



HAL
open science

Bayesian modelling of species-environment relationships for partially observed data.

Bastien Mourguiart

► **To cite this version:**

Bastien Mourguiart. Bayesian modelling of species-environment relationships for partially observed data.. General Mathematics [math.GM]. Université de Pau et des Pays de l'Adour, 2022. English. NNT : 2022PAUU3023 . tel-04085564

HAL Id: tel-04085564

<https://theses.hal.science/tel-04085564>

Submitted on 29 Apr 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



UNIVERSITÉ DE PAU ET DES PAYS DE L'ADOUR
ÉCOLE DOCTORALE 211

Laboratoire de Mathématiques et de leurs Applications de Pau
(LMAP)

**Bayesian modelling of species-environment
relationships for partially observed data**

Bastien Mourguiart

A thesis submitted for the degree of
Doctor of Philosophy in Mathematics

Thesis defense planned for the 10th November 2022 in front of the thesis
committee composed by:

Benoit Lique	Professor	UPPA	Supervisor
Kerrie Mengersen	Distinguished Professor	QUT	Supervisor
Marie-Pierre Etienne	Assistant Professor	IRMAR	Reviewer
Olivier Gimenez	Director of Research	CNRS	Reviewer
Gurutzeta Guillera-Arroita	Research scientist	CSIC	Examiner
Etienne Prevost	Director of Research	INRAE	Examiner

All we have to decide is what to do with the time that is given us.

Gandalf

Acknowledgments

First, I would like to acknowledge my two supervisors, **Benoit Liquet** and **Kerrie Mengersen**, for giving me their trust and allowing me to pursue my research questions.

Benoit, I am very grateful to have been lucky enough to meet you six years ago and learn from you all this time. Working with you has always been a pleasure, and I hope it will continue. It was also nice to discuss sports and, above all, injuries with you ;).

Kerrie, I wish that the circumstances were different and that our collaboration had not been 100% remote. Anyway, I wanted to thank you for the time that you addressed me through our constructive weekly discussions.

Second, I want to thank all the collaborators who significantly helped me in this work.

Thank the Orthoptera team, especially **Aurélien Besnard** and **Thibaut Couturier**. You always have been reactive, attentive and resourceful. I am much obliged to you for that.

I am also grateful to the “team Palourde”, especially to **Nathalie Caill-Milly**, **Mathieu Chevalier** and **Martin Marzloff**. Thank you **Nathalie** for your kindness, availability and our discussions around the interesting topic of Manila clam’s ecology. You have always offered me your help, and I deeply appreciate it.

Thanks to **Martin** and **Mathieu** for your constructive advice through the Manila clam project and your humanity. I am looking forward to continuing working and spending time with you in the following months.

Even if I spent a lot of time complaining about the time I spent alone during this thesis, I’ve also met and spent time with great people I would like to thank.

Thanks to regular participants of the coffee/tea breaks at Ifremer Anglet. It has always been a pleasure to spend time with you, **Claire**, **Florence**, **Gilles**, **Marie-Noëlle**, **Maud**, **Muriel**, **Nathalie** and **Yann**.

I wish to express all my gratitude to the LEBCO team. Huge thanks for your welcome, the four months spent with you have been a breath of fresh air for me. A special thanks to the open space members: **Lou**, **Lyndsay**, **Léa**, **Alex**, **Laure**, **Clément**, **Mathisse** and **Marion**. It was really nice to meet you!

Naturally, I am grateful to the PhD students of the ‘Team Benoit’: **Teo**, **Sebastien**, **Aurélien**, and **Floren** for our discussions and mutual cheering up during difficult moments. A special thanks to **Claire Kermorvant**, who has been my only colleague most of the time. Thank you for your help and encouragement all along this work.

Last but not least, huge thanks to all my relatives who have always been there for me.

Je remercie toute ma famille : oncles, tantes, cousins, cousine, papi, mamie, et tout particulièrement mes parents, mon frère et mes sœurs, pour leur amour inconditionnel sans lequel je ne serais pas la personne que je suis. Une énorme pensée pour toi maminotte qui m’a accompagné et continue de le faire malgré ton départ.

Un immense merci également à ma seconde famille, mes amis. Si tout se passe bien et que je vous force à m'appeler docteur Bastoun, sachez que vous n'y êtes pas pour rien. Ça a été crucial de pouvoir compter sur vous pour un "p'tit" verre dans les moments difficiles.

Popi, je ne pourrais jamais te remercier assez pour tout ce que tu m'apportes, mais pour commencer, merci pour tout ce que tu as fait pour moi pendant cette thèse. Elle n'aurait pu aboutir si à maintes reprises, tu n'avais pas su trouver les mots pour m'encourager, me faire relativiser et me détendre. **Milesker Popi !!**.

Abstract

Characterising how a species responds to its environment is of central interest in ecology. Species-environment relationships (SERs) are studied, for instance, in community ecology, for species distribution modelling, and to guide conservation or management actions. Statistical models that link species distribution data (e.g., presence/absence or counts) to environmental data (e.g., temperature) are often used to estimate SERs. Standard statistical models assume that the data are representative of the SER. However, in many cases the available data represent only a partial description of the SER. In this work, we investigated the effects on modelling SERs of three kinds of partially observed data:

1. Partially observed response data, e.g., species sampled occurrences may only represent a partial observation of the occupancy status due to missing species present (i.e., imperfect detection).
2. Partially observed environmental data, e.g., environmental descriptors may represent averaged conditions at a coarser spatial scale than the one at which the SER is studied (i.e., area-to-point spatial misalignment).
3. Partially observed relationship, e.g., the gradient of environmental conditions that describe the SER are not entirely surveyed (i.e., truncated gradient).

Hierarchical Bayesian models, allowing multi-species inferences and disentangling ecological from observational processes, have been developed and tested in three case studies, each involving a particular type of partially observed data. In the first case study, we emphasized that even a robust sampling design that involves multiple sampling replicates and detection techniques can lead to species detection probabilities lower than one in an insect community. We then advocated for the use of Multi-Species Occupancy Models to account for imperfect detection in insect studies. In the second case study, we showed how using area-to-point misaligned covariate can flatten SERs estimated by generalized linear models and how fitting a Berkson error model can lower the bias. In the third case study, we developed a hierarchical model that explicitly estimates optimum shifts. By constraining estimated SERs to concave shapes (following ecological theory), the new model improved estimates relative to past methods, especially in the case of truncated gradients. The methods and insights developed in these case studies contribute new knowledge to both the statistical and ecological research communities. They can also be used to inform ecological practice.

Key-words: species-environment relationships, hierarchical Bayesian model, partially observed data, gradient truncation, imperfect detection, spatial scales, spatial misalignment

Résumé

Décrire comment les organismes sont affectés par l'environnement des habitats qu'ils occupent est un des principaux sujets d'étude de l'écologie. Par exemple, l'étude des relations espèces-environnements permet une meilleure compréhension de la structuration des communautés, de la répartition spatiotemporelle des espèces et des effets des changements globaux sur la biodiversité. Les modèles statistiques sont souvent utilisés pour estimer les relations espèces-environnements à partir de données décrivant la réponse d'une espèce, par exemple des données de présence/absence ou d'abondance, le long de gradients environnementaux, comme des gradients de températures. La qualité de l'estimation repose néanmoins sur l'hypothèse de représentativité de l'échantillonnage, les données doivent décrire correctement la relation espèce-environnement étudié. Cependant, du fait de la difficulté d'échantillonner la totalité de la relation, celle-ci est souvent partiellement observée.

Dans ce travail, nous avons étudié les effets de trois types de données partiellement observées sur l'estimation des relations espèces-environnements:

1. La réponse de l'espèce est partiellement observée en cas de détection imparfaite: les données décrivant la présence d'une espèce représentent seulement une portion des vraies occurrences du fait de potentielles non-détections.
2. Les données environnementales ne sont pas spatialement alignées avec les données de réponse: l'environnement décrit représente les conditions moyennes sur une surface plus grande que celle où la réponse de l'espèce a été échantillonnée, et non les conditions locales responsables de la réponse observée.
3. La relation espèce-environnement est partiellement observée: le gradient de conditions environnementales décrivant la relation n'est pas entièrement échantillonné, la relation observée est tronquée.

Des modèles hiérarchiques Bayésiens, qui permettent la modélisation de plusieurs espèces et de distinguer les processus écologiques des processus d'observations, ont été développés pour étudier les relations espèces-environnements de trois cas d'études, chacun étant confronté aux problèmes engendrés par un des trois types de données partiellement observées. Le premier cas d'étude nous a permis de mettre en évidence la nécessité de prendre en compte d'éventuels problèmes de détection lors de l'étude des relations espèces-environnements au sein de communautés d'insectes, via l'utilisation de "modèle d'occupation multi-espèces", même lorsque les données sont récoltées en suivant un protocole maximisant la détection des espèces. Le second cas d'étude a permis de caractériser les biais que pouvait engendrer un non-alignement spatial des données sur les

estimations de relations espèces-environnements. Un modèle linéaire hiérarchique pouvant limiter ces biais a également été présenté. Le dernier cas d'étude présente un nouveau modèle permettant d'estimer explicitement les déplacements d'optimums le long de gradients environnementaux, y compris lorsque ces derniers sont partiellement observés.

Mots-clefs: relations espèces-environnement, modèles hiérarchiques Bayésien, données partiellement observées, réponse tronquée, détection imparfaite, échelle spatiale, non-alignement spatial

Contents

Acknowledgments

Abstract	i
Résumé	ii
Acronyms	v
1 Introduction	1
1.1 What are species-environment relationships?	1
1.2 Why are species-environment relationships studied?	1
1.3 How to describe species-environment relationships?	2
1.4 Partially observed species-environment relationships	3
1.4.1 Imperfect detection	4
1.4.2 Area-to-point spatial misalignment	5
1.4.3 Truncated relationships	6
1.5 Hierarchical Bayesian Modelling of SERs	7
1.6 Organisation of the Thesis	9
2 Chapter 1: Multi-species occupancy models: an effective and flexible framework for studies of insect communities	11
2.1 Synopsis	11
2.2 Publication	11
2.3 Conclusion	41
3 Chapter 2: Modelling species-environment relationships using coarse scale and spatially misaligned environmental data	42
3.1 Synopsis	42
3.2 Publication	43
3.3 Conclusion	89
4 Chapter 3: A new method to explicitly estimate the shift of optimum along gradients in multispecies studies	90
4.1 Synopsis	90
4.2 Publication	91
4.3 Conclusions	139
5 Discussion	140
5.1 Contribution	141
5.2 Perspectives	142
5.2.1 Making more of the Bayesian framework	142
5.2.2 Effects of species traits on SERs	143
5.2.3 Room for more complexity	144
Bibliography	146

Acronyms

BEM Berkson error model.

EHMOS Explicit hierarchical model of optimum shifts.

GAM Generalised additive model.

GLM Generalised linear model.

GLMM Generalised linear mixed model.

HBM Hierarchical Bayesian model.

HM Hierarchical model.

HOF Huisman-Olf-Fresco.

MCMC Markov chain Monte Carlo.

MSOM Multi-species occupancy model.

SDM Species distribution model.

SER Species-environment relationship.

SRC Species response curve.

1 Introduction

1.1 What are species-environment relationships?

Any living organism, whether a bacteria, a plant, or an animal, can survive only a finite range of environmental conditions due to physiological constraints. At the species level, this means that a species can occupy only a limited number of habitats in which environmental conditions fall in the species' physiological tolerance range. Within this set of environmental conditions, often called an environmental or ecological niche (Hutchinson, 1957; McNerny & Etienne, 2012), certain conditions might better suit the species than others. Hence, the species' response (e.g., presence/absence or the number of individuals) to the environment will vary along the gradient of conditions. How it varies is referred to as the Species-Environment Relationship (SER). When one environmental variable is studied, a curve, called a Species Response Curve (SRC), can represent the SER. Ecologists often assume that SRCs can take five different shapes depending on the SER studied (Figure 1); Austin (2002); Oksanen & Minchin (2002); Huisman *et al.* (1993)): 1) a flat curve, i.e., the environment does not influence species presence; 2) a sigmoidal curve, i.e., a monotone increase (or decrease) in species occupancy probabilities along the environmental gradient; 3) a sigmoidal curve with a plateau; 4) a symmetric unimodal curve, i.e., bell-shaped curve with a peak and occupancy probabilities decreasing at the same rate at both side of the peak; and 5) an asymmetric unimodal curve, i.e., occupancy probabilities decreasing at different rates depending on the side of the optimum. In addition, bimodal responses can occur in particular conditions (e.g., exclusion of the species on one part of its physiological tolerance range by competitors) but are seldom considered (see Jansen & Oksanen, 2013 for a counter example).

