Smith and Santos 2020) or result in poor(er) estimates of the shapes of species response curves (Coudun and Gégout 2006; Shiroyama et al. 2020; Bazzichetto et al. 2023; Wang and Jackson 2023). In general, all studies agreed that increasing sample size increased a model's predictive performance (keeping the number of predictors fixed), although a plateau in model performance is generally reached (Stockwell and Peterson 2002). According to recent studies, hundreds of presences are needed to reach the plateau where increasing sample size further adds little to the model performance (Liu et al. 2019; Gábor et al. 2020a).

Insufficient attention has so far been devoted to the evaluation of possible effects of the testing dataset sample size on validating SDMs' predictive performances. Generally, small validation datasets can lead to inaccurate assessment of model performance (Hallman and Robinson 2020). Recently, Jiménez-Valverde (2020) showed that 30 presence-absence records (i.e., 15 presences and 15 absences) are a (minimum) adequate sample size for a validation dataset to estimate the predictive performance of presence-absence models. Nevertheless, their conclusions are based on simulations, and it should be, therefore, pointed out that studies using real data are essential to generalise these results. In addition, the minimal sample size of a validation dataset has not yet been evaluated in the case of presence-background data; since these carry less information than presence-absence data, the validation set should be correspondingly larger (Collart and Guisan 2023). While the importance of a

sufficient validation sample is intuitive, the impact of validation dataset sample size on model performance and validation accuracy urgently needs to be further tested.

**Table 1.** Overview of studies testing the role of the number of presences or presences and absences for model performance. PA - presences-absences.

Study	Number of species	Training sample	Testing sample	Study extent / resolution	No. predictors	No. obs. suggested
Stockwell and Peterson (2002)	130 birds	1 - 100	1000; presence-background	Mexico / 3 × 3 minutes	8	at least 50 presences
Kadmon et al. (2003)	192 plants	10 - 200	96 plots; presence-absence	Israel / 1 × 1 km	3	50 - 75 presences
Hernandez et al. (2006)	18 animals	5 - 100	50 presences	California / 1 × 1 km	10	50 - 75 presences
Wisn et al. (2008)	46 plants, animals	10 - 100	presence-absence data	five regions / 100 × 100 m; 1 × 1 km	11-13	at least 30 presences
Mateo et al. (2010)	2 plants	9 - 60	compared to maps created with full datasets	Ecuador / 1 × 1 km	19	at least 20 presences
Feeley and Silman (2011)	65 plants	25 - 150	compared to maps created with full datasets	tropical South America / 5 × 5 km	3	Larger than evaluated
Hanberry et al. (2012)	16 trees	30 - 2500	Presence samples not used for training	46,000 km <sup>2</sup> / 310,000 polygons	16	at least 200 presences
Proosdij et al. (2016)	6 virtual	3 - 50	Compared with actual virtual species distribution	18,000,000 km <sup>2</sup> / 5 × 5 minutes	15	14-25 presences
Liu et al. (2019)	1800 virtual	20 - 640	3000 presences and absences of virtual species distribution	62,500 km <sup>2</sup> / 1 × 1 km	6	a few hundred presences
Støa et al. (2019)	30 insects	5 - 320	Compared to maps created with full datasets	Norway / 1 × 1 km	2	10-15 presences
Smith and Santos (2020)	1 virtual	8 - 1024	400 presences and absences of virtual species distribution	Virtual landscape / 1024 × 1024 cells	1	at least 128 presences
McPherson et al. (2004)	7 birds	50 - 500	500 presences-absences	South Africa / 0.25 × 0.25 degrees	61	300 PA
Coudun and Gégout (2006)	54 virtual	50 - 5000	Not used	Not relevant	1	at least 50 PA
Jiménez-Valverde et al. (2009)	1 virtual	182 - 182,288	Compared with actual virtual species distribution	6,576 km <sup>2</sup> / 0.04 × 0.04 degrees	4	at least 70 PA
Shiroyama et al. (2020)	Bluegill	50 - 900	110 presences absences	Seven rivers in Kanto region, Japan.	4	at least 400 PA
Bazzichetto et al. (2023)	2 virtual	200 - 500	Compared with actual virtual species distribution	10 794 km <sup>2</sup> / 1 × 1 km	2	at least 200 PA
Wang and Jackson (2023)	16 virtual	50 - 800	50 presences absences	140 000 km <sup>2</sup> / 4 × 4 km	2	at least 100 PA

## 2.2. On the relationships between sample size, species ecology and model complexity

The association between model performance and sample size depends largely on the species' ecology. Studies have repeatedly indicated that for a given sample size, SDMs better predict species with restricted geographical distributions (i.e., low range size, prevalence, or relative occurrence area), as well as specialist species with strict ecological requirements (i.e., narrow ecological niche) than species with wide geographic ranges and generalist (i.e., wider ecological niche) species (Stockwell and Peterson 2002; Seoane et al., 2005; Hernandez et al. 2006; Tsoar et al. 2007; Mateo et al. 2010; Tessarolo et al. 2014; Proosdij et al. 2016; Hallman and Robinson 2020; Arenas-Castro et al. 2022; Wang and Jackson 2023). The association between model performance, sample size, and species ecology can be explained by niche completeness (i.e., the proportion of a species' niche covered by the sampling). For example, if a species has a restricted ecological niche (or range), the niche may likely be well represented by a low number of occurrences. On the other hand, large sample size does not necessarily mean a complete coverage of the entire ecological niche for widespread species (Bazzichetto et al. 2023; Boyd et al. 2023). This is further related to model complexity. The minimum required sample size increases with the number of variables or parameters and the complexity of the assumed species response curves (e.g., quadratic response curves or statistical interactions among predictors; Austin 2002; Barry and Elith 2006; Maggini et al. 2006; Ficetola et al. 2014; Merow et al. 2014; Bell and Schlaepfer 2016; Carretero and Sillero, 2016). However, even large sample sizes can result in low accuracy in the estimation of model parameters if the model is overly complex (i.e., includes too many parameters or interactions, e.g., Wisz et al. 2008; Moreno-Amat et al. 2015).

## 2.3. Recommendations associated with sample size

The above-mentioned studies showed that SDMs can perform relatively well even with small sample sizes (**Table 1**). However, the studies mentioned in Table 1 are difficult to directly compare due to the use of different species, differences in the used modelling algorithms, numbers of parameters, or spatial resolutions and geographical extents. Whether the sample size is considered small or sufficient depends largely on the number of predictors in the model and the complexity and nature of the species–environment relationships (e.g., Merow et al. 2014; Smith and Santos 2020; Bazzichetto et al. 2023). Hence, given how context-dependent these relationships are, we cannot recommend a specific threshold of what a 'small' or 'large' sample is but we provide a series of steps that researchers should consider when preparing SDMs:

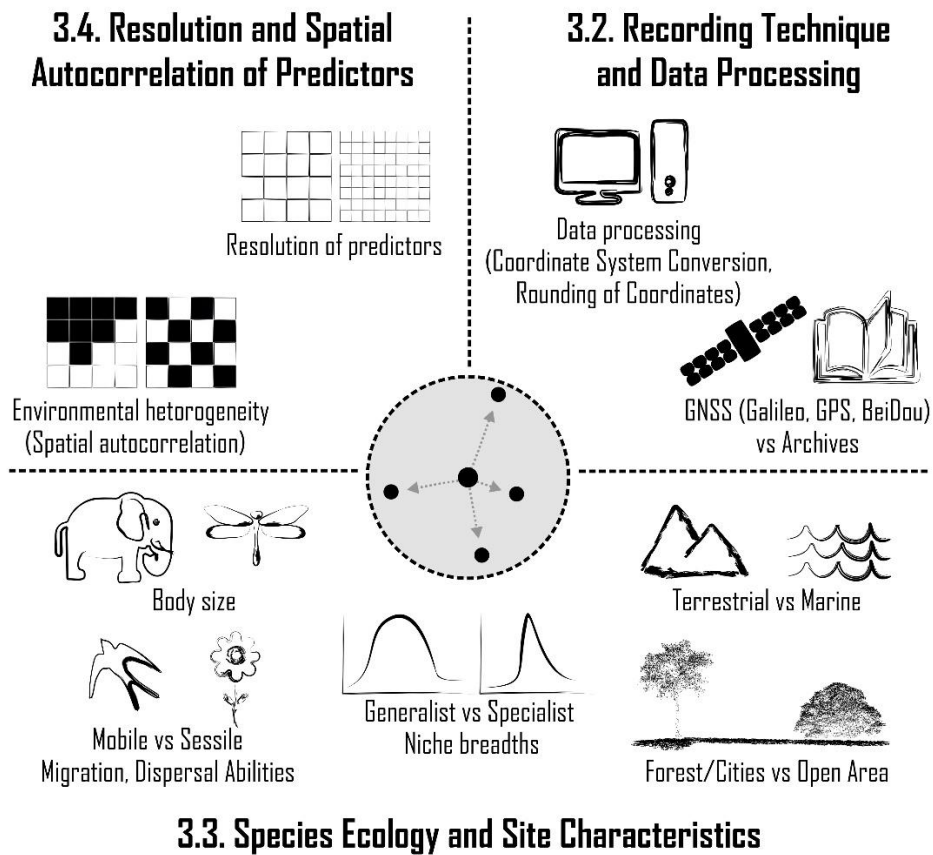
- First of all, the sample size required for a particular analysis requires careful consideration of the purpose of the study (Foody, 2011). On the one hand, models based on low sample sizes can help identify potential knowledge gaps and optimise the allocation of funds for field surveys (e.g. to pinpoint areas with a high potential for discovering unknown populations of the studied species; Raxworthy et al. 2003; Fois et al. 2018). On the other hand, it is crucial to emphasise that predictions derived from models with small sample sizes should not be employed as guidelines for applications such as modelling species ranges, predicting responses to climate change, or planning conservation efforts (Loiselle et al. 2008; Feeley and Silman 2011; Duputié et al. 2014).

- Second, species' ecology has to be considered as SDMs better predict specialist species with narrow ecological niches than generalist species with wider ecological niches (e.g. Tsoar et al. 2007).
- Third, researchers should consider the number of predictors investigated. As the ability to differentiate between influential and uninfluential variables decreases with decreasing sample sizes, the challenge lies in the identification of variables that genuinely influence species distribution (Smith and Santos 2020). Studies that include a small number of variables selected based on expert opinion will generally require a smaller sample size than studies that select variables from a large pool using automated algorithms (Ficetola et al. 2014).
- Fourth, the complexity of the shape of species response curves must be taken into account as models based on small sample sizes result in less accurate estimates of these shapes (Bazzichetto et al. 2023). Models aiming at generating simple response curves (e.g., linear, hinge or step) can be developed with relatively low sample sizes. However, models identifying more complex shapes such as logistic, gaussian or even non-parametric smooth functions of variable flexibility require much larger sample sizes. Adding interactions between variables increases the requirements for the sample size even more.
- We cannot suggest a minimum number of presences (presences-absences) as a rule of thumb but if you are unsure whether your sample size is sufficient given the objectives and complexity of your model, we recommend testing the effect of sample size in your dataset. Start with the most comprehensive model you think is appropriate in your particular case and progressively increase the sample size until you reach your possible maximum (i.e., all presences you have), and see if your model performance is reaching a plateau. If no plateau is reached, it is likely that more presences are necessary. In such a case, a reduction in the number of variables or the complexity of the response curves should be considered. Remember to set aside at least 30 presence-absence records for model validation, as recommended by Jiménez-Valverde (2020).
- Finally, while it is possible to design accurate SDMs with a well-balanced sampling of as few as 50 presences (**Table 1**), most observational data are too *ad-hoc* and far from being representative of spatial variations in species-environment associations due to confounding effects of data limitations such as positional errors (**Section 3**), or sampling bias (**Section 4**). Hence, researchers should also consider these data limitations before attempting to build a model based on a small sample.

### 3. POSITIONAL ERROR

Species occurrence data are always prone to positional error, i.e. the difference between the actual and recorded location of a species in the coordinate reference system of the dataset. The magnitude of the positional error associated with species observations can range from a few centimetres up to tens of kilometres. Under high positional error, SDMs using environmental layers at spatial resolutions finer than the magnitude of the positional error (e.g., environmental layers at a 10 m resolution and a 50 m positional uncertainty of species observations) can estimate erroneous/misleading species-environment relationships. The potential effect of positional error on SDMs performance is determined by several interacting factors (**Figure 2**). Therefore, positional error should be assessed

before calibrating and validating SDMs, as it can negatively affect training and testing datasets as well as modelling decisions, such as the spatial resolution of environmental variables.



**Figure 2.** Three groups of interacting factors that determine the magnitude and potential impact of positional error on SDMs performance can be specified: the recording technique and data processing (**Section 3.2**); species ecology and characteristics of the site (**Section 3.3**); and the spatial resolution and degree of spatial autocorrelation of the predictors (**Section 3.4**).

### 3.1. How to address positional error in training and testing datasets

Several studies have examined the impact of positional error on SDMs performance by simulating shifts in species presences (**Table 2**). These studies typically compare SDMs outcomes based on data with high positional accuracy against results obtained using the same data but affected by positional error of different magnitudes. Findings from these studies have been somewhat mixed: some found little effect of positional error and reported that SDMs were relatively robust to it (Graham et al. 2008; Fernandez et al. 2009); others concluded that species occurrence data with positional error generally lead to less accurate SDMs (Johnson and Gillingham 2008; Osborne and Leitão 2009; Mitchell et al. 2017).

In real-world applications, a mix of high- and low-accuracy distribution data is the most common situation and analysts usually have to find a compromise between positional error and sample size (Smith et al. 2023). However, studies focusing on this issue yielded somewhat conflicting results.



Reside et al. (2011) warned that increasing the sample size by incorporating historic species occurrence data with inaccurate positions can reduce SDMs performances. On the other hand, Smith et al. (2023) showed that the removal of data with high positional error can excessively reduce the sample size and, thus, the model accuracy (Smith et al. 2023). Furthermore, Gábor et al. (2023b) showed that even models affected by positional error in species distribution data can be ecologically interpretable. Another study investigating the effect of positional error concluded that models with small sample sizes were more affected by positional error than models based on larger sample sizes (Mitchell et al. 2017).

The role of positional error is seldom accounted for during SDMs evaluation/testing. Surprisingly, most SDMs studies dealing with positional error only focused on the training dataset, while ignoring the (potential) effect of inaccurately georeferenced data in the validation dataset. The ultimate consequence of positional error in species data lies in an erroneous identification of the presence or absence in a given cell. In this regard, Foody (2011) demonstrated that validation data should be error-free (i.e., correctly distinguish between presences and absences) as even a small amount of error could result in misidentification of presences/absences and substantial misestimation of model performance. Therefore, data correctly labelled as the presence or absence of species (i.e., with minimal positional error) is essential for assessing model performance. More recently, Moudrý et al. (2017) showed that the inclusion of potentially erroneous presences (in this case ambiguous breeding bird categories used in the breeding bird atlases, i.e. possible and probable breeding) severely affected models' performance metrics when added to the validation dataset, while it had a relatively minor effect on model performance when added to the training dataset. Therefore, we suggest relying on large sample size possibly including observations with low positional accuracy (i.e. with higher positional error than the spatial resolution of predictors) for model calibration, while preserving high-accuracy data for model validation.

Alternatively, Moudrý and Šimová (2012) suggested that knowing the positional error of the occurrences allows balancing high and poor-quality data in both training and testing datasets and including a covariate predictor in the model (even as a factor variable with a few levels) to be tested or to up/down weight the importance of observations (see Velásquez-Tibatá et al. 2016 for such approach using Bayesian models). This allows preserving most of the data and offsetting the potential negative effect of high positional error. On the other hand, if the covariate has many levels and few observations, it might be better to subselect the data to retain only those of the best quality. If only a small sample size is available, we recommend considering the use of methods to mitigate positional error (Hefley et al. 2014; Zhang et al. 2018; Smith et al. 2023). Note, however, that the existing approaches typically either require knowledge of the magnitude of the error and their use is limited to data with relatively small positional error (Zhang et al. 2018), or they require that at least part of the dataset is recorded with minimal positional error (Hefley et al. 2014; Smith et al. 2023). Although recent literature is favouring the inclusion of observations with reasonable positional error rather than reducing sample size, we recommend careful consideration of their trade-off. Whether it is preferable to maintain the sample size or to minimize the adverse effect of positional error remains a very current and unanswered question.

**Table 2.** Studies analysing the influence of positional error in species occurrence data on SDMs.

Study	Species data	Resolution of environmental var.	Range of shifting occurrences	
			Distance	Cells
Graham et al. (2008)	Observed	100 × 100 m	0 - 5 km	0 - 50 cells
Johnson and Gillingham (2008)	Observed	30 × 30 m	50 - 1000 m	1 - 34 cells
Osborne and Leitão (2009)	Observed	1 × 1 km	0 - 5 km	0 - 5 cells
Fernandez et al. (2009)	Observed	1 × 1 km	5 - 50 km	1 - 50 cells
Naimi et al. (2011)	Virtual	artificial data	Not valid	1 – 30 cells
Mitchell et. al. (2017)	Observed	2.5 × 2.5 m	5 - 400 m	1 - 160 cells
Velásquez-Tibatá et al. (2016)	Virtual	150 × 150 cells	Not valid	5 - 15 cells
Gábor et al. (2020b)	Virtual	5 × 5 m	5 – 500 m	1 – 100 cells
Gábor et al. (2023b)	Virtual	50 × 50 m	50 - 1500 m	1 - 30 cells
Gábor et al. (2023b)	Observed	200 × 200 m	1 - 30 km	1 - 30 cells

### 3.2. Role of recording technique and data processing

Old datasets, such as historical observations archived in museums, atlases and natural history collections that were retrospectively georeferenced, are usually thought to be more prone to relatively higher positional error than new ones (Graham et al. 2004; Wiczorek et al. 2004; Newbold 2010; Bloom et al. 2018, Marcer et al. 2022). However, positional error affects any dataset, including those georeferenced using modern technologies such as the global navigation satellite systems (GNSS). Indeed, several factors can degrade GNSS positional accuracy, including the number and position of satellites, and the characteristics of the study site (e.g. beneath a dense forest canopy vs. an open grassland). The use of a low number of satellites to georeference species data may be due to the use of outdated technology, such as the use of a device that relies only on the United States' Global Positioning System (GPS), instead of using all currently available systems (e.g. Galileo, Glonass, and Beidou). Even when the above-mentioned challenges are overcome, species occurrence data may still be impacted by errors introduced during data processing (e.g. wrong transformations among coordinate reference systems, rounding of coordinates, or lack of error correction procedures (e.g. post-differential correction; Sillero and Seco 2014). Unfortunately, the positional error of species records is often undocumented (Moudrý and Devillers 2020; Marcer et al. 2022).

### 3.3. Relationships between positional error, species ecology and ecosystem characteristics

It is usually impossible to accurately georeference positions for non-sessile species (unless they are equipped with transmitters) due to environmental barriers (for example, it's impossible to get close to the species in some habitats) and/or species characteristics (e.g. size, mobility and behaviour) (Frair et al. 2010). Besides, georeferencing species' location using GNSS in a dense forest or at the bottom of a narrow and deep ravine may be difficult due to the poor reception of the satellite signal. In addition, buildings, walls and trees in the proximity of an antenna can reflect the signal from satellites, thereby further reducing the positioning accuracy (a phenomenon known as multipath; Kos et al. 2010). Besides, GNSS does not work underwater; in effect, the positioning of species observations in marine and freshwater environments is based on acoustic positioning, which leads to a decrease in accuracy with the water column depth, or simply on recording a position at the surface of water and disregarding movements of the sampling gear in the water column (Rattray et al. 2014, Mitchell et al. 2017). As a result, data for mobile animals can have a positional uncertainty of tens to hundreds of metres. The distance between an animal and the observer is positively associated with the species' body size and, therefore, big animals are typically less accurately georeferenced as they move a lot or can be even dangerous, which leads to recording their location from a distance (Zhang et al. 2018).