1.2 Why are species-environment relationships studied?

SERs give essential knowledge of the ecological characteristics of species. Describing a particular SER provides information about the range of conditions that the species can occupy (i.e., the environmental range/tolerance/breadth/width; Heegaard (2002); Hernandez *et al.* (2006); Clavel *et al.* (2011)). Species environmental tolerance can be used to assess the sensitivity of species to some disturbance gradient (Russell *et al.*, 2009) or to expected changes in environmental conditions (e.g., changes induced by climate change, Bellard *et al.*, 2012), and then guide conservation actions (Watson *et al.*, 2012). Furthermore, it can be related to the ecological specialization type (i.e., generalist or specialist species), which can inform about community structure and functioning (Devictor *et al.*, 2008). SERs also form the foundation of species distribution models (SDMs) that have been used for many purposes, including guiding conservation actions (Guisan *et al.*, 2013; Zurell *et al.*, 2021), assessing biological invasion risks (Elith *et al.*, 2010) and evaluating global

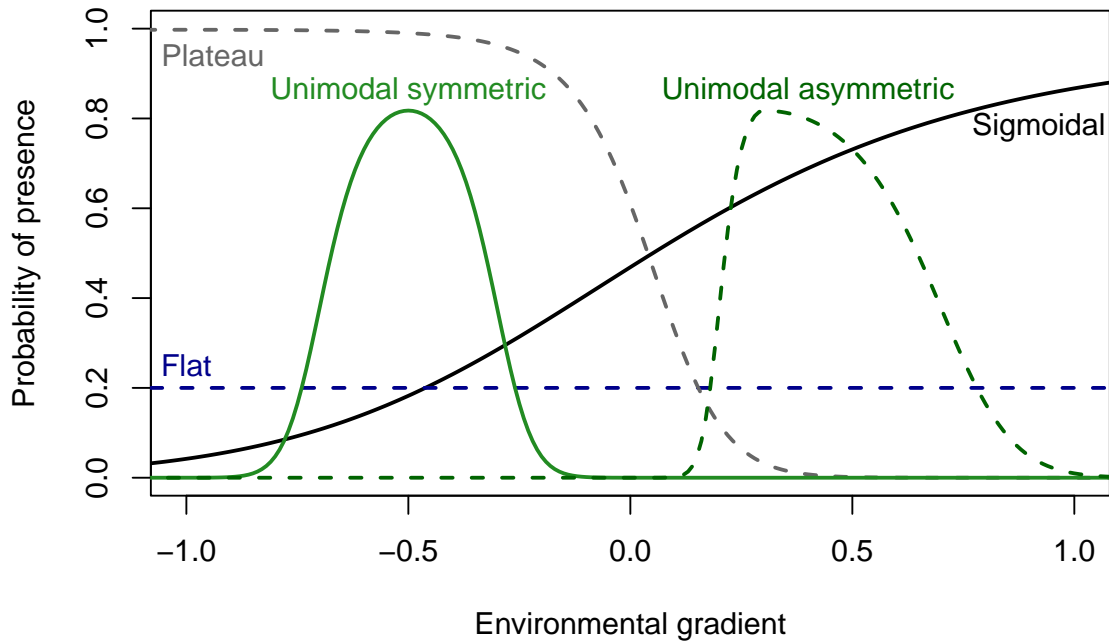


Figure 1: Examples of the five potential shapes of species response to environmental gradient mentioned in the text.

change impacts (Guillera-Arroita, 2017). In addition to indirectly leading to the assessment of global change impacts through species distribution modelling, the study of SERs and their changes along geographical (e.g., latitudinal or elevational) gradients has allowed direct investigations into climate change-induced range shifts (Shoo *et al.*, 2006; Lenoir *et al.*, 2008). Hence, the study of SERs can improve our knowledge of global change effects on biodiversity, which is one of the principal challenges tackled by modern ecology (Bellard *et al.*, 2012; Mouquet *et al.*, 2015; Pecl *et al.*, 2017).

1.3 How to describe species-environment relationships?

As in many ecological studies, the description of SER is based on a sample of data from the studied system. A statistical model is then used to fit the sampled data and make inference about the SER. However, accuracy of inference is conditioned by the quality of the sample. Thus, proper sampling designs have to be developed (Yoccoz *et al.*, 2001; Kermorvant *et al.*, 2019). Ideally, the design should address three questions: “why?”, “what?”, and “how?” (Yoccoz *et al.*, 2001). When studying SER, the answer to “why?” should be to describe how the species respond to the environmental gradient. Thus the aim is define the species environmental range (i.e., the environmental limits distinguishing suitable from unsuitable conditions) and how the species response fluctuates within this range (i.e., the shape of the SRC). The “what?” refers to the type of data to collect to describe both response and explanatory (i.e., environmental) variables. The latter is directly induced by the SER under study. The former depends on the ecological quantity

of interest (Guillera-Arroita *et al.*, 2015). Here, we consider studies in which the response is quantified by presence probabilities or abundance. Thus, presence/absence or count data should be recorded (see Guillera-Arroita *et al.*, 2015 for an explanation on why presence-only data should be avoided when studying probability of species presence). Finally, the “how?” raises many questions, for instance relative to statistical assumptions about sample properties (Kermorvant *et al.*, 2019). These questions can not be all presented here but, specifically for SER studies particular attention should be paid to define: 1) control for detection issues (Yoccoz *et al.*, 2001; Guillera-Arroita *et al.*, 2010; Banks-Leite *et al.*, 2014); 2) the spatial scale at which to describe the SER and thus at which resolution to measure the environment (Dormann, 2007; Connor *et al.*, 2018); and, 3) the length of the environmental gradient to sample (Thuiller *et al.*, 2004; Faurby & Araújo, 2018; Chevalier *et al.*, 2021). Hence, an ideal sampling design should lead to the sampling of multiple spatial units distributed along a gradient of environmental conditions, long enough to describe both suitable and unsuitable conditions, in which species presence/absence or counts are collected without error (perfect detection).

1.4 Partially observed species-environment relationships

In the previous section it was argued that an ideal survey should be designed to answer “why?”, “what?” and “how?”. Failing to answer those questions can lead to unrepresentative data (Yoccoz *et al.*, 2001), thus bias in inference, even if the model is well-specified. However, answering those questions is not straightforward in SER studies, especially the “how?” which requires extensive knowledge about the system. We present below three situations that can prevent answering this question:

- species detection patterns are unknown, i.e., the probability of detecting a species when it is present can vary with unknown external factors, which prevents controlling for detection issues (Guillera-Arroita *et al.*, 2014);
- measuring the environment at each sampling location is impossible, constraining to use of existing environmental data, e.g., when interests lie in the effect of climate on species response (Lembrechts *et al.*, 2019);
- the entire range of environmental conditions that species can survive is not available for sampling, e.g., a species is adapted to conditions that are unavailable on earth at the time of sampling (Faurby & Araújo, 2018).

Each of these situations can lead to what we call partially observed data. Imperfect detection can arise from the first situation, which leads to false-absence and thus partially observed presence. Relying on existing environmental data can force the modeller to describe the environment by

average conditions available at a coarser resolution than the ecological resolution, leading to partially observed environmental data. Finally, not sampling the entire range of environmental conditions can lead to a partially observed relationship.

1.4.1 Imperfect detection

Decades of research on species detectability Devarajan *et al.* (2020) have shown that individuals and species are seldom perfectly detected. Imperfect detection leads to false absences (i.e., incorrectly considering a species absent when it is present but undetected; Tyre *et al.* (2003)). In such a situation, the partially observed presence/absence data no longer allow inference on the presence probability but instead on the probability of observing the species, i.e., the joint probability of presence and detection events (Guillera-Arroita *et al.*, 2015). This confusion can lead to significant bias in SER estimates (Tyre *et al.*, 2003; Lahoz-Monfort *et al.*, 2014), especially if species detectability varies along the environmental gradient in different ways than the occupancy (Lahoz-Monfort *et al.*, 2014).

In the early 2000s, site-occupancy models were developed to deal with imperfect detection of species (MacKenzie *et al.*, 2002; Tyre *et al.*, 2003) and the number of species individuals (i.e., N-mixture models; Royle (2004)). These models distinguish the ecological process (i.e., species occupancy or abundance at a site) from the observational process (i.e., species or individual detection given its presence), allowing for simultaneous estimates of species occupancy (or abundance) and detection probabilities. They can be used to model the effects of environmental covariates on occupancy probabilities (i.e., estimate SER) while accounting for detection issues. Estimating detection probabilities requires replicated sampling at the surveyed locations (Bailey *et al.*, 2007). Replication can involve multiple visits, observers, observation techniques, or spatial replicates (Guillera-Arroita, 2017).

Despite the quantity of literature on imperfect detection issues and the development of models to deal with them, imperfect detection often continues to be overlooked in SER studies (Kellner & Swihart, 2014; Devarajan *et al.*, 2020). Furthermore, geographical and taxonomic biases have been observed in the use of site-occupancy models (Kellner & Swihart, 2014; Devarajan *et al.*, 2020). Several explanations exist. One is the belief that a robust sampling design can prevent imperfect detection issues (Welsh *et al.*, 2013), which can be true if detection patterns are well-known. However, such knowledge is often lacking (Kellner & Swihart, 2014), precluding the control of unknown external factors that could influence detectability (Guillera-Arroita *et al.*, 2014). Moreover, the extra sampling effort required for sampling replication can be perceived as not worthwhile (Welsh *et al.*, 2013), preventing the use of site-occupancy models. Finally, a lack of studies in specific domains (e.g., entomology or marine environment; Kellner & Swihart (2014);

Devarajan *et al.* (2020)) can also constrain the use of site-occupancy models. Thus, broadly promote the use and advantages of site-occupancy models may improve the study of SERs.

1.4.2 Area-to-point spatial misalignment

When studying a SER, ideally, environmental data should be collected at the spatial resolution at which it acts on the species response (i.e., at its scale of effect; Chandler & Hepinstall-Cymerman (2016)). Indeed, using environmental covariates at scales that differ from the scale of effect can lead to bias in SER estimates (Connor *et al.*, 2018). However, at least two common situations can prevent the sampling of environmental covariates at the appropriate scale:

- The scale of effect is not known beforehand. The presence of a species at specific locations can be explained by the environmental conditions at those locations, but also by particular features of the surrounding landscape (De Knecht *et al.*, 2010; Chandler & Hepinstall-Cymerman, 2016).
- Depending on the environment of interest, it can be infeasible to collect environmental data at all sampling units. For instance, studying effects of temperature on species occupancy implies recording temperature data during months or years. Indeed, temperatures collected during species sampling alone will not represent the range of conditions that the species actually experienced. Deployment of sensors at all sampling units seem unreasonable.

Thus, many studies have used existing/available environmental data to describe species-environment relationships.

Many environmental data used to estimate SERs come from outputs of numerical global climate models (e.g., WorldClim, Hijmans *et al.*, 2005). In these situations, the spatial resolution of environmental covariates is imposed by the resolution of the model outputs. This often leads to a scale (i.e., a resolution) mismatch between environmental and species response data (Potter *et al.*, 2013). A particularly common situation, referred as *area-to-point spatial misalignment* in the geostatistics literature (Gotway & Young, 2002), occurs when environmental conditions are described at coarser resolutions than the species response data (Latimer *et al.*, 2006; McInerney & Purves, 2011; Potter *et al.*, 2013). Area-to-point spatial misalignment is known to induce bias in regression estimates (Gotway & Young, 2002). However, for a long time, ecologists did not consider area-to-point spatial misalignment as a problem in SER studies (Potter *et al.*, 2013). They assumed that large-scale climate was the main driver of species response (Pearson & Dawson, 2003) and thus, that coarsely-resolved climate data represented the scale of effect. While this might be true when SER is studied at large-scale or for mobile species (Guisan *et al.*, 2007), it can be otherwise (Meineri & Hylander, 2017; Lembrechts *et al.*, 2019). For instance, micro-climate

variability is known to play an important role for fine-scale species distribution, especially in forest ecosystems (Zellweger *et al.*, 2020) or for sessile species (Chauvier *et al.*, 2022). Thus, fine-scale environment might only be partially described by coarse-scale environment which represents averaged conditions less variable than the true conditions (McInerny & Purves, 2011).

Partial description of fine-scale environmental variability induced by area-to-point spatial misalignment is also recognized to bias SRC (Latimer *et al.*, 2006; McInerny & Purves, 2011; Martínez-Minaya *et al.*, 2018). However, there have been only a few attempts to address this problem (Latimer *et al.*, 2006; McInerny & Purves, 2011). For instance, McInerny & Purves (2011) proposed to use a Berkson measurement-error model (BEM), in which the observed covariate w is considered as an error-prone representation of an unobserved error-free variable x . The observed covariate is a less variable version of the error-free variable: $x = w + u$, with u the amount of error (i.e., the variance error). McInerny & Purves (2011) showed through simulations that BEM improve estimates of SRC compared to a model omitting the error, but assuming a known variance error. This information is, however, unavailable in most studies. Latimer *et al.* (2006) advocated for the use of spatial generalized linear models to account for area-to-point misalignment but without providing any external evidence of accuracy or predictive power of this method. Thus, no “general” method is currently used to address or even detect area-to-point misalignment.

1.4.3 Truncated relationships

In some studies, sampling design is defined geographically and not environmentally (e.g., in SDM). Instead of capturing the entire gradient of values that a species could occupy, environmental gradients can be only partially observed (Thuiller *et al.*, 2004) or partially available during the survey (Faurby & Araújo, 2018). Such situations, lead to a partial observation of the SER. Such a partially observed relationship is also referred as truncated niche (Chevalier *et al.*, 2021) or truncated response (Austin, 2007).

Partially observed SER can lead to severe bias in SRC estimates (Thuiller *et al.*, 2004; Citores *et al.*, 2020; Chevalier *et al.*, 2021). For instance, simulation studies have shown that regression methods can fit a sigmoidal curve instead of a unimodal curve when optimum is near the end of the observed gradient (e.g., generalised linear models (GLMs): ter Braak & Looman, 1986; Coudun & Gégout, 2006; generalised additive models (GAMs): Citores *et al.*, 2020; Huisman-Olff-Fresco (HOF) models: Jansen & Oksanen, 2013). Such bias can lead to inaccuracies in predictions outside the sampled gradient (Thuiller *et al.*, 2004). In addition, optimum estimates can not be derived from the SRCs in those situations (ter Braak & Looman, 1986; Coudun & Gégout, 2006). This is particularly prejudicial when optimum is the ecological parameter of interest, e.g.,

in climate change studies in which optimum shifts are estimated (Lenoir *et al.*, 2008).

In optimum shift modelling, accurately estimate optimum shifts for edge species is crucial as such species can be close to their environmental limits and thus particularly sensitive to change (Freeman *et al.*, 2018). However, the current most two widely used methods, GLM-based (e.g., Lenoir *et al.*, 2008) and mean comparison based approaches (e.g., Shoo *et al.*, 2006), can not accurately estimate optima near gradient edges (ter Braak & Looman, 1986). Hence, a method allowing accurate optimum shift estimates for edge species can enhance studies of SRC optimum shifts.

1.5 Hierarchical Bayesian Modelling of SERs

GLM is a very popular method for studying SERs (Guisan *et al.*, 2002; Austin, 2007). GLMs are commonly used to link presence/absence or count data to a linear combination of predictors by means of a link-function (Guisan *et al.*, 2002). Usually, for the study of SER, GLMs include a combination of linear and quadratic terms to model unimodal symmetrical response curves expected to occur for many SERs (Austin, 2007). Sometimes, asymmetric responses are expected, and more flexible approaches (e.g., GAMs or HOF models) are preferred (Oksanen & Minchin, 2002). However, flexibility comes with a loss of simplicity in interpretation and high-quality data requirements (Merow *et al.*, 2014). As in SER studies, the main goal is often ecological interpretation and that partially observed data might lead to low-quality data, so we chose to consider only GLM-based approaches in this work.