The effect of positional error on SDMs may also depend on the species mobility, expressed as the daily dispersal range or migration ability. Many birds, fishes and big predators are very mobile, and the accurate georeferencing of their location may play a smaller role in SDM performance than in the case of sessile species (see **Figure 2** for an overview of the factors that may interact with the magnitude of positional error when building SDMs). In this regard, Gábor et al. (2023b) showed that the performance of a band-tailed pigeon SDM only slightly decreased with an increasing positional error, while virtual species simulations that did not consider species mobility showed a rapid decrease. Although positional error seems to depend on species characteristics, its role in affecting SDMs for different groups (such as insects vs. big mammals; mobile organisms like birds vs. sessile organisms like plants, corals, etc.) is understudied. Among the few studies that analysed the interaction between positional error and species ecology, Velásquez-Tibatá et al. (2016) and, more recently, Gábor et al. (2020b), showed that positional error has a greater impact on SDMs' performances for specialists (i.e., species with a narrow niche breadth) than for generalist species (i.e., those with a wide niche breadth). This is due to occurrences of specialist species being more susceptible to a shift into unsuitable environments.

### 3.4. Relationship between positional error, spatial resolution and autocorrelation

The spatial resolution of predictors used in SDMs is another critical factor determining the impact of positional error on model performance. Previous studies on positional error considered shifts from 5 m up to 50 km. Such a range of errors results in a less impactful shift of species data over raster cells (and across environmental conditions) in a coarse-resolution set of environmental layers (e.g., 1 × 1 km) than in a fine-resolution set of environmental layers (e.g., 10 × 10 m). Note that the recent studies investigated shifts of the species occurrence data by up to 160 pixels (which is almost threefold compared to older studies) thanks to the reduced pixel sizes in the current environmental data (see **Table 2** for the combinations of adopted resolution and positional error in existing studies). Indeed, with today's availability of high spatial resolution predictors, misuse of positionally inaccurate species

occurrences is increasingly likely, with the risk of exacerbating the negative effect of positional error on SDMs' performances.

To reduce the effect of positional error, multiple studies suggested adjusting spatial resolution so that the largest positional error associated with occurrence data is lower than the spatial resolution of the predictors (Engler et al. 2004; Moudrý and Šímová 2012; Keil et al. 2014; Vollering et al. 2016; Sillero et al. 2021a). However, coarsening the spatial resolution of the environmental variables may degrade information on fine-scale heterogeneity in environmental variables, eventually reducing their explanatory power for predicting species distribution (Mertes and Jetz 2018). In addition, spatial resolution can be coarsened to a level that is too far from the relevant ecological scale (Lecours et al. 2015, Moudrý et al. 2023). Recently, Gábor et al. (2022) showed that coarsening the spatial resolution to compensate for positional error does not improve model performance. However, they used a relatively simple virtual species approach, so more studies preferably using "real" species are needed to validate their results. Whether maintaining the spatial resolution of the response variable close to the ecological scale is more important than minimizing the adverse effect of positional error (or whether the opposite is true) remains a very current and unanswered question.

It is crucial to recognize that shifting species records in the geographic space does not necessarily translate to an equivalent shift in the environmental space. High positional error can lead to mischaracterizing the conditions under which a species occurs, especially in regions characterised by steep ecological gradients, such as mountainous areas or heavily fragmented landscapes. Indeed, the impact of positional error is related to the spatial autocorrelation of environmental variables. Naimi et al. (2011) found that the impact of positional error on SDMs' prediction performance decreased with increasing spatial autocorrelation in the environmental variables. In this regard, examining the degree of spatial autocorrelation in environmental variables was suggested as a way to a priori assess the impact of positional error on SDMs predictions (Naimi et al. 2011; 2014).

### **3.5. Recommendations associated with positional error**

It is crucial to consider data quality and to carefully assess the implications of using inaccurate data in either the training or validation process. Such considerations will yield more reliable assessments of model performance and improve the accuracy of SDMs.

- First of all, we recommend "cleaning" the dataset and removing aberrant errors (e.g., records with switched latitude and longitude, or records located at ZOOS or botanical gardens). This can be performed using automated methods such as the *CoordinateCleaner* package (Zizka et al. 2019).
- Second, researchers should quantify the positional error of the remaining input data, for example, using attributes specifying positional error. If such assessment is limited by metadata availability, for example in the case of historical data, it is recommended to at least approximate the positional error based on known information, such as the collection methodology or the number of decimals recorded with coordinates (e.g., Watcharamongkol et al. 2018; Moudrý and Devillers 2020).
- Third, we recommend that researchers carefully weigh the trade-offs between the positional error and resolution of environmental variables, with greater emphasis on the use of a

resolution as close to the ecological scale as possible (Gábor et al. 2022; Moudrý et al. 2023). Preferably, the positional error should be lower than the spatial resolution of the environmental variables (Moudrý and Šímová 2012). We suggest that the spatial resolution should be at least twice the positional error to reduce the risk of miscalculation of species-environment relationships. However, this may not always be achievable. In such a case, it is important to consider the following aspects to estimate and acknowledge the potential impact of positional error on the performance of the model.

- Fourth, we suggest considering positional error in the light of the particular species' ecology as some groups of species, such as mobile species, might be less affected by positional error than others (Gábor et al. 2020b).
- Fifth, researchers should examine the spatial autocorrelation in predictors to gain insight into whether predictions are likely to be affected by the positional error (Naimi et al. 2011; 2014). This may include testing the impact of various resolutions on model performance.
- Finally, we recommend considering the use of methods to mitigate positional error (Hefley et al. 2014; Zhang et al. 2018; Smith et al. 2023). Alternatively, knowing the positional error of the occurrences allows the inclusion of covariate predictors in the model to be tested or to up/downweight the importance of observations (Moudrý and Šímová 2012; Velásquez-Tibatá et al. 2016). For new surveys, we suggest using measurement techniques that minimize positional error, such as differential GNSS (e.g. Sillero et al. 2021b), and providing an estimate of the measurement accuracy (as is increasingly common in global databases).

#### 4. SAMPLING BIAS

Sampling bias poses a significant challenge in SDMs, leading to models that provide a partial or distorted view of species distribution or ecological niche (Kadmon et al. 2004; Leitao et al. 2011; Bean et al. 2012; Beck et al. 2014; Bardon et al. 2021). Despite advances, our knowledge of species distributions still remains limited for most taxa due to the variations in the sampling intensity over time and huge regions of the world remaining poorly sampled (Isaac and Pocock 2015; Menegotto and Rangel 2018; Hughes et al. 2021; Daru and Rodriguez 2023). Typically, positive sampling biases have been reported towards easily accessible areas (e.g. proximity to roads, rivers and urban settlements, Kadmon et al. 2004), protected areas (Boakes et al. 2010; Girardello et al. 2019), more populated areas (Geldmann et al. 2016), and charismatic species (Troudet et al. 2017), leading to spatial and taxonomic biases (Hughes et al. 2021). Uneven data-sharing practices further exacerbate this issue (Meyer et al. 2015).

Various methods have been proposed to compensate for sampling bias in species occurrence records, aiming to create models with quality comparable to models developed with unbiased data. Prevalent approaches for bias compensation include adjusting background samples (target-group background, TGB, approach; Phillips et al. 2009) in presence-background models or filtering (thinning) presences (Veloz 2009) (**Table 3**). The rationale behind the TGB is to select background data with the same sampling bias as for the set of presence records (i.e. to bias background locations towards areas where the presences were sampled; Phillips et al. 2009). The filtering approach was designed to reduce the negative effect of sampling bias by reducing the number of presences in oversampled regions in the geographic space (Veloz 2009) or oversampled environmental conditions in the environmental space

(Varela et al. 2014). Both geographic and environmental filtering use a distance between presences to determine the filter size. However, while geographic filtering uses distances in the geographic space (e.g., latitude and longitude) environmental filtering uses the range between values of multiple environmental variables (Varela et al. 2014; Castellanos et al. 2019). Another strategy carried out in the environmental space is to use presence data (i.e., their position in the environmental space) to identify and filter out background points associated with suitable habitats (Da Re et al. 2023). Many studies have evaluated the performance of these methods, simulating the bias by sub-sampling the original data (i.e. a presumably complete dataset without any bias) or by addressing bias already present in the datasets (**Table 3**). Such assessments require independent evaluation data containing both presence and absence records or comparison against models based on the unbiased dataset before sub-sampling simulation.

#### **4.1. Should the bias be assessed in the geographic or environmental space?**

There is an ongoing debate about whether bias should be assessed in the geographic or environmental space, or both (e.g., Varela et al. 2014; Moudrý 2015; Cosentino and Maiorano 2021). Indeed, the challenge in estimating species-environment relationships lies not only in the spatial bias within the geographic space where the bias originates but also in how this bias is reflected in the environmental space (i.e. the ecological niche space). All SDMs are not purely spatial methods (like interpolation, for instance), and the calculations actually occur within the environmental space defining the species' ecological niche. Therefore, addressing bias within the environmental space directly tackles the model calibration.

Sampling bias is influenced by the sampling design (Hirzel and Guisan 2002; Tessarolo et al. 2014; Mateo et al. 2018; Bazzichetto et al. 2023). A fundamental assumption underlying presence-background methods is that environmental conditions are sampled in proportion to their actual availability (Hastie and Fithian 2013). If not, clustered occurrences (i.e. geographic bias) may lead to the overestimation of the environmental suitability for the respective species in environments that have been sampled more intensively and underestimated for those surveyed less intensively (Barry and Elith 2006; Guillera-Arroita et al. 2015). For instance, fully random draws of species' presence in the geographic space may introduce a bias towards the most widespread environmental conditions, which possibly leads to uneven sampling of the species' realized niche within the environmental space (Bazzichetto et al. 2023). This issue is associated with another underlying assumption: that the species' niche is comprehensively sampled across the entire spectrum of environmental conditions in which it occurs (Phillips et al. 2009). Failing to meet this assumption, which can happen when there is a lack of knowledge about a species' tolerance to abiotic conditions (i.e. environmental bias), may cause a poor estimation of the actual niche occupied by the species (Hortal et al. 2008). If the ecological niche of the species is truncated (i.e. the complete niche of the species is not captured by the occurrences), it is not possible to extrapolate a reliable model into different spatial or temporal dimensions (Chevalier et al. 2022). Therefore, representative sampling of the environmental space should in principle give better results, regardless of their bias in the geographic space (Tessarolo et al. 2014; Sabatini et al. 2021; Bazzichetto et al. 2023).

Geographic and environmental spaces are communicating vessels, and so correcting one component (geographic or environmental) may have a detrimental effect on the other. For example, geographical

filtering could unwittingly exclude occurrences in the environmental space with unique environmental conditions (Varela et al. 2014). On the other hand, environmental filtering (down-weighting repeated species occurrences in similar environmental conditions) identifies grid cells within marginal habitats to be equally suitable as the cells representing core habitats. For example, if the species probability of occurrence is 0.1 at one site and 0.7 at another, such sites will be occupied in one and seven out of 10 cases, respectively. If we disregard the presences at the latter site, we lose the ability to discern the conditions favoured by the species (Moudrý et al. 2015). Indeed, it is impossible to use presence-background data to determine whether species observed in particular environments result from a larger sampling effort or ecological preferences (Guillera-Arroita et al. 2015), and removing bias without the information on the sampling effort becomes quixotic (Rocchini et al. 2023).

We recommend considering both geographic and environmental spaces in the assessment of sampling bias (e.g. Tessarolo et al. 2014; Cosentino and Maiorano 2021). In areas of high geographic and high environmental bias, and particularly in undersampled environments, further sampling efforts are required. Alternatively, bias correction based on the TGB method or geographic filtering can be a suitable option (e.g. Inman et al. 2021). However, a bias in the geographic space does not necessarily lead to a bias in the environmental space. If the geographic bias is high but the environmental bias is low, no corrections are needed, and the data can be used 'as is' for modelling. For example, Kadmon et al. (2004) and more recently McCarthy et al. (2012) showed that the road network provided a good sample of environmental gradients, and allowed uncovering of true species-environment relationships. In the case of low geographic but high environmental bias, it is reasonable to consider directly a correction in environmental space using environmental filtering (Varela et al. 2014; Cosentino and Maiorano 2021).

#### **4.2. How sampling bias (and correction methods) interact with species ecology**

Several studies have reported that there was no improvement or even detrimental effects on SDMs performance after filtering out biased samples (e.g. Chefaoui and Serrão 2017; Ranc et al. 2017; Gábor et al. 2020a), and it has been suggested that this might be related to species ecology (e.g. Bystrakova et al. 2012). For example, Ranc et al. (2017) showed that range size was the most important factor driving species vulnerability to sampling bias and that widespread species were more affected by sampling bias and more likely to benefit from bias correction than species with narrow geographic ranges. Similarly, Baker et al. (2022) showed that species type has a notable effect on model performance, with models generally being more robust to the effects of sampling bias for specialist (narrow environmental niches) than for generalist (wide environmental niches) species. In addition, a few studies highlighted that bias correction was detrimental for species with narrow ranges (Ranc et al. 2017), narrow niches (Inman et al. 2021), or low prevalence (Gábor et al. 2020a) and yielded worse models than without bias correction. It is evident that different species are differently affected by sampling bias and respond differently to bias correction. Therefore, species ecology should be considered when correcting for sampling bias. We recommend not to use bias correction methods for specialist species.

**Table 3.** Studies that evaluated the effect of sampling bias and the effectiveness of methods proposed to compensate for sampling bias on model performance.

Study	Number of species	Bias type	Evaluation approach	Bias correction	Main conclusion
Phillips et al. (2009)	226	Existing	Independent data	TGB	Bias correction improve models
Bystriakova et al. (2012)	5 plants ( <i>Asplenium spp.</i> )	Existing	Independent data (but only presences)	TGB	Bias correction improve models
Kramer Schadt et al. (2013)	Malay civet, two virtual species	Existing, Simulated	Simulated data	Geographic filtering, TGB	Geographic filter is preferred relative to TGB
Syfert et al. (2013)	Tree fern	Existing	Independent data	TGB	Bias correction improve models
Fourcade et al. (2014)	Turtle, salamander, virtual species	Simulated	Original model based on unbiased data	Five methods	Variable efficiency, further research needed
Varela et al. (2014)	Virtual	Simulated	Original model based on unbiased data	Environmental and geographic filtering	Recommend environmental filtering
Ranc et al. (2017)	Virtual	Simulated	True distribution of simulated species	TGB	Bias correction is detrimental for some species
Castellanos et al. (2019)	Virtual	Simulated	True distribution of simulated species	Environmental and geographic filtering	Recommend environmental filtering
Gábor et al. (2020a)	Virtual	Simulated	True distribution of simulated species	Environmental filtering	Filtering is not necessarily helpful
Chauvier et al. (2021)	1,900 plants	Existing	Independent data	Bias covariate correction, and environmental bias correction	Combining both methods might be the best choice
Inman et al. (2021)	Virtual	Simulated	True distribution of simulated species	TGB, geographic and environmental filtering	Bias correction is detrimental for some species
Baker et al. (2022)	Virtual	Simulated	True distribution of simulated species	Geographic filtering	More mechanistic understanding of how sampling biases arise is needed

### 4.3. Recommendations associated with sampling bias

Complete elimination of spatial bias from the modelling procedure is impossible without proper knowledge of all the processes generating it (Rocchini et al. 2023), and it is unrealistic to assume that sampling bias in biodiversity data can be eliminated, even with the development of automated observation technologies. Hence, SDMs need to explore and acknowledge the inherent biases associated with the data in both the geographic and environmental space (Cosentino and Maiorano 2021; Rocchini et al. 2023).

- First, researchers should quantify the sampling bias of their input data in the geographic space. For example, the *sampbias* algorithm (Zizka et al. 2021) can be used for such purposes.

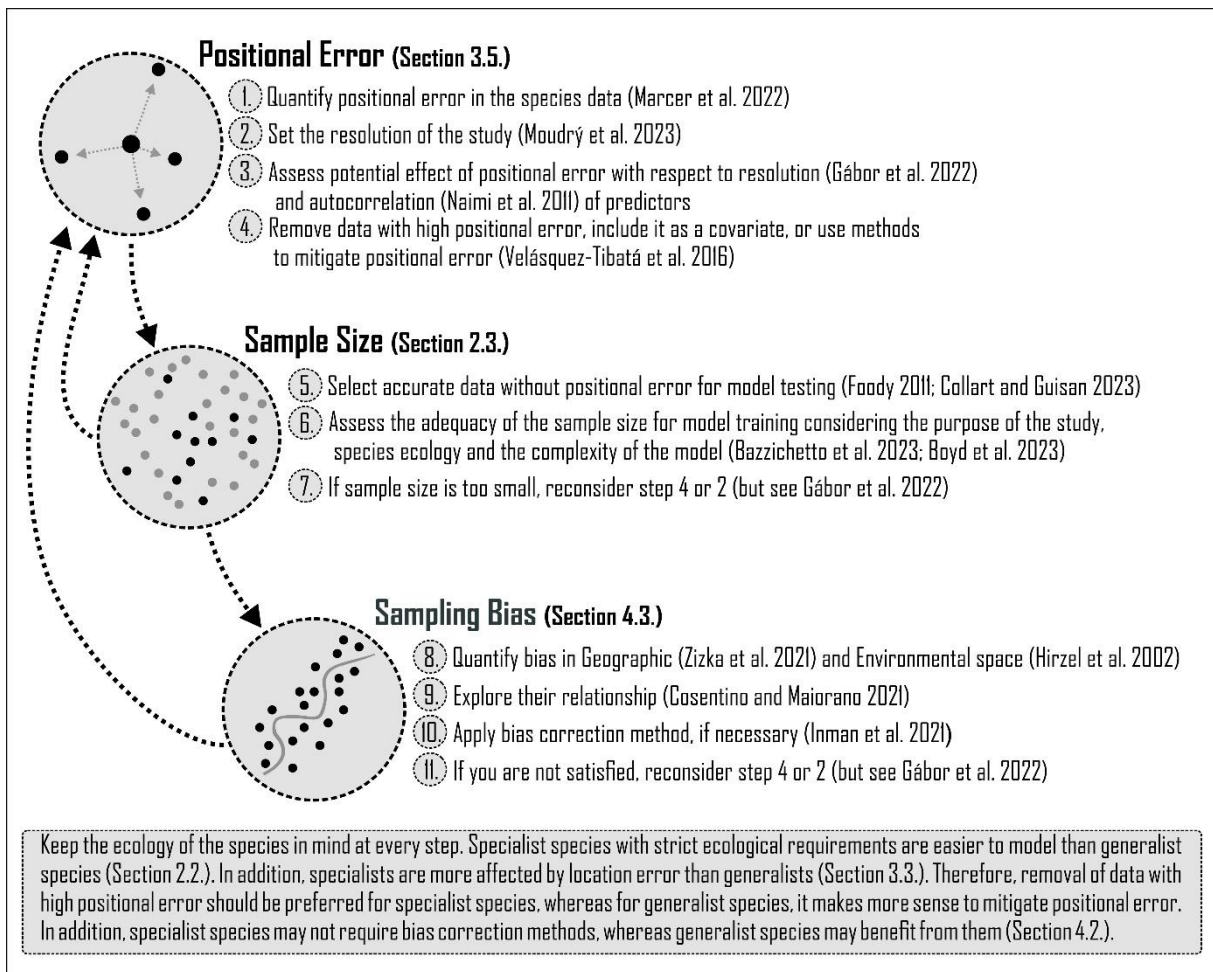


- Second, bias should also be evaluated in the environmental space by comparing the distribution of the cells where the focal species was present to all cells in the study area in a gridded environmental space of ecological predictors. This can be done, for example, by using Ecological Niche Factor Analysis (Hirzel et al. 2002); *hypervolume* R package (Blonder et al. 2014); or principal component analysis in the *ecospat* R package (Di Cola et al. 2017).
- The relationship between geographic and environmental bias should be further explored using local indicators of spatial association (LISA; Anselin 1995) and the results of such an assessment should be used as a basis for the selection of bias-correction methods (Cosentino and Maiorano 2021; Rocchini et al. 2023). This quantification can also assist researchers in effectively directing their further sampling efforts.
- The next step lies in the application of the bias-correction method, if necessary. Filtering or the TGB approach are possible options, but caution is needed as it could result in lower model performance in particular cases. This requires consideration of species' ecology as specialist species typically do not benefit from bias correction or can even be negatively affected by it (Gábor et al. 2020a; Inman et al. 2021; Baker et al. 2022). In addition, it is important to notice that filtering will inevitably reduce the number of presences available for modelling (but see Da Re et al. 2023 for filtering of background points). Therefore, if the sample size is relatively small, the TGB approach might be a preferred method.