GLMs can be too simple and unsuitable for most ecological data's complex structure. As mentioned earlier, in the case of partially observed SER, data collected represent a tangle between ecological and observational processes rather than the description of the SER (e.g., in case of imperfect detection, Tyre *et al.*, 2003). In such situations, GLMs are expected to produce biased estimates of the SER (Lahoz-Monfort *et al.*, 2014) (see also a simulated example in Figure 2). However, they can be extended into Hierarchical Models (HMs; also called mixed models (Bolker *et al.*, 2009) or multi-level models (McElreath, 2016)) that allow the decomposition of complex systems into hierarchical sequences of simpler models (Wikle, 2003). Such decomposition allows disentangling ecological processes from observational processes (Cressie *et al.*, 2009), making HMs valuable to study SER with partially observed data (Royle & Dorazio, 2008; Hefley & B. Hooten, 2016). In addition, HMs provide a flexible approach to modelling multiple SERs simultaneously (Ovaskainen *et al.*, 2017; Poggiato *et al.*, 2021). Multi-species HMs assume that species within a community respond similarly but not equally to the environment (Bolker *et al.*, 2009), i.e., specific SERs are described by a combination of shared community-level parameters

and species-level variability (McElreath, 2016; Pedersen *et al.*, 2019). Hence, data information is shared among species which can improve SER estimates, especially for data-poor species (Ovaskainen & Soininen, 2011).

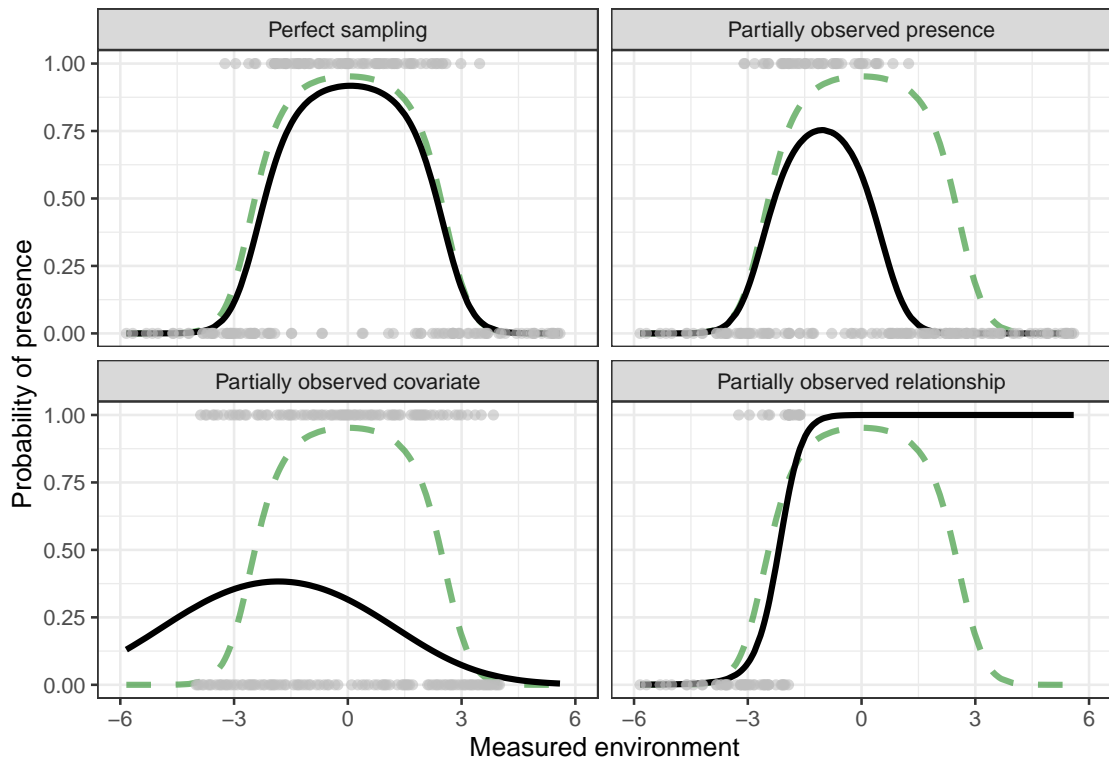


Figure 2: Effects of partially observed data on modelling of species-environment relationships. In all panels, the dashed green line shows the true relationship between the probability of presence and the environment. Black lines represent the relationships predicted by fitting a generalized linear model on different types of observed data: a) fully observed data, i.e., data come from an error-free sampling; b) partially observed response data, i.e., species imperfectly detected and detection probability decreases across the gradient; c) partially observed environmental data, i.e., measured environment describes average conditions at a coarser resolution than conditions the species experience; d) partially observed relationship, i.e., environmental gradient not entirely sample.

Hierarchical models can be fitted with either Bayesian or frequentist approaches (Cressie *et al.*, 2009). There are pros and cons associated with both paradigms (Lele & Dennis, 2009; Dorazio, 2016). For example, among the cons against Bayesian analysis, there are : computational costs of Markov chain Monte Carlo (MCMC) methods used for estimation, and the subjectivity induced by the choice of priors (Lele & Dennis, 2009). Among the pros for Bayesian analysis, there are: clear and intuitive estimates of modelling uncertainties (Ellison, 2004) and applicability for data with low sample size (Dorazio, 2016). The elements that led to my choice to adopt Bayesian methods are subjective and linked to: my conception of modelling and how I have been taught statistics as an ecologist. Bayesian inference provides a probabilistic measure of a model being true given the data observed ($P(model|data)$). This better fits my point of view that a model is always a simplified and subjective representation of the reality chosen by a modeller.

Model selection methods may seem objective, however, they still depend on the models and on the comparison techniques chosen by the modeller. Then, it seems more appropriate to have a probability statement about the degree of support for a model given the observed data, rather than how likely it is to observe the data given the model. Besides, in my experience, ecologists are trained to use frequentist methods by the prism of specific software and the functions within. Thus, an ecologist can be constrained by the models available in the software she/he is using. In contrast, once an ecologist has acquired a little background in Bayesian statistics, MCMC samplers and flexible software such as WinBUGS (Lunn *et al.*, 2000) or JAGS (Plummer, 2003), which is feasible thanks to some great books (Royle & Dorazio, 2008; e.g., Kéry, 2010; Kéry & Royle, 2015, 2021), she/he can build any sort of hierarchical model.

1.6 Organisation of the Thesis

This thesis investigates potential bias and related issues that can arise when using partially observed data to estimate species-environment relationships (SERs). Specifically, we focus on particular case studies involving three types of partially observed data:

- 1) partially observed response data produced by imperfect detection of species;
- 2) partially observed covariate data related to a specific case of spatial misalignment;
- 3) partially observed relationship in the particular context of optimum shift modelling.

We are also interested in potential solutions that hierarchical Bayesian models (HBMs) can offer in the presence of such partially observed data.

This work is organised in 3 chapters, each dedicated to a specific type of partially observed data that may lead to bias in SER modelling.

The first chapter of this thesis highlights the need to account for imperfect detection in the study of insect communities. While it is known that imperfect detection often occurs during field surveys, it is widely overlooked by entomologists. Through a case study on an Orthoptera communities, we show that even with a robust sampling design specific detection probabilities are likely to be less than one. Thus, we advocate for using an HBM, called multi-species occupancy model (MSOM), that accounts for imperfect detection by estimating both species occupancy and detection probabilities. In addition, we test a method to investigate sampling efficiency and the corresponding effects of potential optimisation to overcome reservations about the relatively high sampling effort required for the use of MSOM. This chapter is composed of an article published in *Ecological Entomology*.

The second chapter of this work investigates the effects of spatial misalignment on models'

performance using covariates averaged across scales that are coarser than the scale of the ecological process. We compare the accuracy of three models (a GLM, a spatial GLM and a Berkson measurement Error Model (BEM)) through simulations and application to a real case study of the Manila clam in Arcachon Bay. We find that the GLM and the spatial GLM produce biased SER estimates, with bias increasing with the coarsening of the environmental data. Thus, we advocate for not interpreting GLM estimates if environmental data are available at a resolution that may be coarser than the species response scale. In contrast, the BEM give promising results, with almost no bias in SER estimates even for very coarse environmental data, but needs further investigations for broad applications. This chapter is composed of an article in preparation for submission in *Journal of Applied Ecology*.

In the third chapter, we propose a new ecological formulation of a hierarchical Bayesian model to explicitly estimate optimum shifts along environmental gradients for multiple species. We demonstrate through simulations that this model, called Explicit Hierarchical Model of Optimum Shifts (EHMOS), is more accurate than a GLM-based approach and the mean comparison method (two widely used methods). Furthermore, the ecological formulation of EHMOS allows accurate optimum shift estimates for edge species (i.e., species having a partially observed SER) in opposition to the mean comparison method and the GLM-based approach. In addition, EHMOS has better accuracy under unbalanced sampling design relative to the mean comparison method. Finally, we discuss further developments of the method to investigate ecological mechanisms underlying observed optimum shift variability among species. This chapter is composed of an article under review in *Journal of Biogeography*.

2 Chapter 1: Multi-species occupancy models: an effective and flexible framework for studies of insect communities

The core of this chapter is a paper published in *Ecological Entomology*: Mourguiart, B., Couturier, T., Braud, Y., Mansons, J., Combrisson, D., & Besnard, A. (2021). Multi-species occupancy models: an effective and flexible framework for studies of insect communities. *Ecological Entomology*, 46(2), 163-174. I also made a Shiny application available online (https://bastien-mourguiart.shinyapps.io/shiny_MSOM/) to facilitate species-specific results dissemination, especially for entomologists with low statistical knowledge.

2.1 Synopsis

Studies of insect communities often aim to estimate species distributions, community composition, or species-richness patterns. False absences (i.e., noted an absence while the species was present but missed due to imperfect detection) can, however, bias estimates of models that do not account for imperfect detection (Tyre *et al.*, 2003; Lahoz-Monfort *et al.*, 2014; Tingley *et al.*, 2020). Multi-species occupancy models (MSOMs) seem to afford a flexible solution to cover the main topics of ecological entomology while dealing with imperfect detection issues (Mata *et al.*, 2014; Tingley *et al.*, 2020). However, MSOMs are still rarely used in insect studies (Devarajan *et al.*, 2020), and imperfect detection issues are often overlooked (Kellner & Swihart, 2014).

MSOMs require specific sampling designs to distinguish the ecological process (i.e., species occupancy) from the observational process (i.e., species detection). Sampling has to occur at multiple sites to estimate occupancy probabilities and sampling on each site has to be replicated to estimate detection probabilities (Guillera-Arroita, 2017). Replication could be conducted by visiting the site during multiple visits, at different locations inside it, by multiple observers or by using multiple detection techniques (Guillera-Arroita, 2017). This sampling replication can be costly, which could explain the little use of MSOMs in insect studies.

In this chapter, we developed an MSOM to estimate species-specific occupancy and detection probabilities of Orthoptera species. In addition, we used the MSOM to estimate species richness and inventory completeness. The chapter's goals were to highlight (1) the need to account for imperfect detection in insect studies, even for easily detectable species such as Orthoptera sampled following a robust sampling design, and (2) the use of MSOM to investigate sampling efficiency and potential sampling optimisation.

2.2 Publication

METHODS

Multi-species occupancy models: an effective and flexible framework for studies of insect communities

BASTIEN MOURGUIART,^{1,2} THIBAUT COUTURIER,²
YOAN BRAUD,³ JÉRÔME MANSONS,⁴ DAMIEN COMBRISSON⁵
and AURÉLIEN BESNARD²

¹CNRS/UNIV PAU PAYS ADOUR/E2S UPPA, Laboratoire de Mathématiques et de leurs Applications de Pau – MIRA, UMR5142, Anglet, France, ²Centre d'Ecologie Evolutive et Fonctionnelle, UMR 5175, Univ Montpellier, CNRS, EPHE-PSL University, Univ Paul Valéry Montpellier 3, Montpellier, France, ³Bureau d'études Entomia, Col de Clans, Vaumeilh, France, ⁴Parc national du Mercantour, Nice, France and ⁵Parc national des Ecrins, domaine de Charance, Gap, France

Abstract. 1. Entomological studies often aim to estimate species distribution, community composition, or species-richness patterns. False absences can, however, bias these estimates and should consequently not be overlooked in insect studies. Multi-species occupancy models (MSOMs) afford a flexible solution to cover the main topics in ecological entomology while dealing with detectability issues.

2. We sampled Orthoptera communities at 81 mountain grasslands sites in France, using three sampling techniques: sighting, listening, and sweep netting. Five plots were sampled per site. This sampling design allowed MSOMs to be used to estimate richness, occupancy, and detection probabilities while accounting for the effect of covariates. We also used MSOMs to evaluate the efficiency of the survey design and to assess the effects of sampling optimisation.

3. The estimates obtained for altitudinal distribution were reliable, with known species distributions confirming the relevance of MSOMs to model the effects of covariates on Orthoptera communities. The species-specific detection probability was often less than one and varied with the detection technique used and the grass height, confirming the need to deal with detection issues in orthopteran studies.

4. We estimated an inventory completeness superior to 0.80 for 93% of the sites, and an overall detection probability superior to 0.95 for 52% of the species, suggesting the sampling design was suitable for studying occupancy in Orthoptera communities. We also found that the sweep netting step may be omitted or the number of plots reduced without affecting species detectability or inventory completeness. Those recommendations may help to optimise future sampling strategies.

Key words. Hierarchical model, imperfect detection, orthoptera communities, sampling optimisation, sampling efficiency, species distribution modelling.

Introduction

The study of species distribution and its determinants is of central interest in theoretical and applied ecology (Rush-ton *et al.*, 2004; Guisan & Thuiller, 2005; Vaughan &

Ormerod, 2005). By acquiring information on species occupancy patterns in a set of sampling locations and/or sampling periods (Guillera-Arroita, 2017), ecologists are, for instance, able to make inferences about the environmental drivers behind species distribution (Guisan & Thuiller, 2005; Zipkin *et al.*, 2009). This information can in turn be used to plan appropriate management actions given conservation aims or to understand species range dynamics (Moritz *et al.*, 2008; Pecchi *et al.*, 2019).