## 5. GUIDELINES AND FUTURE DIRECTIONS

Despite the increasing number of studies focusing on how various limitations inherent to species data affect the performance of SDMs, there are still gaps in our knowledge and the use of SDMs remains problematic in many contexts. To advance our understanding, future studies should focus on comprehensive analyses that simultaneously consider various issues, such as sample size, sampling bias in the geographic and environmental space, positional error, spatial resolution, and species' ecological characteristics (**Figure 1**). Such studies can help establish the urgently needed guidelines for better-informed modelling choices (e.g. bias correction, removal of data with high positional error) concerning data limitations and species ecology. For instance, the consideration of data limitations becomes particularly important for specialist species of high conservation concern, where SDMs may be the only feasible means of estimating their distribution and responses to environmental changes. Regarding species characteristics, it is important to do such evaluations on characteristics that are easy to specify (i.e. we know them for the majority of species), such as dispersal ability, body size, or trophic group. This way, the assessments can be further used to guide data selection processes in other studies.

Finally, it is crucial to transparently report any potential biases and errors in the data used for modelling. Whenever possible, rigorous tests should be conducted to examine the impact of these biases and errors on model performance. Until more comprehensive assessments are available, we strongly recommend remaining vigilant about data limitations and following the basic guidelines for a critical assessment of spatial data to be used in SDMs shown in **Figure 3**.



**Figure 3.** Workflow for a critical assessment of spatial data to be used in SDMs. For more information on the individual steps, we refer the reader to the Recommendations subsections at the end of each main section.

## GLOSSARY BOX

**Ecological niche:** Hutchinsonian niche, which is defined as a hypothetical hypervolume spanned by the eco-physiological responses of a species to all environmental factors affecting its fitness.

**Model complexity** refers to the level of intricacy and flexibility in the representation of a species' ecological niche. It reflects how well the model can capture the underlying relationships between predictors and species distribution. The choice of model complexity depends on the nature of the problem, the amount and quality of available data, the number of model parameters, and the available computational resources. Finding the right balance between a model's ability to capture patterns and its potential for overfitting is a key challenge in building effective models.

**Model performance:** Here intended in a broad sense as a model capacity of recovering the underlying species-environment relationship using available data ('explanatory' performance), while also being able to extend (predict) out of the sample used for training/calibration ('predictive' performance).

**Model training** is the process of teaching a machine learning or statistical model to make predictions based on data. It is a crucial step in building and developing predictive models. Model training involves using a dataset with known outcomes to enable the model to learn the underlying patterns and relationships in the data.

**Model testing**, also known as model evaluation, is the process of assessing the performance and effectiveness of a machine learning or statistical model using a separate (independent) dataset that the model has not seen during training. The primary purpose of model testing is to determine how well the trained model generalizes to new, unseen data and to assess its predictive accuracy and reliability.

**Positional error** in species occurrence data refers to inaccuracies or uncertainty in the recorded coordinates of where a species was observed or collected. This error can result from factors such as imprecise GNSS measurements, data entry mistakes, or a lack of accurate location information.

**Spatial resolution or grain** refers to the level of detail or granularity at which data is collected, represented, or analysed in a spatial context. It can also be thought of as the size of the smallest spatial unit in a dataset (i.e. pixel size).

**Sampling design** refers to the approach used to collect species occurrence data. The sampling design is a crucial aspect of SDM, as it should in principle ensure that the data include all relevant information to represent the ecological niche of the species and the environmental conditions in the study area. The quality and representativeness of the data collected directly impact the accuracy and reliability of the model.

**Sample size:** The size of the data sample used to train and validate the model. Here, we define sample size as the total number of presences and absences (i.e. presence-absence data). When discussing studies based on presence-background data, we will refer specifically to the number of presences.

**Sampling bias:** Species occurrence records typically exhibit spatial bias, wherein some locations or environmental conditions are more intensively sampled than others. People sample accessible locations more intensively than remote or unpopular ones. This type of bias means that the available data used as the response variable fail to represent the complete niche of the species.

## ACKNOWLEDGEMENTS

Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them. VM, RR, AFC, VB, DR, RCG, PS were supported by the Horizon Europe project EarthBridge (Grant agreement No 101079310). MB acknowledges funding from the European Union's Horizon Europe research and innovation programme under the Marie Skłodowska-Curie grant agreement No 101066324. RGM was funded by project grants Connect2restore (TED2021-129589B-I00, funded by MCIN/AEI/10.13039/501100011033 and by the European Union NextGenerationEU/PRTR), and NextDive (PID2021-124187NB-I00, funded by MCIN/AEI/10.13039/501100011033 and by ERDF, a way of making Europe). AZA was supported by a Margarita Salas Contract financed by the European Union-NextGenerationEU, Ministerio de Universidades y Plan de Recuperacion, Transformacion y Resiliencia, Spain. MM, JW and JP were supported by the Czech Academy of Sciences (project RVO 67985939). FL was funded by the European Union (ERC, BEAST, 101044740). JIL was supported by BiodivERsA+ project ASICS (BiodivMon-call 2021-2022). NS is supported by a CEEC2017 contract (CEECIND/02213/2017) from FCT - Fundação para a Ciência e a Tecnologia, Portugal. MT was partially funded by the European Union's Horizon 2020 research and innovation program under grant agreement No 862480 (SHOWCASE).

## REFERENCES

- Anselin, L. (1995) Local indicators of spatial association—LISA. *Geographical analysis*, 27, 93-115.
- Araújo, M. B., Anderson, R. P., Márcia Barbosa, A., Beale, C. M., Dormann, C. F., Early, R., ... & Rahbek, C. (2019) Standards for distribution models in biodiversity assessments. *Science advances*, 5, eaat4858.
- Arenas-Castro, S., Regos, A., Martins, I., Honrado, J., & Alonso, J. (2022) Effects of input data sources on species distribution model predictions across species with different distributional ranges. *Journal of Biogeography*, 49, 1299-1312.
- Austin, M. P. (2002) Spatial prediction of species distribution: an interface between ecological theory and statistical modelling. *Ecological modelling*, 157, 101-118.
- Baker, D. J., Maclean, I. M., Goodall, M., & Gaston, K. J. (2022) Correlations between spatial sampling biases and environmental niches affect species distribution models. *Global Ecology and Biogeography*, 31, 1038-1050.
- Bardon, L. R., Ward, B. A., Dutkiewicz, S., & Cael, B. B. (2021) Testing the skill of a species distribution model using a 21st century virtual ecosystem. *Geophysical Research Letters*, 48, e2021GL093455.
- Barry, S., & Elith, J. (2006) Error and uncertainty in habitat models. *Journal of Applied Ecology*, 43, 413-423.
- Bazzichetto, M., Lenoir, J., Da Re, D., Tordoni, E., Rocchini, D., Malavasi, M., ... & Sperandii, M. G. (2023) Sampling strategy matters to accurately estimate response curves' parameters in species distribution models. *Global Ecology and Biogeography*, 32, 1717-1729.
- Bazzichetto, M., Massol, F., Carboni, M., Lenoir, J., Lembrechts, J. J., Joly, R., & Renault, D. (2021) Once upon a time in the far south: Influence of local drivers and functional traits on plant invasion in the harsh sub-Antarctic islands. *Journal of Vegetation Science*, 32, e13057.
- Bean, W. T., Stafford, R., & Brashares, J. S. (2012) The effects of small sample size and sample bias on threshold selection and accuracy assessment of species distribution models. *Ecography*, 35, 250-258.
- Beck, J., Böller, M., Erhardt, A., & Schwanghart, W. (2014) Spatial bias in the GBIF database and its effect on modeling species' geographic distributions. *Ecological Informatics*, 19, 10-15.

- Bell, D. M., & Schlaepfer, D. R. (2016) On the dangers of model complexity without ecological justification in species distribution modeling. *Ecological Modelling*, 330, 50-59.
- Blonder, B., Lamanna, C., Violle, C., & Enquist, B. J. (2014) The n-dimensional hypervolume. *Global Ecology and Biogeography*, 23, 595-609.
- Bloom, T. D., Flower, A., & DeChaine, E. G. (2018) Why georeferencing matters: Introducing a practical protocol to prepare species occurrence records for spatial analysis. *Ecology and Evolution*, 8, 765-777.
- Boakes, E. H., McGowan, P. J., Fuller, R. A., Chang-qing, D., Clark, N. E., O'Connor, K., & Mace, G. M. (2010) Distorted views of biodiversity: spatial and temporal bias in species occurrence data. *PLoS biology*, 8, e1000385.
- Botella, C., Deneu, B., Marcos, D., Servajean, M., Estopinan, J., Larcher, T., ... & Joly, A. (2023) The GeoLifeCLEF 2023 Dataset to evaluate plant species distribution models at high spatial resolution across Europe. arXiv preprint arXiv:2308.05121.
- Boyd, R. J., Harvey, M., Roy, D. B., Barber, T., Haysom, K. A., Macadam, C. R., ... & Pescott, O. L. (2023) Causal inference and large-scale expert validation shed light on the drivers of SDM accuracy and variance. *Diversity and Distributions*.
- Bystrakova, N., Peregrym, M., Erkens, R. H., Bezsmertna, O., & Schneider, H. (2012) Sampling bias in geographic and environmental space and its effect on the predictive power of species distribution models. *Systematics and biodiversity*, 10, 305-315.
- Carretero, M. A., & Sillero, N. (2016) Evaluating how species niche modelling is affected by partial distributions with an empirical case. *Acta Oecologica*, 77, 207-216.
- Castellanos, A. A., Huntley, J. W., Voelker, G., & Lawing, A. M. (2019) Environmental filtering improves ecological niche models across multiple scales. *Methods in Ecology and Evolution*, 10, 481-492.
- Chauvier, Y., Zimmermann, N. E., Poggiato, G., Bystrova, D., Brun, P., & Thuiller, W. (2021) Novel methods to correct for observer and sampling bias in presence-only species distribution models. *Global Ecology and Biogeography*, 30, 2312-2325.
- Chevalier, M., Zarzo-Arias, A., Guélat, J., Mateo, R. G., & Guisan, A. (2022) Accounting for niche truncation to improve spatial and temporal predictions of species distributions. *Frontiers in Ecology and Evolution*, 10.
- Chefaoui, R. M., & Serrão, E. A. (2017) Accounting for uncertainty in predictions of a marine species: integrating population genetics to verify past distributions. *Ecological Modelling*, 359, 229-239.
- Collart, F., & Guisan, A. (2023) Small to train, small to test: Dealing with low sample size in model evaluation. *Ecological Informatics*, 75, 102106.
- Collart, F., Broennimann, O., Guisan, A., & Vanderpoorten, A. (2023) Ecological and biological indicators of the accuracy of species distribution models: lessons from European bryophytes. *Ecography*, e06721.
- Cosentino, F., & Maiorano, L. (2021) Is geographic sampling bias representative of environmental space?. *Ecological Informatics*, 64, 101369.
- Coudun, C., & Gégout, J. C. (2006) The derivation of species response curves with Gaussian logistic regression is sensitive to sampling intensity and curve characteristics. *Ecological modelling*, 199, 164-175.
- Da Re, D., Tordoni, E., Lenoir, J., Vanwambeke, S. O., Rocchini, D., Bazzichetto, M., & SoilTemp Consortium. (2023) USE it: uniformly sampling pseudo-absences within the environmental space for applications in habitat suitability models.
- Daru, B. H., & Rodriguez, J. (2023) Mass production of unvouchered records fails to represent global biodiversity patterns. *Nature Ecology & Evolution*, 1-16.

- Di Cola, V., Broennimann, O., Petitpierre, B., Breiner, F. T., d'Amen, M., Randin, C., ... & Guisan, A. (2017) ecospat: an R package to support spatial analyses and modeling of species niches and distributions. *Ecography*, 40, 774-787.
- Duputié, A., Zimmermann, N. E., & Chuine, I. (2014) Where are the wild things? Why we need better data on species distribution. *Global Ecology and Biogeography*, 23, 457-467.
- Ehrlén, J., & Morris, W. F. (2015) Predicting changes in the distribution and abundance of species under environmental change. *Ecology letters*, 18, 303-314.
- Elith, J., Burgman, M. A., & Regan, H. M. (2002) Mapping epistemic uncertainties and vague concepts in predictions of species distribution. *Ecological modelling*, 157, 313-329.
- Elith, J., & Leathwick, J. R. (2009) Species distribution models: ecological explanation and prediction across space and time. *Annual review of ecology, evolution, and systematics*, 40, 677-697.
- Engler, R., Guisan, A., & Rechsteiner, L. (2004) An improved approach for predicting the distribution of rare and endangered species from occurrence and pseudo-absence data. *Journal of applied ecology*, 41, 263-274.
- Feeley, K. J., & Silman, M. R. (2011) Keep collecting: accurate species distribution modelling requires more collections than previously thought. *Diversity and distributions*, 17, 1132-1140.
- Feng, X., Park, D. S., Walker, C., Peterson, A. T., Merow, C., & Papeş, M. (2019) A checklist for maximizing reproducibility of ecological niche models. *Nature Ecology & Evolution*, 3, 1382-1395.
- Fernandez, M., Blum, S., Reichle, S., Guo, Q., Holzman, B., & Hamilton, H. (2009) Locality uncertainty and the differential performance of four common niche-based modeling techniques. *Biodiversity Informatics*, 6.
- Ferrier, S., Jetz, W., & Scharlemann, J. (2017) Biodiversity modelling as part of an observation system. *The GEO handbook on biodiversity observation networks*, 239-257.
- Ficetola, G. F., Bonardi, A., Múcher, C. A., Gilissen, N. L., & Padoa-Schioppa, E. (2014) How many predictors in species distribution models at the landscape scale? Land use versus LiDAR-derived canopy height. *International Journal of Geographical Information Science*, 28, 1723-1739.
- Fois, M., Cuenca-Lombraña, A., Fenu, G., & Bacchetta, G. (2018) Using species distribution models at local scale to guide the search of poorly known species: Review, methodological issues and future directions. *Ecological Modelling*, 385, 124-132.
- Foody, G. M. (2011) Impacts of imperfect reference data on the apparent accuracy of species presence-absence models and their predictions. *Global Ecology and Biogeography*, 20, 498-508.
- Fourcade, Y., Engler, J. O., Rödder, D., & Secondi, J. (2014) Mapping species distributions with MAXENT using a geographically biased sample of presence data: a performance assessment of methods for correcting sampling bias. *PloS one*, 9, e97122.
- Fourcade, Y., Besnard, A. G., & Secondi, J. (2018) Paintings predict the distribution of species, or the challenge of selecting environmental predictors and evaluation statistics. *Global Ecology and Biogeography*, 27, 245-256.
- Frair, J. L., Fieberg, J., Hebblewhite, M., Cagnacci, F., DeCesare, N. J., & Pedrotti, L. (2010) Resolving issues of imprecise and habitat-biased locations in ecological analyses using GPS telemetry data. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365, 2187-2200.
- Gábor, L., Moudrý, V., Barták, V., & Lecours, V. (2020a) How do species and data characteristics affect species distribution models and when to use environmental filtering?. *International Journal of Geographical Information Science*, 34, 1567-1584.

- Gábor, L., Moudrý, V., Lecours, V., Malavasi, M., Barták, V., Fogl, M., ... & Václavík, T. (2020b) The effect of positional error on fine scale species distribution models increases for specialist species. *Ecography*, 43, 256-269.
- Gábor, L., Jetz, W., Lu, M., Rocchini, D., Cord, A., Malavasi, M., ... & Moudrý, V. (2022) Positional errors in species distribution modelling are not overcome by the coarser grains of analysis. *Methods in Ecology and Evolution*, 13, 2289-2302.
- Gabor, L., Cohen, J., & Jetz, W. (2023a) Assessing the impact of binary land cover variables on species distribution models: A North American study on water birds. *bioRxiv*, 2023-07.
- Gábor, L., Jetz, W., Zarzo-Arias, A., Winner, K., Yanco, S., Pinkert, S., ... & Moudrý, V. (2023b) Species distribution models affected by positional uncertainty in species occurrences can still be ecologically interpretable. *Ecography*, e06358.
- Geldmann, J., Heilmann-Clausen, J., Holm, T. E., Levinsky, I., Markussen, B. O., Olsen, K., ... & Tøttrup, A. P. (2016) What determines spatial bias in citizen science? Exploring four recording schemes with different proficiency requirements. *Diversity and Distributions*, 22, 1139-1149.
- Girardello, M., Chapman, A., Dennis, R., Kaila, L., Borges, P. A., & Santangeli, A. (2019) Gaps in butterfly inventory data: A global analysis. *Biological conservation*, 236, 289-295.
- Graham, C. H., Ferrier, S., Huettman, F., Moritz, C., & Peterson, A. T. (2004) New developments in museum-based informatics and applications in biodiversity analysis. *Trends in ecology & evolution*, 19, 497-503.
- Graham, C. H., Elith, J., Hijmans, R. J., Guisan, A., Townsend Peterson, A., Loisele, B. A., & NCEAS Predicting Species Distributions Working Group. (2008) The influence of spatial errors in species occurrence data used in distribution models. *Journal of Applied Ecology*, 45, 239-247.
- Guillera-Arroita, G., Lahoz-Monfort, J. J., Elith, J., Gordon, A., Kujala, H., Lentini, P. E., ... & Wintle, B. A. (2015) Is my species distribution model fit for purpose? Matching data and models to applications. *Global ecology and biogeography*, 24, 276-292.
- Guisan, A., Zimmermann, N. E., Elith, J., Graham, C. H., Phillips, S., & Peterson, A. T. (2007) WHAT MATTERS FOR PREDICTING THE OCCURRENCES OF TREES: TECHNIQUES, DATA, OR SPECIES CHARACTERISTICS?. *Ecological monographs*, 77, 615-630.
- Guisan, A., Tingley, R., Baumgartner, J. B., Naujokaitis-Lewis, I., Sutcliffe, P. R., Tulloch, A. I., ... & Buckley, Y. M. (2013) Predicting species distributions for conservation decisions. *Ecology letters*, 16, 1424-1435.
- Haesen, S., Lenoir, J., Gril, E., De Frenne, P., Lembrechts, J. J., Kopecký, M., ... & Van Meerbeek, K. (2023) Microclimate reveals the true thermal niche of forest plant species. *Ecology Letters*.
- Hallman, T. A., & Robinson, W. D. (2020) Deciphering ecology from statistical artefacts: Competing influence of sample size, prevalence and habitat specialization on species distribution models and how small evaluation datasets can inflate metrics of performance. *Diversity and Distributions*, 26, 315-328.
- Hanberry, B. B., He, H. S., & Dey, D. C. (2012) Sample sizes and model comparison metrics for species distribution models. *Ecological Modelling*, 227, 29-33.
- Hastie, T., & Fithian, W. (2013) Inference from presence-only data; the ongoing controversy. *Ecography*, 36, 864-867.
- Hefley, T. J., Baasch, D. M., Tyre, A. J., & Blankenship, E. E. (2014) Correction of location errors for presence-only species distribution models. *Methods in Ecology and Evolution*, 5, 207-214.
- Heikkinen, R. K., Luoto, M., Araújo, M. B., Virkkala, R., Thuiller, W., & Sykes, M. T. (2006) Methods and uncertainties in bioclimatic envelope modelling under climate change. *Progress in Physical Geography*, 30, 751-777.