Correspondence: Bastien Mourguiart, CNRS/UNIV PAU PAYS ADOUR/E2S UPPA, Laboratoire de Mathématiques et de leurs Applications – MIRA, UMR5142, 1 allée du parc Montaury, 64600 Anglet, France. E-mail: bastien.mourguiart@etud.univ-pau.fr

Over the last decade, species distribution models (SDMs), relying on the modelling of occurrence data together with environmental covariates, have become the central tool used to study species distribution range or dynamics (Guisan & Thuiller, 2005). Yet most classically used SDMs are based on presence-only data (Guisan & Zimmermann, 2000; Guillera-Aroita *et al.*, 2015). They should be used with caution as their results could be highly biased with sampling effort (Phillips *et al.*, 2009) or with imperfect detection (Guillera-Aroita *et al.*, 2015). Besides, these models allow inference about indices of species occurrence and not about probability of species presence, which is the common purpose of ecologists (Guillera-Aroita *et al.*, 2015). SDMs based on presence-absence data are considered as more flexible, but inferences about probability of species presence depend on the assumptions that the species detectability is almost perfect and remains constant between sites (Guillera-Aroita *et al.*, 2015). Violating those assumptions may induce high bias, especially for species characterised by a limited detectability (Lahoz-Monfort *et al.*, 2014). However, decades of studies on detection issues in ecological monitoring have shown that species probabilities of detection are often less than one and may vary with environmental features (Tyre *et al.*, 2003; Lahoz-Monfort *et al.*, 2014).

Detection issues generate false absences that in turn may result in strong bias when modelling species distribution (Tyre *et al.*, 2003; Lahoz-Monfort *et al.*, 2014). For instance, detection issues can lead to the systematic underestimation of the distribution range (MacKenzie *et al.*, 2002) or can overestimate occupancy turnover (MacKenzie *et al.*, 2003). In the early 2000s, so-called ‘site-occupancy models’ were specifically developed to solve this issue by simultaneously estimating detection and species occupancy probability (MacKenzie *et al.*, 2002; Tyre *et al.*, 2003). Since then, they have been rapidly extended to allow modelling the effects of covariates on occupancy and detection probabilities (MacKenzie *et al.*, 2006), to model range dynamics (Moritz *et al.*, 2008), or to take advantage of information on species states at study locations (Nichols *et al.*, 2007). Initially developed to model occupancy of one focal species, site-occupancy models have been extended to study the occupancy patterns of several species simultaneously (Dorazio & Royle, 2005). These models, called multi-species occupancy models (MSOMs), increase precision in occupancy estimations compared to single-species models, especially for rare species, by borrowing information from data-rich species (Zipkin *et al.*, 2009; Ovaskainen & Soininen, 2011). The hierarchical structure of these models also allows making inferences about the true species richness at study locations, a result that cannot be achieved by species-by-species analysis (Dorazio & Royle, 2005; Guillera-Aroita *et al.*, 2019). MSOMs thus appear to have potential as robust tools for biodiversity analysis or biological assessment (Mata *et al.*, 2014; Devarajan *et al.*, 2020; Tingley *et al.*, 2020). However, despite these models were developed 15 years ago, only 106 published studies have relied on their use, according to a recent review (Devarajan *et al.*, 2020). Among them, MSOMs have mainly been used for vertebrates, and only marginally in other taxa such as plants (see Roth *et al.*, 2018) or insects (Devarajan *et al.*, 2020; but see Mata

et al., 2014; Brodie *et al.*, 2019; Dorazio *et al.*, 2006), even though these taxa may pose specific detection issues because of their ecology (phenology, for instance). Accordingly, the relevance of MSOMs appears to be still overlooked for a large range of taxa.

As has been previously highlighted by other authors (Mata *et al.*, 2014; Brodie *et al.*, 2019), since many entomological studies look for changes in insect communities, it is crucial that they take into account imperfect detection. In this context, MSOMs present some clear advantages. First, this approach makes it possible to estimate the occupancy probability even for data-poor taxa, considering for instance the numerous rare, cryptic, or elusive species in insect communities, often difficult to detect in the field (Coddington *et al.*, 2009; Silva *et al.*, 2019). These taxa are usually excluded in common analyses, such as multivariate methods, in which species found at less than 5% of sites are usually removed (Ter Braak & Smilauer, 2002; Pierik *et al.*, 2017). Secondly, both the species- and community- levels can be simultaneously studied with MSOMs (Mata *et al.*, 2014), thanks to the hierarchical structure of their models. The biodiversity estimators provided by MSOMs, such as species richness, have the advantage of accounting explicitly for the effects of survey-, site-, and species-level covariates that may affect detectability, in contrast with most usual estimators (Tingley *et al.*, 2020). This is of particular interest in entomological studies, since the activity rate and the density of insects—and thus detection probability—are strongly affected by survey and site conditions (Wolda, 1988; Bale *et al.*, 2002).

Orthoptera is an insect group intensively studied in ecology and regularly used as ecological indicator (Marini *et al.*, 2009; Bazelet & Samways, 2011). Detectability issues are expected to occur in orthopteran field surveys because of their small size and the strong variations in abundance related to their phenology (Badenhausser *et al.*, 2009). Detectability is also expected to vary greatly among species (Badenhausser *et al.*, 2009) due to the high diversity in their ecological traits (e.g. mobility, singing activity, mimicry, etc.) and in their habitat preferences (e.g. closed forests vs open grasslands). Orthoptera detection also strongly depends on the sampling techniques used, the effectiveness of which is often influenced by species-specific traits. For instance, highly mobile Orthoptera may flush out when the observer approaches, becoming easily detectable by sight, but hardly detectable using the sweep net or the box quadrat methods. Conversely, sweep netting and box quadrats may increase the detectability of cryptic Orthoptera living close to the ground, often hardly detectable simply by sight. Despite these detection issues, only two studies conducted on single orthopteran species have explicitly dealt with imperfect detection by using site-occupancy models (MacKenzie *et al.*, 2003; Veran *et al.*, 2015). A third research applied site-occupancy models on orthopteran communities, but using single-species site-occupancy models for each species, without using a standardised sampling design (Malinowska *et al.*, 2014).

Site-occupancy models require a specific survey design, usually based on temporal replication conducted on a set of sampling sites (MacKenzie *et al.*, 2006). Replication at the site scale can be obtained through repeated visits, but also by using spatial replicates, multiple observers, or multiple sampling techniques

depending on the study (Guillera-Arroita, 2017). A constraint of this method is that this replication increases the sampling effort and the associated costs, precluding the use of this approach by practitioners, notably when the monitoring budget is limited (Field *et al.*, 2005). On the other hand, this replication is needed to explicitly deal with detection issues and thus to develop robust monitoring (Yoccoz *et al.*, 2001). Hence, there is a crucial need to develop monitoring that optimises the trade-off between sampling effort, techniques, and effectiveness when designing site-occupancy surveys (Field *et al.*, 2005).

Most existing grasshopper-sampling protocols are based on estimates of the abundance index or raw presence–absence (Gardiner *et al.*, 2005). To our knowledge, none were designed to use multispecies site-occupancy models. In this study, we therefore investigated the effectiveness of MSOMs in estimating the occupancy probability of Orthoptera at community level, while accounting for imperfect detection. We also evaluated the ability of MSOMs to assess the efficiency of the survey design under different sampling optimisation scenario.

Materials and methods

Sampling method and design

Field surveys were conducted between 9th August and 5th October 2018 in the Mercantour National Park, located in the southern French Alps. We selected 81 sampling sites among 179 locations already studied between 1983 and 1988 (Gueguen, 1990), in order to encompass all the altitudinal and exposure gradients, ranging from 928 to 2614 m (mean = 1869 m) above the sea level. Each sampling site consisted in a circle (70 m radius) placed in relatively homogeneous grasslands (Fig. 1), and distanced at least 30 m from woodland areas in order to avoid edge effects (Bieringer & Zulka, 2003).

Within each sampling site, we defined five spatial replicates (hereafter ‘plots’), placed on two lines perpendicular to the

mountain slope and spaced at least 30 m from each other, to avoid potential double counts among plots (Fig. 1). In few cases ($N = 5$ sites), this spatial design had to be slightly adapted depending on the sampling site configuration, notably when the grassland area was too small. Thus, some replicates were placed less than 30 m apart or from the forest edge. However, as a minimum distance of 20 m was maintained between plots and from the forest edge, we considered close plots as independent and no effect of the forest proximity on species composition.

Each plot consisted of a 30 m² area, measured by means of a cord, circular or rectangular in shape, depending on ground cover and slope steepness within the sampling site. In particular, circular plots were preferred in sites characterised by gentle slopes and/or shrubby vegetation, while rectangular plots were surveyed on steep slopes and/or in short-sward sites.

Samplings were carried out by a single trained observer (Y.B.), when the weather conditions were optimal for diurnal Orthoptera activity, i.e. no rain, low to moderate wind speed, and sunshine (or temperature exceeding 18°C if cloudy). Species identification was conducted in the field for almost all species, except for *Anonconotus occidentalis* Carron & Wermeille, 2002 and *Anonconotus ligustinus* Galvagni, 2002 which cannot be distinguished without the examination of genitalia morphology. Hence, individuals that could be both *A. occidentalis* and *A. ligustinus* were captured and identified in the laboratory, only *A. occidentalis* was detected in our inventories. In each plot, orthopteran assemblages were surveyed following three successive steps: (1) 1 min of listening to species stridulating in the plot by standing close to its edge, (2) 6 min of sighting species by walking across the entire plot, and (3) two 45-s sweep netting sessions across the entire plot. We chose this execution order for the different sampling techniques because we expected that it would optimise the number of species encountered. Beginning by the listening step leads to record singing species before disturbing them. Then, we expected that walking across the plot will make the mobile species flush away

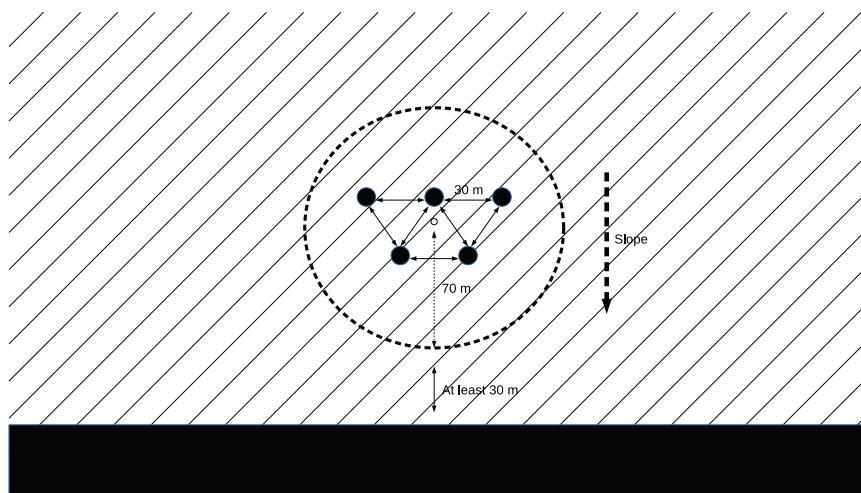


Fig. 1. Schematic representation of the sampling design conducted at each sampling site (dashed circle), placed in relatively homogenous grasslands (hatched area) and spaced at least 30 m from the closest forested habitat (black area). Orthoptera were sampled in five plots (solid circles) within each site using three sampling techniques (sighting, listening, and sweep netting).

making them easily detectable by sight. Finally, sweep netting sessions aimed to capture the remnant less mobile species that did not flush away in step 2. By means of this sampling design, detection/non-detection data were available for each sampling technique in each plot at each site.

Multi-species occupancy modelling

We then modelled the detection/non-detection data of the 15 replicates per sampling site (i.e. 5 plots \times 3 sampling techniques) using site-occupancy models (MacKenzie *et al.*, 2002; Tyre *et al.*, 2003). These models estimate the occupancy probability while modelling imperfect detection through the use of sampling replication conducted at each site. The sampling replication at one site can be achieved with multiple visits, multiple detection methods or multiple observers, as well as with spatial replicates if samplings occur at different locations within a site in a single visit (MacKenzie *et al.*, 2006). In addition, different types of replication are available, depending on the characteristics of the target species, the study area and the objectives (Guillera-Arroita, 2017). Among these options, the spatial replicates method was selected in this study, due to field constraints related to the mountain environment (which required long travel to access the sampling sites). When based on spatial replicates, the model assumes that the occupancy status of sampling sites does not change between site replicates (in our case, plots), i.e. closure assumption (Kendall & White, 2009). Because the plots were closed in space and set up in sites of homogeneous vegetation cover, we considered that this assumption was met. We also assumed no false positives, i.e. only zero can be reported for a species at a site where it is absent, considering the observer's identification skills.

Site-occupancy models disentangle the ecological process, i.e. the true occupancy state for a species in a site, from the observational process, i.e. the detection/non-detection of a species at a site given it is present. In order to distinguish a true absence from a non-detection, we modelled the raw data for species i at replicate k of site j , denoted $X_{i,j,k}$, as the outcome of a Bernoulli random variable, defined by:

$$X_{i,j,k} \sim \text{Bernoulli}(p_{i,j,k} \times Z_{i,j})$$

where $p_{i,j,k}$ is the detection probability of species i at replicate k of site j , and $Z_{i,j}$ is a binary variable corresponding to the occupancy state of site j by species i (latent state). The model for occurrence is specified as:

$$Z_{i,j} \sim \text{Bernoulli}(\psi_{i,j}) \quad (1)$$

where $\psi_{i,j}$ is the probability that species i occurs at site j .

We estimated the occupancy ($\psi_{i,j}$) and the detection ($p_{i,j,k}$) probabilities of all the species using a multi-species site-occupancy model (MSOM) with single-species occupancy models as building blocks (Dorazio & Royle, 2005; Zipkin *et al.*, 2009). The species-specific intercepts and slopes from occupancy and detection models were drawn from a shared, community-level distribution via random effects. Through such hierarchical structure, species with less data

“borrow” information from other species that are data-rich, which improves precision in estimates (Zipkin *et al.*, 2009; Ovaskainen & Soininen, 2011). The linking of species via random effects also allows inferences to be made about the number of species N_j present at each site j , including species never detected (Dorazio *et al.*, 2006; Guillera-Arroita *et al.*, 2019). In this modelisation process, we followed the ‘data augmentation’ method described by Royle *et al.* (2007), also including N_0 hypothetical species to the dataset, all with zero detection.

The species-specific effect on the occupancy and the detection probabilities was integrated into the model using the logit link function. Furthermore, covariates supposed to influence the occupancy and the detection probabilities were also considered (altitude and grass height), and the occurrence probability for species i at site j was modelled by incorporating site-specific characteristics. We developed a relatively simple model including just few covariates in order to show MSOMs potential rather than explaining ecologically Orthoptera distribution or detectability. In particular, linear and quadratic effects of altitude, varying across species, were included to study the altitudinal distribution of Orthoptera. We could assume that grass height affects Orthoptera occurrence, but we did not add this occupancy covariate as the results and the interpretations were virtually identical, and it complexified the model (see ESM S1 for results comparison). The occupancy model was defined as:

$$\text{logit}(\psi_{i,j}) = \alpha_{0_i} + \alpha_{1_i} \times \text{altitude}_j + \alpha_{2_i} \times \text{altitude}_j^2$$

where α_{0_i} is the species-level intercept and $(\alpha_{1_i}, \alpha_{2_i})$ are the species-specific covariate effects.

The detection probability for species i was assumed to vary depending on the detection technique used (sighting, listening or sweep netting). We coded the sampling technique covariates as dummy variables, with the intercept corresponding to the sighting technique. We also added the linear effect of grass height on the probability of detecting species i using the sighting technique. The grass height was standardised to have mean equal to zero:

$$\begin{aligned} \text{logit}(p_{i,j,k}) = & \beta_{0_i} + \beta_{1_i} \times \text{listening}_{j,k} + \beta_{2_i} \times \text{netting}_{j,k} \\ & + \beta_{3_i} \times \text{height}_{j,k} \times \text{sighting}_{j,k} \end{aligned}$$

with β_{0_i} , β_{1_i} and β_{2_i} the species-specific effects of the sampling techniques, and β_{3_i} the species-specific covariate effects.

Each random parameter (α 's and β 's) was modeled as drawn from a normal distribution described by the community mean (μ) and the variance between species (σ^2):

$$\alpha_{0_i} \sim N(\mu_{\alpha_0}, \sigma_{\alpha_0}^2), \alpha_{1_i} \sim N(\mu_{\alpha_1}, \sigma_{\alpha_1}^2), \dots$$

A latent variable W_i was incorporated in the model to estimate overall species richness. It represents whether species i belongs or not to the community of N species. This binary variable was modeled as the outcome of a Bernoulli random variable, defined by:

$$W_i \sim \text{Bernoulli}(\Omega)$$

where Ω describes the probability of belonging to the community. Species detected at least once during the study belong

to the community ($W_i = 1$, with i from 1 to N_{obs}), but species never encountered (i from N_{obs} to $N_{obs} + N_0$) could occupy the sampling area and remain undetected ($W_i = 1$ and $\sum X_{i,j,k} = 0$) or not belong to the community ($W_i = 0$). Therefore, the variable W_i is incorporated in the occupancy model (1):

$$Z_{i,j} \sim \text{Bernoulli}(\psi_{i,j} \times W_i)$$

We implemented the model in a Bayesian framework using the BUGS language and running it in JAGS (Plummer *et al.*, 2003), through the *jagsUI* package (Kellner, 2018) in the R software (R Core Team, 2018). The code is available in Electronic Supplementary Material S2 (ESM S2). Given the lack of prior knowledge of a parameter's true value, parameters and hyper-parameters were implemented with non-informative priors, following common practice. We used uniform distributions from 0 to 1 for the community level parameter Ω , and for the species-level intercepts of occurrence and detection probabilities (α_0 and β_0). We used wide normal priors (with mean 0 and variance 1000) for the means of hyper-distributions of the site-specific and survey-specific effects (the μ_α 's and μ_β 's). We used inverse-gamma priors (*Inv-gamma*(0.1, 0.1)) for the community variances (the σ^2 's) of all these parameters. We ran the analysis for three chains of 15 000 iterations with a burn-in of 15 000 iterations and a thinning rate of 15. Convergence was assessed by examining the Gelman-Rubin statistic (\hat{R}) for each parameter estimate, with $\hat{R} > 1.1$ suggesting a lack of convergence (Gelman & Hill, 2006). Model fit was checked graphically and using the Bayesian P -value (Kéry & Schaub, 2011).

Sampling effectiveness and optimisation

May we reduce the number of spatial replicates? We investigated the effectiveness of the sampling design in reaching inventory completeness using three, four, or five plots. To this end, the inventory completeness ($C_{j,K}$) at site j after K plots ($K = \{3, 4, 5\}$) was calculated as the ratio between the observed number of species ($N_{obs,j,K}$, raw data) and the true number of species estimated by the MSOM (\hat{N}_j). We used the median values of the posterior distributions of the estimated number of species at each sampling site as point estimates for the true species richness. We reduced the number of plots per site by randomly selecting one to four plots among the five samples from the raw data to assess how reducing the number of plots affected the completeness.

We assessed the effectiveness of the sampling design in detecting species through the site-level probability of detecting a species given it is present. We computed for each iteration of the Markov chain Monte Carlo (MCMC) sample the overall detection probability ($P_{i,K}$) for species i to be detected at least once in K plots using the three detection techniques:

$$P_{i,K} = \Pr\left(\sum_{k=1}^K X_{i,j,k} > 0 \mid Z_{i,j} = 1\right)$$

$$P_{i,K} = 1 - (1 - p_i(\text{sighting}))^K \times (1 - p_i(\text{listening}))^K \times (1 - p_i(\text{netting}))^K$$

where $p_i(\text{technique})$ is the plot-level detection probability for species i using the technique cited. We also calculated the site-level detection probability for an "average" species, using estimates of the community-level parameters.

Are the three detection techniques necessary? The completeness and the overall detection probability (i.e. combining the five plots and the three techniques) were also used to investigate if the sampling techniques were complementary or redundant. Therefore, in order to assess if the sampling protocol could be optimised, the species-specific detection probabilities at site level for each technique were calculated, also verifying the effect of omitting the sweep netting technique. The detection probability without sweep netting was obtained from the expression of $P_{i,K}$ above, in which we removed the term involving sweep netting. We did not try to investigate the effects of removing the listening technique because this is necessary to identify certain species that are tricky to distinguish visually, even for Orthoptera experts (Walker, 1964), e.g. *Chorthippus* sp. Fieber, 1852.

Results

As a result of field surveys, we collected 2222 presence data at the plot level (from the three sampling techniques combined) in 405 plots within the 81 sampling sites, belonging to 56 Orthoptera species (ESM S3, Table S1). Eight species represented more than 50% of the presence data, while almost half ($N = 26$) of the observed Orthoptera were detected on less than 5% of the plots. Besides, the true number of species estimated by the model was 75 ($CrI_{95\%} = [61, 98]$).

Species-by-species results concerning estimations of occupancy and detection probabilities can be consulted on the shinyapps web application (see details in ESM S4).

Effects of environmental parameters on the occupancy and the detection probabilities

The MSOM method allowed us to investigate the effects of environmental covariates on the occupancy and detection probabilities at both community and species level (Fig. 2; see details in ESM S4; ESM S3, Tables S2 and S3). The overall trends at the community level are represented by the average species response curves (Fig. 2, top panels). They informed us that species tend to have their distribution optimum in the middle altitudinal range (Fig. 2a) and their detectability influenced positively, but not significantly, by the grass height (Fig. 2d). MSOM enabled us to characterise species response along the altitudinal gradient (Fig. 2, left panels; see details in ESM S4), giving us predictions of the altitudinal optimum of each species. Although the precision decreased for species with low presence data (wider credible intervals: Fig. 2c), the prediction was still informative in terms of optimum position. In addition, the effect of grass height on the probability of detecting a species by sight was also estimated (Fig. 2, right panels; ESM S4 for details), highlighting variations among species, with some Orthoptera more detectable in tall grass (Fig. 2e), and others less detectable with increasing grass height (Fig. 2f).

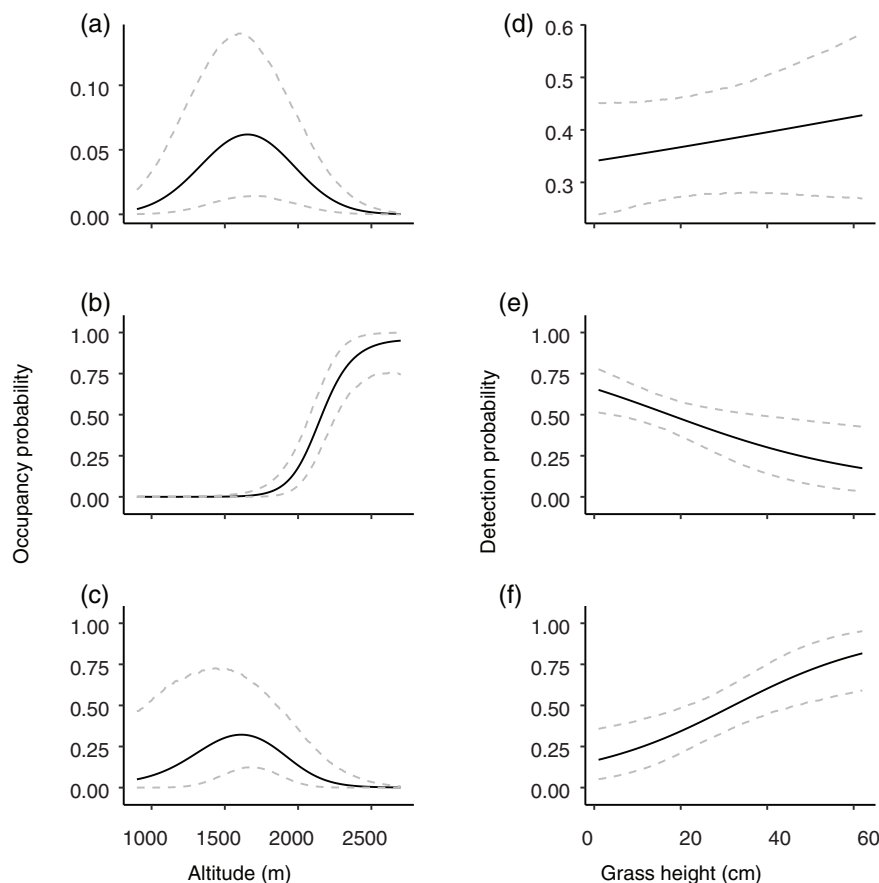


Fig. 2. Effect of altitude on the occupancy probability for (a) an average species at the community level, (b) *Gomphocerus sibiricus sibiricus*, and (c) *Antaxius pedestris*; and the effect of grass height on the probability of sighting for (d) an average species at the community level, (e) *Podisma dechambrei*, and (f) *Euthystira brachyptera*. The solid lines represent the posterior means, and the dashed lines correspond to the 95% credible intervals.

Sampling effectiveness and optimisation

Considering all the detection techniques, the average number of species observed at site level was 9.54 ($CrI_{95\%} = [8.77, 10.32]$), varying from 2 to 17 among sampling sites. Besides, the estimated site-level species richness obtained from the median of posterior distributions of each site was 10.31 ($CrI_{95\%} = [9.46, 11.16]$) on average. The completeness (the ratio between the observed number of species and the estimated number of species) increased with the number of plots (Fig. 3). With five plots, 76 sites (94%) had completeness superior to 80%, while considering four plots, this decreased to 70 sites (86%), and it further decreased to 41 sites (51%) when a survey in only three plots was simulated.

The overall detection probability of an ‘average’ species with all the detection techniques rose sharply when increasing from one to three plots, from 0.45 ($CrI_{95\%} = [0.34, 0.56]$) to 0.83 ($CrI_{95\%} = [0.71, 0.92]$), then grew slightly from 0.91 ($CrI_{95\%} = [0.81, 0.96]$) to 0.95 ($CrI_{95\%} = [0.87, 0.98]$) when adding a fourth and fifth plot, respectively. The detection probabilities estimated at the site level for each detection technique highlighted the importance of sighting in overall detection and the marginality of the two other techniques (see

ESM S4 for details). The detection probability of an ‘average’ species at the site level (five plots) was 0.89 ($CrI_{95\%} = [0.79, 0.95]$) with the sighting technique only. The average detection probabilities at the site level when considering only the listening or the sweep netting technique were much lower: 0.22 ($CrI_{95\%} = [0.10, 0.38]$) and 0.40 ($CrI_{95\%} = [0.26, 0.54]$), respectively. Hence, the detection probability of an ‘average’ species without the sweep netting step (0.91, $CrI_{95\%} = [0.82, 0.97]$) was only slightly lower than the detection probability combining the three techniques.