- Hernandez, P. A., Graham, C. H., Master, L. L., & Albert, D. L. (2006) The effect of sample size and species characteristics on performance of different species distribution modeling methods. *Ecography*, 29, 773-785.
- Hirzel, A., & Guisan, A. (2002) Which is the optimal sampling strategy for habitat suitability modelling. *Ecological modelling*, 157, 331-341.
- Hirzel, A. H., Hausser, J., Chessel, D., & Perrin, N. (2002) Ecological-niche factor analysis: how to compute habitat-suitability maps without absence data?. *Ecology*, 83, 2027-2036.
- Hortal, J., Jiménez-Valverde, A., Gómez, J. F., Lobo, J. M., & Baselga, A. (2008) Historical bias in biodiversity inventories affects the observed environmental niche of the species. *Oikos*, 117, 847-858.
- Hortal, J., de Bello, F., Diniz-Filho, J. A. F., Lewinsohn, T. M., Lobo, J. M., & Ladle, R. J. (2015) Seven shortfalls that beset large-scale knowledge of biodiversity. *Annual Review of Ecology, Evolution, and Systematics*, 46, 523-549.
- Hughes, A. C., Orr, M. C., Ma, K., Costello, M. J., Waller, J., Provoost, P., ... & Qiao, H. (2021) Sampling biases shape our view of the natural world. *Ecography*, 44, 1259-1269.
- Hughes, A., Dorey, J., Bossert, S., Qiao, H., & Orr, M. (2023) Big data-big problems? How to circumvent problems in biodiversity mapping and ensure meaningful results.
- Inman, R., Franklin, J., Esque, T., & Nussear, K. (2021) Comparing sample bias correction methods for species distribution modeling using virtual species. *Ecosphere*, 12, e03422.
- Isaac, N. J., & Pocock, M. J. (2015) Bias and information in biological records. *Biological Journal of the Linnean Society*, 115, 522-531.
- Jansen, J., Woolley, S. N., Dunstan, P. K., Foster, S. D., Hill, N. A., Haward, M., & Johnson, C. R. (2022) Stop ignoring map uncertainty in biodiversity science and conservation policy. *Nature Ecology & Evolution*, 6, 828-829.
- Jiménez-Valverde, A., Lobo, J., & Hortal, J. (2009) The effect of prevalence and its interaction with sample size on the reliability of species distribution models. *Community Ecology*, 10, 196-205.
- Jiménez-Valverde, A. (2020) Sample size for the evaluation of presence-absence models. *Ecological Indicators*, 114, 106289.
- Johnson, C. J., & Gillingham, M. P. (2008) Sensitivity of species-distribution models to error, bias, and model design: an application to resource selection functions for woodland caribou. *Ecological Modelling*, 213, 143-155.
- Kadmon, R., Farber, O., & Danin, A. (2003) A systematic analysis of factors affecting the performance of climatic envelope models. *Ecological Applications*, 13, 853-867.
- Kadmon, R., Farber, O., & Danin, A. (2004) Effect of roadside bias on the accuracy of predictive maps produced by bioclimatic models. *Ecological Applications*, 14, 401-413.
- Keil, P., Wilson, A. M., & Jetz, W. (2014) Uncertainty, priors, autocorrelation and disparate data in downscaling of species distributions. *Diversity and Distributions*, 20, 797-812.
- Kos, T., Markezic, I., & Pokrajcic, J. (2010, September) Effects of multipath reception on GPS positioning performance. In Proceedings ELMAR-2010 (pp. 399-402). IEEE.
- Kramer-Schadt, S., Niedballa, J., Pilgrim, J. D., Schröder, B., Lindenborn, J., Reinfelder, V., ... & Wilting, A. (2013) The importance of correcting for sampling bias in MaxEnt species distribution models. *Diversity and distributions*, 19, 1366-1379.
- Lecours, V., Devillers, R., Schneider, D. C., Lucieer, V. L., Brown, C. J., & Edinger, E. N. (2015) Spatial scale and geographic context in benthic habitat mapping: review and future directions. *Marine Ecology Progress Series*, 535, 259-284.



- Leitão, P. J., Moreira, F., & Osborne, P. E. (2011) Effects of geographical data sampling bias on habitat models of species distributions: a case study with steppe birds in southern Portugal. *International Journal of Geographical Information Science*, 25, 439-454.
- Liu, C., Newell, G., & White, M. (2019) The effect of sample size on the accuracy of species distribution models: considering both presences and pseudo-absences or background sites. *Ecography*, 42, 535-548.
- Loiselle, B. A., Jørgensen, P. M., Consiglio, T., Jiménez, I., Blake, J. G., Lohmann, L. G., & Montiel, O. M. (2008) Predicting species distributions from herbarium collections: does climate bias in collection sampling influence model outcomes?. *Journal of Biogeography*, 35, 105-116.
- Maggini, R., Lehmann, A., Zimmermann, N. E., & Guisan, A. (2006) Improving generalized regression analysis for the spatial prediction of forest communities. *Journal of biogeography*, 33, 1729-1749.
- Marcer, A., Chapman, A. D., Wieczorek, J. R., Xavier Picó, F., Uribe, F., Waller, J., & Ariño, A. H. (2022) Uncertainty matters: ascertaining where specimens in natural history collections come from and its implications for predicting species distributions. *Ecography*, 2022, e06025.
- Mateo, R. G., Felicísimo, Á. M., & Muñoz, J. (2010) Effects of the number of presences on reliability and stability of MARS species distribution models: the importance of regional niche variation and ecological heterogeneity. *Journal of Vegetation Science*, 21, 908-922.
- Mateo, R. G., Gaston, A., Aroca-Fernández, M. J., Saura, S., & García-Viñas, J. I. (2018) Optimization of forest sampling strategies for woody plant species distribution modelling at the landscape scale. *Forest Ecology and Management*, 410, 104-113.
- Mccarthy, K. P., Fletcher Jr, R. J., Rota, C. T., & Hutto, R. L. (2012) Predicting species distributions from samples collected along roadsides. *Conservation biology*, 26, 68-77.
- McPherson, J. M., & Jetz, W. (2007) Effects of species' ecology on the accuracy of distribution models. *Ecography*, 30, 135-151.
- McPherson, J. M., Jetz, W., & Rogers, D. J. (2004) The effects of species' range sizes on the accuracy of distribution models: ecological phenomenon or statistical artefact?. *Journal of applied ecology*, 41, 811-823.
- Menegotto, A., & Rangel, T. F. (2018) Mapping knowledge gaps in marine diversity reveals a latitudinal gradient of missing species richness. *Nature communications*, 9, 4713.
- Merow, C., Smith, M. J., Edwards Jr, T. C., Guisan, A., McMahon, S. M., Normand, S., ... & Elith, J. (2014) What do we gain from simplicity versus complexity in species distribution models?. *Ecography*, 37, 1267-1281.
- Mertes, K., & Jetz, W. (2018) Disentangling scale dependencies in species environmental niches and distributions. *Ecography*, 41, 1604-1615.
- Meyer, C., Kreft, H., Guralnick, R., & Jetz, W. (2015) Global priorities for an effective information basis of biodiversity distributions. *Nature communications*, 6, 1-8.
- Mitchell, P. J., Monk, J., & Laurenson, L. (2017) Sensitivity of fine-scale species distribution models to locational uncertainty in occurrence data across multiple sample sizes. *Methods in Ecology and Evolution*, 8, 12-21.
- Moreno-Amat, E., Mateo, R. G., Nieto-Lugilde, D., Morueta-Holme, N., Svenning, J. C., & García-Amorena, I. (2015) Impact of model complexity on cross-temporal transferability in Maxent species distribution models: An assessment using paleobotanical data. *Ecological Modelling*, 312, 308-317.
- Moudrý, V., & Devillers, R. (2020) Quality and usability challenges of global marine biodiversity databases: An example for marine mammal data. *Ecological Informatics*, 56, 101051.