Detection probability varied consistently among species, ranging from 0.07 ($CrI_{95\%} = [0.01, 0.22]$) to 0.98 ($CrI_{95\%} = [0.96, 1]$) at the plot level and from 0.29 ($CrI_{95\%} = [0.07, 0.71]$) to 1 at the site level when all the techniques were used (Fig. 4; see ESM S4 for details). Some species (38%) had a high probability of being detected after having surveyed a first plot ($P_{i,1} > 0.60$). For these species, the overall detection probability followed an asymptotic curve approaching 1 after three or four plots. In contrast, there were species (27%) with low detection probability in one sampling plot ($P_{i,1} < 0.30$) for which each additional plot sampled sharply increased the overall detection probability at the site level. For almost all species, the differences in detection probability with or without the netting step were not significant,

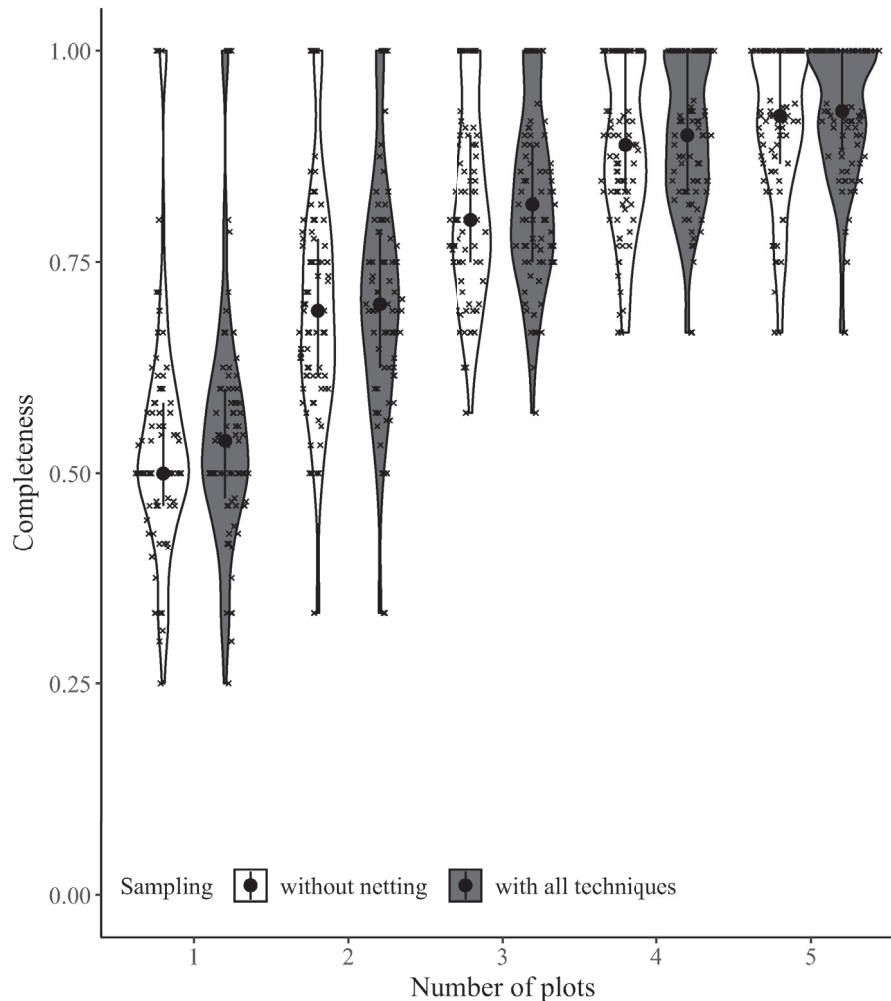


Fig. 3. Inventory completeness at site level according to the number of plots sampled and the detection techniques used. Black dots represent the medians, and black segments represent the first and third quartiles.

except for *Oecanthus pellucens* (Scopoli, 1763) and *Leptophyes punctatissima* (Bosc, 1792). In these two species, the detection probability decreased from 0.95 ($CrI_{95\%} = [0.87, 0.99]$) to 0.66 ($CrI_{95\%} = [0.41, 0.88]$) for *O. pellucens*, and from 0.61 ($CrI_{95\%} = [0.34, 0.87]$) to 0.45 ($CrI_{95\%} = [0.20, 0.73]$) for *L. punctatissima* when removing the sweep netting technique.

Discussion

Using the data from our survey, we were able to estimate the occupancy probabilities of 56 Orthoptera species along an elevation gradient in the Mercantour National Park, while accounting for imperfect detection through a multi-species occupancy model. The species-specific detection probabilities varied widely between species, from 0.29 to 1 at the site level. This could be affected, positively or negatively, or unaffected by the grass height, depending on the species. The inventory completeness was more than 0.80 for 94% of the sites, and the overall detection probability at the community level was

0.95 ($CrI_{95\%} = [0.87, 0.98]$) when using all of a site's five plots and the three sampling techniques. These values slightly decreased when we hypothetically reduced the sampling effort by omitting the netting step or by removing one plot, suggesting that the sampling effort could be reduced with minimal impact on estimate quality.

Reliability of MSOM estimates

Reliability of MSOMs estimates first relies on the respect of the two major assumptions implied by the model: site-closure and no false-presence. The closure assumption indicates that if a plot is occupied, all plots within the site are also occupied. Violating this assumption involves underestimation of detection probabilities and then an overestimation of occupancy probabilities (Kendall & White, 2009). In our study, as we sampled homogeneous grasslands, we were confident about the respect of this assumption. False presences due to misidentification also induce overestimation of occupancy probabilities if not addressed in the

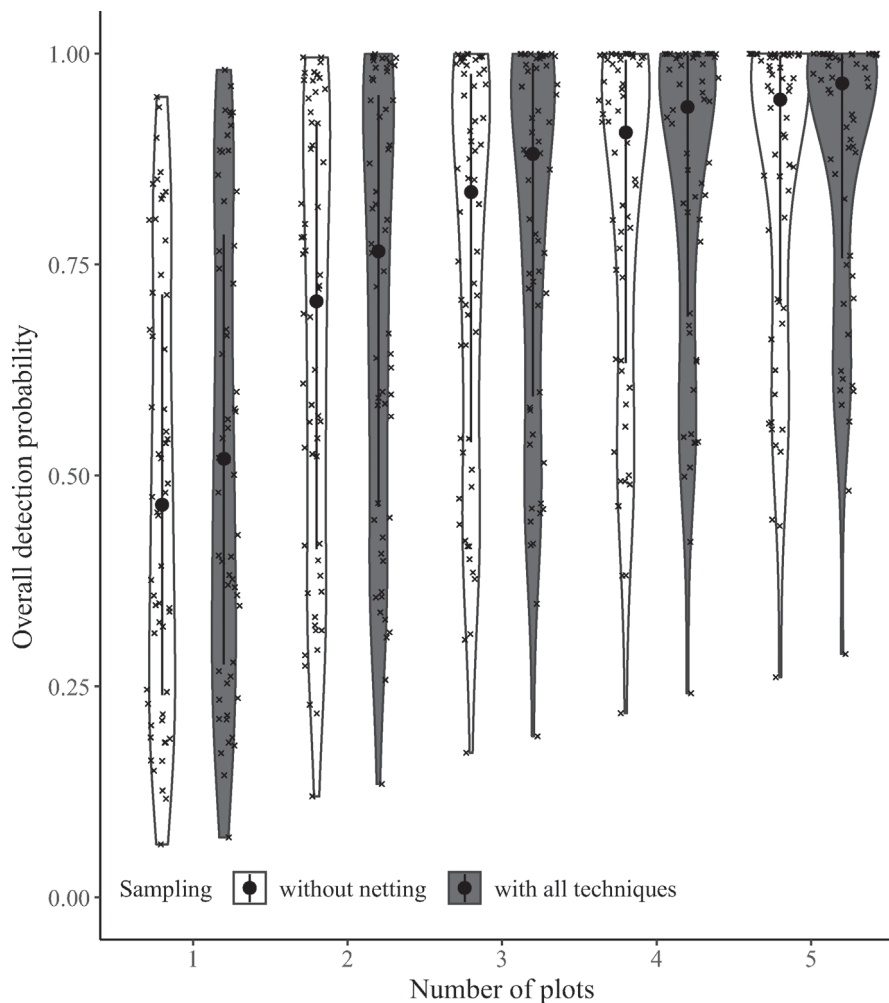


Fig. 4. Distribution of the species-specific overall detection probabilities at the site level depending on the number of sampled plots and the detection techniques used. Black dots represent the medians, and black segments represent the first and third quartiles.

model (Royle & Link, 2006). This assumption is likely to be violated in unexperienced observer. Yet in our case, the observer is highly skilled. Hence, we were confident about the respect of the assumptions, but readers should remember those assumptions when planning to use MSOM.

As expected, the distribution of Orthoptera in our study was structured according to the elevation with (i) maximum in occupancy probabilities of thermophile species such as *Pezotettix giornae* (Rossi, 1794) estimated at low elevations, (ii) wide estimated distribution for generalist species such as *Stauroderus scalaris* (Fischer von Waldheim, 1846), and (iii) arctic-alpine species such as *Gomphocerus sibiricus* (Linnaeus, 1767) having their estimated elevation optimum at high elevations. These results are consistent with what is known about Orthoptera species distribution in the Mercantour National Park area (Gueguen, 1990; Lemonnier, 1999; Braud, com. Pers.). This reliability in the estimated distribution range confirms the relevance of MSOMs to model the effects of biotic or abiotic factors on Orthoptera communities.

MSOM estimates also proved to be useful to assess the inventory completeness reached at each site, referring in particular to the true number of species present, which is rarely known despite its importance in biological studies. However, in some cases, MSOMs may produce unreliable species richness estimates (Guillera-Arroita *et al.*, 2019), especially when detection and/or occupancy are low, inducing a lack of observation and a large number of missing species. Aside from these cases, MSOMs seem to produce reliable estimates in most scenarios and often outperform commonly used estimators such as *iChao2* or *Jack-knife* (Tingley *et al.*, 2020). In this study, the total species richness estimated by MSOM (75 species, $CrI_{95\%} = [61, 98]$) is consistent with the 95 Orthoptera species known to be present above an altitude of 900 m in the Mercantour National Park in grasslands and ecotone habitats (Braud com. Pers.). We chose to integrate ecotone species in the list of known species as we encountered some of those during sampling, such as *Nemobius sylvestris* (Bosc, 1792) or *Pholidoptera griseoptera* (De Geer, 1773). We are thus confident about the reliability of the MSOM estimates at the site level and of the derived inventory completeness.

Importance of detectability in orthopteran studies

Even with the inclusion of the full sampling process (all five plots per sampling site and all three sampling techniques), the overall detection probabilities we estimated were less than one for most species. In some cases, it was quite low: for example, 0.29 for *Eupholidoptera chabrieri* (Charpentier, 1825) or 0.48 for *Calliptamus italicus* (Linnaeus, 1758), confirming the importance of using methods that explicitly correct for imperfect detection. Our results also indicate that detection probability is affected by certain environmental covariates. For instance, we found that Orthoptera detectability by sight may vary with grass height, and that this relationship differs between species. Such a correlation was expected as less mobile species or those living close to the ground, such as *Podisma dechambrei* Chopard, 1952, are likely to be less detectable with increasing grass height. In contrast, the abundance of some Orthoptera, such as *Euthystira brachyptera* (Ocskay, 1826), increases with grass height (Gardiner, 2018), which in turn may increase their detectability (McCarthy *et al.*, 2013). A positive relationship between grass height and species detectability could also be explained by an effect of the higher abundance expected at lower elevations, where grass is slightly higher. Such a correlation between detectability and habitat covariate is likely to generate strong bias when studying the distribution of a species and its relationship with habitat if detection is not modelled explicitly (Lahoz-Monfort *et al.*, 2014). When detection is affected by a habitat covariate, a model that does not include the effect of this covariate on the detection probability may incorrectly identify this habitat covariate as affecting the occupancy rate of the species (Lahoz-Monfort *et al.*, 2014). Such a bias could have huge repercussions in comparative approaches (Archaux *et al.*, 2012), which are commonly used for Orthoptera (e.g. Bomar, 2001; Marini *et al.*, 2009; Löffler *et al.*, 2019).

Some authors have questioned the benefits of modelling imperfect detection (Welsh *et al.*, 2013). They argue that in some cases, i.e. when occupancy is low and detectability is high, 'simple models' perform similarly or better than site-occupancy models. However, this is true only in limited scenarios and assumes high a priori knowledge on the detectability and occupancy of the studied species (Guillera-Arroita *et al.*, 2014). Little is known about the detectability of insects, as there are very few studies accounting for imperfect detection (Kellner & Swihart, 2014; Devarajan *et al.*, 2020). The results obtained studying the orthopteran community in the Mercantour National Park show that occupancy probability is not systematically low, detection probability is not systematically high, and detection probability is highly affected by habitat covariates. These results advocate for the systematic use of MSOMs when studying orthopteran distribution.

Sampling effectiveness and optimisation

The proportion of species richness detected by a survey is a metric commonly used as an indicator of inventory completeness (Moreno & Halffter, 2000; Foggo *et al.*, 2003). Foggo *et al.* (2003) used a completeness threshold of 0.8 to indicate

that an inventory is representative of the community composition in a given site. In our study, 94% of the sites exceeded this threshold with five plots sampled and with the three sampling techniques used. Hence, the composition of the Orthoptera community seems to be well described at the site scale with our sampling design. Our results also suggest that sampling effort may be reduced, notably by omitting one plot or by removing the sweep netting step, while still maintaining a completeness higher than 0.8 for 86% of the sites.

Overall detection probability at the species scale may also be seen as an indicator of sampling efficiency (Moore *et al.*, 2014; Smart *et al.*, 2016). According to the usual detection probability threshold of 0.95 (see e.g. Moore *et al.*, 2014; Smart *et al.*, 2016), the sampling design we used was effective for 31 of the 56 species observed (55%). The number of species above this threshold would decline by three species by omitting the sweep netting step or by six species by removing one plot. The overall detection probability of an 'average' species would also decrease, but slightly and not significantly, from 0.95 ($CrI_{95\%} = [0.87, 0.98]$) with complete sampling, to 0.91 ($CrI_{95\%} = [0.82, 0.97]$) without sweep netting, or 0.91 ($CrI_{95\%} = [0.81, 0.96]$) with four plots. These results confirm that reducing the field effort is possible with a weak impact on detectability.

The best way to optimise the sampling effort, either by removing a detection technique or by reducing the number of plots, may depend on the local species composition, the study objectives and the specific characteristics in the field. In our case, whether we chose to remove a plot or the sweep netting step, the loss in detection probability and in inventory completeness was almost the same. Moreover, in each sampling site the time required to perform five sweep netting steps is quite similar to that needed to fulfil a complete survey in a single plot, around 10 minutes each. Hence, in our case there is not really one choice that is better, especially since the costs are associated mainly with the travel time between sampling sites. However, this could be different in other studies. For example, if we had chosen temporal rather than spatial replicates, removing a plot would have been much more worthwhile than omitting the sweep netting. In the same way, while in our study missing some species was not problematic, as our aim was not to obtain an exhaustive inventory of the entire orthopteran community, other study aims may be different. It should be noted that we found that certain species such as *Oecanthus pellucens* may be missed without the sweep netting technique. It is also important to consider that efficiency of detection techniques depends on their execution order. For instance, walking across the plot during the sighting step made individuals flushing away and so reduced the efficacy of sweep netting sessions. Hence, execution order for the different techniques should be chosen accordingly with the behaviour of the species of interest. Thus, we advocate for implementing a pilot study to help identify how sampling can be best optimised depending on the study objectives.

Conclusion

In this paper, we developed a multi-species occupancy model to demonstrate the need to account for imperfect detection in

insect studies and highlight the potential use of MSOM to investigate the sampling efficiency and optimisation. As our aim was not to explain ecologically Orthoptera distribution or detectability, we developed a relatively simple model including just few covariates in order to show MSOMs potential. However, our model could be easily adapted, with a minimal knowledge in Bayesian computation, to other environmental gradients or pressures commonly investigated in entomological studies, such as management practices (Marini *et al.*, 2009), urbanization levels (Penone *et al.*, 2013), land use intensity (Weking *et al.*, 2016), or climate change (Löffler *et al.*, 2019). Similarly, the effects of other potential covariates on detection probability, such as meteorological (temperature, wind, irradiance, etc.) or phenological variables (date), can easily be implemented in MSOMs. Species traits known to influence detectability, such as mobility capacity, could also be incorporated in MSOMs, for example, by grouping species *a priori* (Pacifi *et al.*, 2014). However, adding species traits in the model can be tricky when using 'data augmentation' approach to estimate species richness as the traits are unknown for species never detected. The unobserved species traits have to be integrated through the hierarchical structure of MSOMs as latent variables, which could however be complicated for non-Bayesian experts. MSOM can also be extended in a dynamic approach to estimate the probability of the colonization or extinction of sites (Dorazio *et al.*, 2010) and thus to study temporal variations in occupancy probability at the community scale. Potential extensions of MSOMs are numerous. Existing model selection methods adapted to Bayesian hierarchical models (Hooten & Hobbs, 2015), thus MSOMs (Broms *et al.*, 2016), may help ecologists to choose the most appropriate model. Bayesian computation can become time-consuming when dealing with large datasets or numerous covariates, and can be a constraint to perform model or variable selection. Yet, alternatives, such as indicator variable selection (Kuo & Mallick, 1998), exist to facilitate variable selection in Bayesian regression models (see: O'Hara & Sillanpää, 2009; for a review), thus in MSOMs (Dorazio *et al.*, 2011). Model selection in MSOMs may be done using information criterion (Drouilly *et al.*, 2018), such as deviance information criterion (DIC; Spiegelhalter *et al.*, 2002) or Watanabe-Akaike information criterion (WAIC; Watanabe, 2013), or describing predictive performance, with cross-validation, for example, (Zipkin *et al.*, 2012). Current MSOMs are now highly flexible and offer an effective framework to model the state or dynamics of insect communities. They can also be used to optimise the efficiency of the sampling design, which could be of interest to many entomologists. We advocate for their systematic use in entomological studies.

Acknowledgements

This work was initiated by a collaboration between the 'Office Français de la Biodiversité' and the 'Centre National de la Recherche Scientifique' (CNRS convention number 169022). The 'Interreg ALCOTRA CCLIMATT' program provided financial support.

The authors declare that they have no conflict of interests.

Author contributions

B.M. carried out the data analysis and drafted the manuscript. T.C. participated in designing the study, coordinating the research and drafting the manuscript. Y.B. carried out the field sampling and participated in designing the study. J.M. and D.C. participated in designing the study. A.B. participated in designing the study, coordinating the research and drafting the manuscript. All the authors read and approved the final manuscript.

Data availability statement

The data that support the findings of this study are available from the corresponding author upon reasonable request.

Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

Supplementary Material S1. Description: Results comparison between model with grass height as occupancy covariate and the model without used in the manuscript.

Supplementary Material S2. Description: Multi-species occupancy model JAGS code.

Supplementary Material S3. Description: Estimates of occupancy, detection, effects of altitude on occupancy and effects of grass height on sighting detection of the Orthoptera community at species- and community-levels.

Supplementary Material S4. Description: An interactive and dynamic application to visualize species-by-species: the detection map, the altitudinal distribution, the detection probability for each sampling technique, and the effect of grass height on the sighting detection.

References

- Archaux, F., Henry, P.-Y. & Gimenez, O. (2012) When can we ignore the problem of imperfect detection in comparative studies? *Methods in Ecology and Evolution*, **3**, 188–194.
- Badenhausser, I., Amouroux, P., Lerin, J. & Bretagnolle, V. (2009) Acridid (orthoptera: Acrididae) abundance in western european grasslands: sampling methodology and temporal fluctuations. *Journal of Applied Entomology*, **133**, 720–732.
- Bale, J.S., Masters, G.J., Hodkinson, I.D., Awmack, C., Bezemer, T.M., Brown, V.K. *et al.* (2002) Herbivory in global climate change research: direct effects of rising temperature on insect herbivores. *Global Change Biology*, **8**, 1–16.
- Bazelet, C.S. & Samways, M.J. (2011) Identifying grasshopper bioindicators for habitat quality assessment of ecological networks. *Ecological Indicators*, **11**, 1259–1269.
- Bieringer, G. & Zulka, K.P. (2003) Shading out species richness: edge effect of a pine plantation on the orthoptera (tettigoniidae and acrididae) assemblage of an adjacent dry grassland. *Biodiversity and Conservation*, **12**, 1481–1495.

- Bomar, C.R. (2001) Comparison of grasshopper (orthoptera: Acrididae) communities on remnant and reconstructed prairies in western Wisconsin. *Journal of Orthoptera Research*, **10**, 105–113.
- Brodie, B.S., Popescu, V.D., Iosif, R., Ciocanea, C., Manolache, S., Vanau, G. *et al.* (2019) Non-lethal monitoring of longicorn beetle communities using generic pheromone lures and occupancy models. *Ecological Indicators*, **101**, 330–340.
- Broms, K.M., Hooten, M.B. & Fitzpatrick, R.M. (2016) Model selection and assessment for multi-species occupancy models. *Ecology*, **97**, 1759–1770.
- Coddington, J.A., Agnarsson, I., Miller, J.A., Kuntner, M. & Hormiga, G. (2009) Undersampling bias: the null hypothesis for singleton species in tropical arthropod surveys. *Journal of Animal Ecology*, **78**, 573–584.
- Devarajan, K., Morelli, T.L. & Tenan, S. (2020) Multi-species occupancy models: review, roadmap, and recommendations. *Ecography*, **n/a**, 1612–1624.
- Dorazio, R.M. & Royle, J.A. (2005) Estimating size and composition of biological communities by modeling the occurrence of species. *Journal of the American Statistical Association*, **100**, 389–398.
- Dorazio, R.M., Royle, J.A., Söderström, B. & Glimskär, A. (2006) Estimating species richness and accumulation by modelling species occurrence and detectability. *Ecology*, **87**, 842–854.
- Dorazio, R.M., Kéry, M., Royle, J.A. & Plattner, M. (2010) Models for inference in dynamic metacommunity systems. *Ecology*, **91**, 2466–2475.
- Dorazio, R.M., Gotelli, N.J. & Ellison, A.M. (2011) *Biodiversity loss in a changing planet*, pp. 277–302. InTech Rijeka, Croatia.
- Drouilly, M., Clark, A. & O’Riain, M.J. (2018) Multi-species occupancy modelling of mammal and ground bird communities in rangeland in the karoo: a case for dryland systems globally. *Biological Conservation*, **224**, 16–25.
- Field, S.A., Tyre, A.J. & Possingham, H.P. (2005) Optimizing allocation of monitoring effort under economic and observational constraints. *The Journal of Wildlife Management*, **69**, 473–482.
- Foggo, A., Rundle, S.D. & Bilton, D.T. (2003) The net result: evaluating species richness extrapolation techniques for littoral pond invertebrates. *Freshwater Biology*, **48**, 1756–1764.
- Gardiner, T. (2018) Grazing and orthoptera: a review. *Journal of Orthoptera Research*, **27**, 3–11.
- Gardiner, T., Hill, J. & Chesmore, D. (2005) Review of the methods frequently used to estimate the abundance of orthoptera in grassland ecosystems. *Journal of Insect Conservation*, **9**, 151–173.
- Gueguen, A. (1990) Impact du pâturage ovin sur la faune sauvage: Exemple des orthoptères.
- Guillera-Arroita, G. (2017) Modelling of species distributions, range dynamics and communities under imperfect detection: advances, challenges and opportunities. *Ecography*, **40**, 281–295.
- Guillera-Arroita, G., Lahoz-Monfort, J.J., MacKenzie, D.I., Wintle, B.A. & McCarthy, M.A. (2014) Ignoring imperfect detection in biological surveys is dangerous: a response to ‘fitting and interpreting occupancy models’. *PLoS One*, **9**, e99571.
- Guillera-Arroita, G., Lahoz-Monfort, J.J., Elith, J., Gordon, A., Kujala, H., Lentini, P.E. *et al.* (2015) Is my species distribution model fit for purpose? Matching data and models to applications. *Global Ecology and Biogeography*, **24**, 276–292.
- Guillera-Arroita, G., Kéry, M. & Lahoz-Monfort, J.J. (2019) Inferring species richness using multispecies occupancy modeling: estimation performance and interpretation. *Ecology and Evolution*, **9**, 780–792.
- Guisan, A. & Thuiller, W. (2005) Predicting species distribution: offering more than simple habitat models. *Ecology Letters*, **8**, 993–1009.
- Guisan, A. & Zimmermann, N.E. (2000) Predictive habitat distribution models in ecology. *Ecological Modelling*, **135**, 147–186.
- Hooten, M. & Hobbs, N. (2015) A guide to bayesian model selection for ecologists. *Ecological Monographs*, **85**, 3–28.
- Kellner, K. (2018) JagsUI: A wrapper around “rjags” to streamline “jags” analyses. R package version 1.5.0.
- Kellner, K.F. & Swihart, R.K. (2014) Accounting for imperfect detection in ecology: a quantitative review. *PLoS One*, **9**, e111436.
- Kendall, W.L. & White, G.C. (2009) A cautionary note on substituting spatial subunits for repeated temporal sampling in studies of site occupancy. *Journal of Applied Ecology*, **46**, 1182–1188.
- Kuo, L. & Mallick, B. (1998) Variable selection for regression models. *Sankhyā: The Indian Journal of Statistics. Series B (1960-2002)*, **60**, 65–81.
- Lahoz-Monfort, J.J., Guillera-Arroita, G. & Wintle, B.A. (2014) Imperfect detection impacts the performance of species distribution models. *Global Ecology and Biogeography*, **23**, 504–515.
- Lemonnier, M. (1999) Les peuplements d’Orthoptères (insecta: Orthoptera) du parc national du mercantour (alpes maritimes, alpes-de-haute-provence). *Bulletin de la Société entomologique de France*, **104**, 149–166.
- Löffler, F., Poniatowski, D. & Fartmann, T. (2019) Orthoptera community shifts in response to land-use and climate change—lessons from a long-term study across different grassland habitats. *Biological Conservation*, **236**, 315–323.
- MacKenzie, D.I., Nichols, J.D., Lachman, G.B., Droege, S., Andrew Royle, J. & Langtimm, C.A. (2002) Estimating site occupancy rates when detection probabilities are less than one. *Ecology*, **83**, 2248–2255.
- MacKenzie, D.I., Nichols, J.D., Hines, J.E., Knutson, M.G. & Franklin, A.B. (2003) Estimating site occupancy, colonization, and local extinction when a species is detected imperfectly. *Ecology*, **84**, 2200–2207.
- MacKenzie, D.I., Nichols, J.D., Royle, J.A., Pollock, K.H., Bailey, L. & Hines, J.E. (2006) *Occupancy estimation and modeling: Inferring patterns and dynamics of species occurrence*. Burlington, MA: Elsevier.
- Malinowska, A.H., van Strien, A.J., Verboom, J., WallisdeVries, M.F. & Opdam, P. (2014) No evidence of the effect of extreme weather events on annual occurrence of four groups of ectothermic species. *PLoS One*, **9**, 1–10.
- Marini, L., Fontana, P., Battisti, A. & Gaston, K.J. (2009) Response of orthopteran diversity to abandonment of semi-natural meadows. *Agriculture, Ecosystems & Environment*, **132**, 232–236.
- Mata, L., Goula, M. & Hahs, A. (2014) Conserving insect assemblages in urban landscapes: accounting for species-specific responses and imperfect detection. *Journal of Insect Conservation*, **18**, 885–894.
- McCarthy, M.A., Moore, J.L., Morris, W.K., Parris, K.M., Garrard, G.E., Vesk, P.A. *et al.* (2013) The influence of abundance on detectability. *Oikos*, **122**, 717–726.
- Moore, A.L., McCarthy, M.A., Parris, K.M. & Moore, J.L. (2014) The optimal number of surveys when detectability varies. *PLoS One*, **9**, e115345.
- Moreno, C.E. & Halfiter, G. (2000) Assessing the completeness of bat biodiversity inventories using species accumulation curves. *Journal of Applied Ecology*, **37**, 149–158.
- Moritz, C., Patton, J.L., Conroy, C.J., Parra, J.L., White, G.C. & Beissinger, S.R. (2008) Impact of a century of climate change on small-mammal communities in Yosemite national park, USA. *Science*, **322**, 261–264.
- Nichols, J.D., Hines, J.E., Mackenzie, D.I., Seamans, M.E. & Gutiérrez, R.J. (2007) OCCUPANCY estimation and modeling with multiple states and state uncertainty. *Ecology*, **88**, 1395–1400.
- O’Hara, R.B. & Sillanpää, M.J. (2009) A review of bayesian variable selection methods: what, how and which. *Bayesian Analysis*, **4**, 85–117.