- Moudrý, V., & Šímová, P. (2012) Influence of positional accuracy, sample size and scale on modelling species distributions: a review. *International Journal of Geographical Information Science*, 26, 2083-2095.
- Moudrý, V. (2015) Modelling species distributions with simulated virtual species. *Journal of Biogeography*, 42, 1365-1366.
- Moudrý, V., Komárek, J., & Šímová, P. (2017) Which breeding bird categories should we use in models of species distribution?. *Ecological Indicators*, 74, 526-529.
- Moudrý, V., Keil, P., Cord, A. F., Gábor, L., Lecours, V., Zarzo-Arias, A., ... & Šímová, P. (2023) Scale mismatches between predictor and response variables in species distribution modelling: A review of practices for appropriate grain selection. *Progress in Physical Geography: Earth and Environment*, 03091333231156362.
- Naimi, B., Skidmore, A. K., Groen, T. A., & Hamm, N. A. (2011) Spatial autocorrelation in predictors reduces the impact of positional uncertainty in occurrence data on species distribution modelling. *Journal of biogeography*, 38, 1497-1509.
- Naimi, B., Hamm, N. A., Groen, T. A., Skidmore, A. K., & Toxopeus, A. G. (2014) Where is positional uncertainty a problem for species distribution modelling?. *Ecography*, 37, 191-203.
- Newbold, T. (2010) Applications and limitations of museum data for conservation and ecology, with particular attention to species distribution models. *Progress in physical geography*, 34, 3-22.
- Osborne, P. E., & Leitao, P. J. (2009) Effects of species and habitat positional errors on the performance and interpretation of species distribution models. *Diversity and Distributions*, 15, 671-681.
- Papeş, M., & Gaubert, P. (2007) Modelling ecological niches from low numbers of occurrences: assessment of the conservation status of poorly known viverrids (Mammalia, Carnivora) across two continents. *Diversity and distributions*, 13, 890-902.
- Pearson, R. G., Raxworthy, C. J., Nakamura, M., & Townsend Peterson, A. (2007) Predicting species distributions from small numbers of occurrence records: a test case using cryptic geckos in Madagascar. *Journal of biogeography*, 34, 102-117.
- Phillips, S. J., Dudík, M., Elith, J., Graham, C. H., Lehmann, A., Leathwick, J., & Ferrier, S. (2009) Sample selection bias and presence-only distribution models: implications for background and pseudo-absence data. *Ecological applications*, 19, 181-197.
- Proosdij, A. S. J. van, Sosef, M. S., Wieringa, J. J., & Raes, N. (2016) Minimum required number of specimen records to develop accurate species distribution models. *Ecography*, 39, 542-552.
- Ranc, N., Santini, L., Rondinini, C., Boitani, L., Poitevin, F., Angerbjörn, A., & Maiorano, L. (2017) Performance tradeoffs in target-group bias correction for species distribution models. *Ecography*, 40, 1076-1087.
- Rattray, A., Ierodiaconou, D., Monk, J., Laurenson, L. J. B., & Kennedy, P. (2014) Quantification of spatial and thematic uncertainty in the application of underwater video for benthic habitat mapping. *Marine Geodesy*, 37, 315-336.
- Raxworthy, C. J., Martinez-Meyer, E., Horning, N., Nussbaum, R. A., Schneider, G. E., Ortega-Huerta, M. A., & Townsend Peterson, A. (2003) Predicting distributions of known and unknown reptile species in Madagascar. *Nature*, 426, 837-841.
- Reside, A. E., Watson, I., VanDerWal, J., & Kutt, A. S. (2011) Incorporating low-resolution historic species location data decreases performance of distribution models. *Ecological Modelling*, 222, 3444-3448.
- Rocchini, D., Hortal, J., Lengyel, S., Lobo, J. M., Jimenez-Valverde, A., Ricotta, C., ... & Chiarucci, A. (2011) Accounting for uncertainty when mapping species distributions: the need for maps of ignorance. *Progress in Physical Geography*, 35, 211-226.

- Rocchini, D., Tordoni, E., Marchetto, E., Marcantonio, M., Barbosa, A. M., Bazzichetto, M., ... & Malavasi, M. (2023) A quixotic view of spatial bias in modelling the distribution of species and their diversity. *npj Biodiversity*, 2, 10.
- Sabatini, F. M., Lenoir, J., Hattab, T., Arnst, E. A., Chytrý, M., Dengler, J., ... & Wagner, V. (2021) sPlotOpen—An environmentally balanced, open-access, global dataset of vegetation plots. *Global Ecology and Biogeography*, 30, 1740-1764.
- Santini, L., Benítez-López, A., Maiorano, L., Čengić, M., & Huijbregts, M. A. (2021) Assessing the reliability of species distribution projections in climate change research. *Diversity and Distributions*, 27, 1035-1050.
- Segal, R. D., Massaro, M., Carlile, N., & Whitsed, R. (2021) Small-scale species distribution model identifies restricted breeding habitat for an endemic island bird. *Animal Conservation*, 24, 959-969.
- Segurado, P., & Araujo, M. B. (2004). An evaluation of methods for modelling species distributions. *Journal of biogeography*, 31, 1555-1568.
- Seoane, J., Carrascal, L. M., Alonso, C. L., & Palomino, D. (2005) Species-specific traits associated to prediction errors in bird habitat suitability modelling. *Ecological Modelling*, 185, 299-308.
- Shiroyama, R., Wang, M., & Yoshimura, C. (2020) Effect of sample size on habitat suitability estimation using random forests: a case of bluegill, *Lepomis macrochirus*. *Annales de Limnologie-International Journal of Limnology*, 56.
- Sillero, N. (2011) What does ecological modelling model? A proposed classification of ecological niche models based on their underlying methods. *Ecological Modelling*, 222, 1343-1346.
- Sillero, N., & Barbosa, A. M. (2021) Common mistakes in ecological niche models. *International Journal of Geographical Information Science*, 35, 213-226.
- Sillero, N., & Goncalves-Seco, L. (2014) Spatial structure analysis of a reptile community with airborne LiDAR data. *International Journal of Geographical Information Science*, 28, 1709-1722.
- Sillero, N., Arenas-Castro, S., Enriquez-Urzelai, U., Vale, C. G., Sousa-Guedes, D., Martínez-Freiría, F., ... & Barbosa, A. M. (2021a) Want to model a species niche? A step-by-step guideline on correlative ecological niche modelling. *Ecological Modelling*, 456, 109671.
- Sillero, N., Dos Santos, R., Teodoro, A. C., & Carretero, M. A. (2021b) Ecological niche models improve home range estimations. *Journal of Zoology*, 313, 145-157.
- Smith, A. B., & Santos, M. J. (2020) Testing the ability of species distribution models to infer variable importance. *Ecography*, 43, 1801-1813.
- Smith, A. B., Murphy, S. J., Henderson, D., & Erickson, K. D. (2023) Including imprecisely georeferenced specimens improves accuracy of species distribution models and estimates of niche breadth. *Global Ecology and Biogeography*, 32, 342-355.
- Støa, B., Halvorsen, R., Stokland, J. N., & Gusarov, V. I. (2019) How much is enough? Influence of number of presence observations on the performance of species distribution models. *Sommerfeltia*, 39, 1-28.
- Stockwell, D. R., & Peterson, A. T. (2002) Effects of sample size on accuracy of species distribution models. *Ecological modelling*, 148, 1-13.
- Syfert, M. M., Smith, M. J., & Coomes, D. A. (2013) The effects of sampling bias and model complexity on the predictive performance of MaxEnt species distribution models. *PloS one*, 8, e55158.
- Tessarolo, G., Rangel, T. F., Araújo, M. B., & Hortal, J. (2014) Uncertainty associated with survey design in Species Distribution Models. *Diversity and Distributions*, 20, 1258-1269.

- Tessarolo, G., Ladle, R. J., Lobo, J. M., Rangel, T. F., & Hortal, J. (2021) Using maps of biogeographical ignorance to reveal the uncertainty in distributional data hidden in species distribution models. *Ecography*, 44, 1743-1755.
- Thibaud, E., Petitpierre, B., Broennimann, O., Davison, A. C., & Guisan, A. (2014) Measuring the relative effect of factors affecting species distribution model predictions. *Methods in Ecology and Evolution*, 5, 947-955.
- Troudet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R., & Legendre, F. (2017) Taxonomic bias in biodiversity data and societal preferences. *Scientific reports*, 7, 9132.
- Tsoar, A., Allouche, O., Steinitz, O., Rotem, D., & Kadmon, R. (2007) A comparative evaluation of presence-only methods for modelling species distribution. *Diversity and distributions*, 13, 397-405.
- Varela, S., Anderson, R. P., García-Valdés, R., & Fernández-González, F. (2014) Environmental filters reduce the effects of sampling bias and improve predictions of ecological niche models. *Ecography*, 37, 1084-1091.
- Velásquez-Tibatá, J., Graham, C. H., & Munch, S. B. (2016) Using measurement error models to account for georeferencing error in species distribution models. *Ecography*, 39, 305-316.
- Veloz, S. D. (2009) Spatially autocorrelated sampling falsely inflates measures of accuracy for presence-only niche models. *Journal of biogeography*, 36, 2290-2299.
- Vollering, J., Schuiteman, A., de Vogel, E., van Vugt, R., & Raes, N. (2016) Phytogeography of New Guinean orchids: patterns of species richness and turnover. *Journal of Biogeography*, 43, 204-214.
- Wang, L., & Jackson, D. A. (2023) Effects of sample size, data quality, and species response in environmental space on modeling species distributions. *Landscape Ecology*, 1-23.
- Watcharamongkol, T., Christin, P. A., & Osborne, C. P. (2018) C4 photosynthesis evolved in warm climates but promoted migration to cooler ones. *Ecology Letters*, 21, 376-383.
- Wieczorek, J., Guo, Q., & Hijmans, R. (2004) The point-radius method for georeferencing locality descriptions and calculating associated uncertainty. *International journal of geographical information science*, 18, 745-767.
- Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan, A., & NCEAS Predicting Species Distributions Working Group. (2008) Effects of sample size on the performance of species distribution models. *Diversity and distributions*, 14, 763-773.
- Wüest, R. O., Zimmermann, N. E., Zurell, D., Alexander, J. M., Fritz, S. A., Hof, C., ... & Karger, D. N. (2020) Macroecology in the age of Big Data—Where to go from here?. *Journal of Biogeography*, 47, 1-12.
- Zhang, G., Zhu, A. X., Huang, Z. P., & Xiao, W. (2018) A heuristic-based approach to mitigating positional errors in patrol data for species distribution modeling. *Transactions in GIS*, 22, 202-216.
- Zizka, A., Silvestro, D., Andermann, T., Azevedo, J., Duarte Ritter, C., Edler, D., ... & Antonelli, A. (2019) CoordinateCleaner: Standardized cleaning of occurrence records from biological collection databases. *Methods in Ecology and Evolution*, 10, 744-751.
- Zizka, A., Antonelli, A., & Silvestro, D. (2021) Sampbias, a method for quantifying geographic sampling biases in species distribution data. *Ecography*, 44, 25-32.
- Zurell, D., Franklin, J., König, C., Bouchet, P. J., Dormann, C. F., Elith, J., ... & Merow, C. (2020) A standard protocol for reporting species distribution models. *Ecography*, 43, 1261-1277.