- Ovaskainen, O. & Soininen, J. (2011) Making more out of sparse data: hierarchical modeling of species communities. *Ecology*, **92**, 289–295.
- Pacifici, K., Zipkin, E.F., Collazo, J.A., Irizarry, J.I. & DeWan, A. (2014) Guidelines for a priori grouping of species in hierarchical community models. *Ecology and Evolution*, **4**, 877–888.
- Pecchi, M., Marchi, M., Burton, V., Giannetti, F., Moriondo, M., Bernetti, I. *et al.* (2019) Species distribution modelling to support forest management. A literature review. *Ecological Modelling*, **411**, 108817.
- Penone, C., Le Viol, I., Pellissier, V., Julien, J.-F., Bas, Y. & Kerbiriou, C. (2013) Use of large-scale acoustic monitoring to assess anthropogenic pressures on orthoptera communities. *Conservation Biology*, **27**, 979–987.
- Phillips, S.J., Dudík, M., Elith, J., Graham, C.H., Lehmann, A., Leathwick, J. *et al.* (2009) Sample selection bias and presence-only distribution models: implications for background and pseudo-absence data. *Ecological Applications*, **19**, 181–197.
- Pierik, M.E., Gusmeroli, F., Marianna, G.D., Tamburini, A. & Bocchi, S. (2017) Meadows species composition, biodiversity and forage value in an alpine district: relationships with environmental and dairy farm management variables. *Agriculture, Ecosystems & Environment*, **244**, 14–21.
- Plummer, M., *et al.* (2003) JAGS: A program for analysis of bayesian graphical models using gibbs sampling. *Proceedings of the 3rd international workshop on distributed statistical computing*. Vienna, Austria., p. 10.
- R Core Team (2018) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Roth, T., Allan, E., Pearman, P.B. & Amrhein, V. (2018) Functional ecology and imperfect detection of species. *Methods in Ecology and Evolution*, **9**, 917–928.
- Royle, J.A. & Link, W.A. (2006) Generalized occupancy models allowing false positive and false negative errors. *Ecology*, **87**, 835–841.
- Royle, J.A., Dorazio, R.M. & Link, W.A. (2007) Analysis of multinomial models with unknown index using data augmentation. *Journal of Computational and Graphical Statistics*, **16**, 67–85.
- Rushton, S.P., Ormerod, S.J. & Kerby, G. (2004) New paradigms for modelling species distributions? *Journal of Applied Ecology*, **41**, 193–200.
- Silva, D.P., Andrade, A.F., Oliveira, J.P., Morais, D.M., Vieira, J.E. & Engel, M.S. (2019) Current and future ranges of an elusive north american insect using species distribution models. *Journal of Insect Conservation*, **23**, 175–186.
- Smart, A.S., Weeks, A.R., van Rooyen, A.R., Moore, A., McCarthy, M.A. & Tingley, R. (2016) Assessing the cost-efficiency of environmental dna sampling. *Methods in Ecology and Evolution*, **7**, 1291–1298.
- Spiegelhalter, D.J., Best, N.G., Carlin, B.P. & Van Der Linde, A. (2002) Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **64**, 583–639.
- Ter Braak, C.J. & Smilauer, P. (2002) *CANOCO reference manual and canodraw for windows user's guide: Software for canonical community ordination (version 4.5)*. www.canoco.com.
- Tingley, M.W., Nadeau, C.P. & Sandor, M.E. (2020) Multi-species occupancy models as robust estimators of community richness. *Methods in Ecology and Evolution*, **11**, 633–642.
- Tyre, A.J., Tenhumberg, B., Field, S.A., Niejalke, D., Parris, K. & Possingham, H.P. (2003) Improving precision and reducing bias in biological surveys: estimating false-negative error rates. *Ecological Applications*, **13**, 1790–1801.
- Vaughan, I.P. & Ormerod, S.J. (2005) The continuing challenges of testing species distribution models. *Journal of Applied Ecology*, **42**, 720–730.
- Veran, S., Simpson, S.J., Sword, G.A., Deveson, E., Piry, S., Hines, J.E. *et al.* (2015) Modeling spatiotemporal dynamics of outbreaking species: influence of environment and migration in a locust. *Ecology*, **96**, 737–748.
- Walker, T.J. (1964) Cryptic species among sound-producing ensiferan orthoptera (gryllidae and tettigoniidae). *The Quarterly Review of Biology*, **39**, 345–355.
- Watanabe, S. (2013) A widely applicable bayesian information criterion. *Journal of Machine Learning Research*, **14**, 867–897.
- Weking, S., Kämpf, I., Mathar, W. & Hölzel, N. (2016) Effects of land use and landscape patterns on orthoptera communities in the western siberian forest steppe. *Biodiversity and Conservation*, **25**, 2341–2359.
- Welsh, A.H., Lindenmayer, D.B. & Donnelly, C.F. (2013) Fitting and interpreting occupancy models. *PLoS One*, **8**, 1–21.
- Wolda, H. (1988) Insect seasonality: why? *Annual Review of Ecology and Systematics*, **19**, 1–18.
- Yoccoz, N.G., Nichols, J.D. & Boulinier, T. (2001) Monitoring of biological diversity in space and time. *Trends in Ecology & Evolution*, **16**, 446–453.
- Zipkin, E.F., DeWan, A. & Andrew Royle, J. (2009) Impacts of forest fragmentation on species richness: a hierarchical approach to community modelling. *Journal of Applied Ecology*, **46**, 815–822.
- Zipkin, E.F., Grant, E.H.C. & Fagan, W.F. (2012) Evaluating the predictive abilities of community occupancy models using auc while accounting for imperfect detection. *Ecological Applications*, **22**, 1962–1972.

Accepted 16 November 2020

Associate Editor: Bernard Roitberg

Electronic supplementary material 1

Results comparison between model with grass height as occupancy covariate and model without this covariate.

The aim of the study was to show MSOMs potential in entomological study. Hence, we developed a relatively simple MSOMs, with just few covariates. However, omitting potential important predictors of Orthoptera occupancy or detectability could bias the results. We could assumed that grass height influence Orthoptera occupancy. Therefore, we fitted another model, similar in the detection model, but with the effect of grass height added in the occupancy model:

$$\text{logit}(\psi_{i,j}) = \alpha_0 + \alpha_1 \times \text{altitude}_j + \alpha_2 \times \text{altitude}_j^2 + \alpha_3 \times \text{height}_j$$

with α_3 , the linear effect of grass height on the occupancy probability of species i .

We compared the principal results found with this model to those presented in the manuscript. The relationship between the environmental parameters and the detection and occupancy probabilities did not change meaningfully between the two models (Figure S1). May be because the effect of grass height on the occupancy at the community-level and for most of the species was not significant (Figure S2). The estimated species richness was 77.6 ($IC_{95\%}=[62.975, 102]$), which is not significantly different than the estimate of the simplest model ($\hat{N}=74.62$, $IC_{95\%}=[61, 98]$). At the site level, the completeness estimates

were very similar between the two models (Figure S3). With the simplest model, 76 sites have completeness superior to 80%, while the second model estimated 77 sites above this threshold. The overall detection probabilities were also very similar (Figure S4). 31 species had an overall detection probability at the site level upper than 95% when considering that grass height influence both detection and occupancy, against 30 species when accounting only for grass height effect on detectability. Results did not change meaningfully, neither did our inferences and interpretations. Hence, we kept the simplest model to facilitate readers understanding.

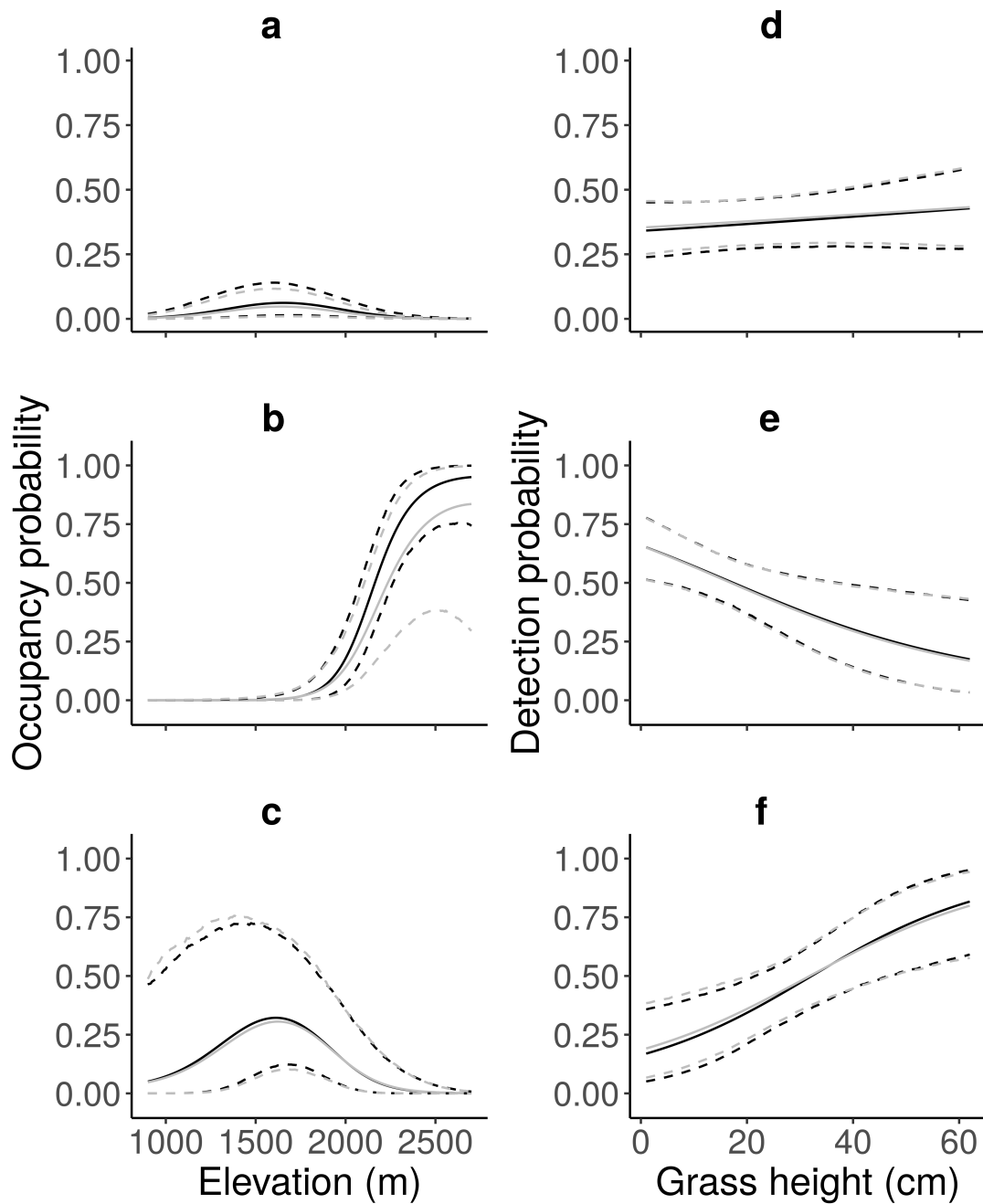


Figure S1: Effect of the altitude on the occupancy probability for (a) an average species at the community-level, (b) *Gomphocerus sibiricus sibiricus*, (c) *Antaxius pedestris*, and effect of the grass height on the probability of sighting (d) an average species at the community-level, (e) *Podisma dechambrei*, and (f) *Euthystira brachyptera*. The colors correspond to

model specification, with one model with grass height as occupancy covariate (grey lines) and one without (black lines). The solid lines represent the posterior mean, and the dashed lines correspond to the 95% credible interval.

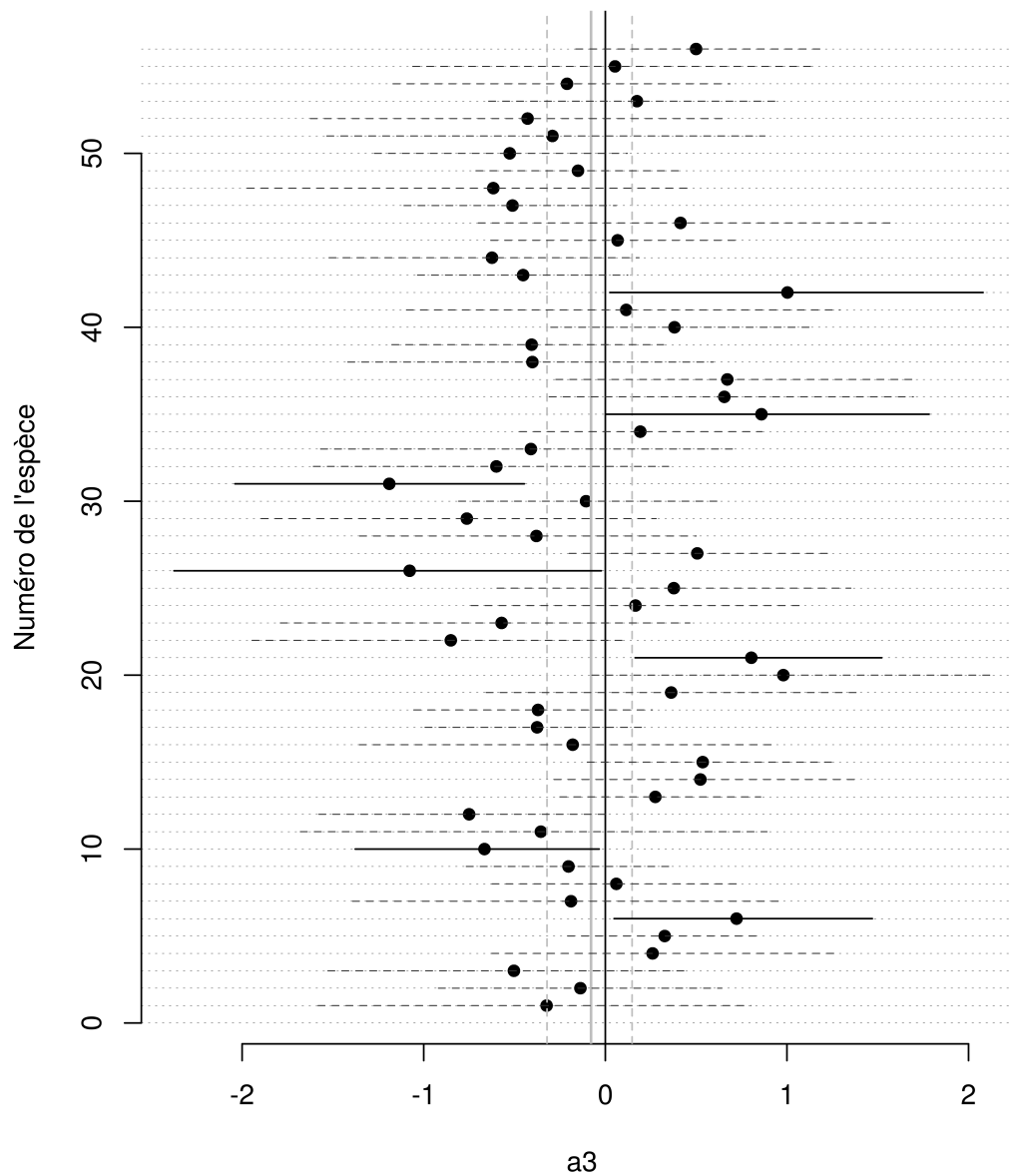


Figure S2: Estimates on the logit scale of the linear grass height effect on species-specific occupancy probabilities for the 56 Orthoptera species detected. Black points correspond to medians of posterior distributions, and the segment to the credible interval associated. Segments with solid lines represent significant effects.

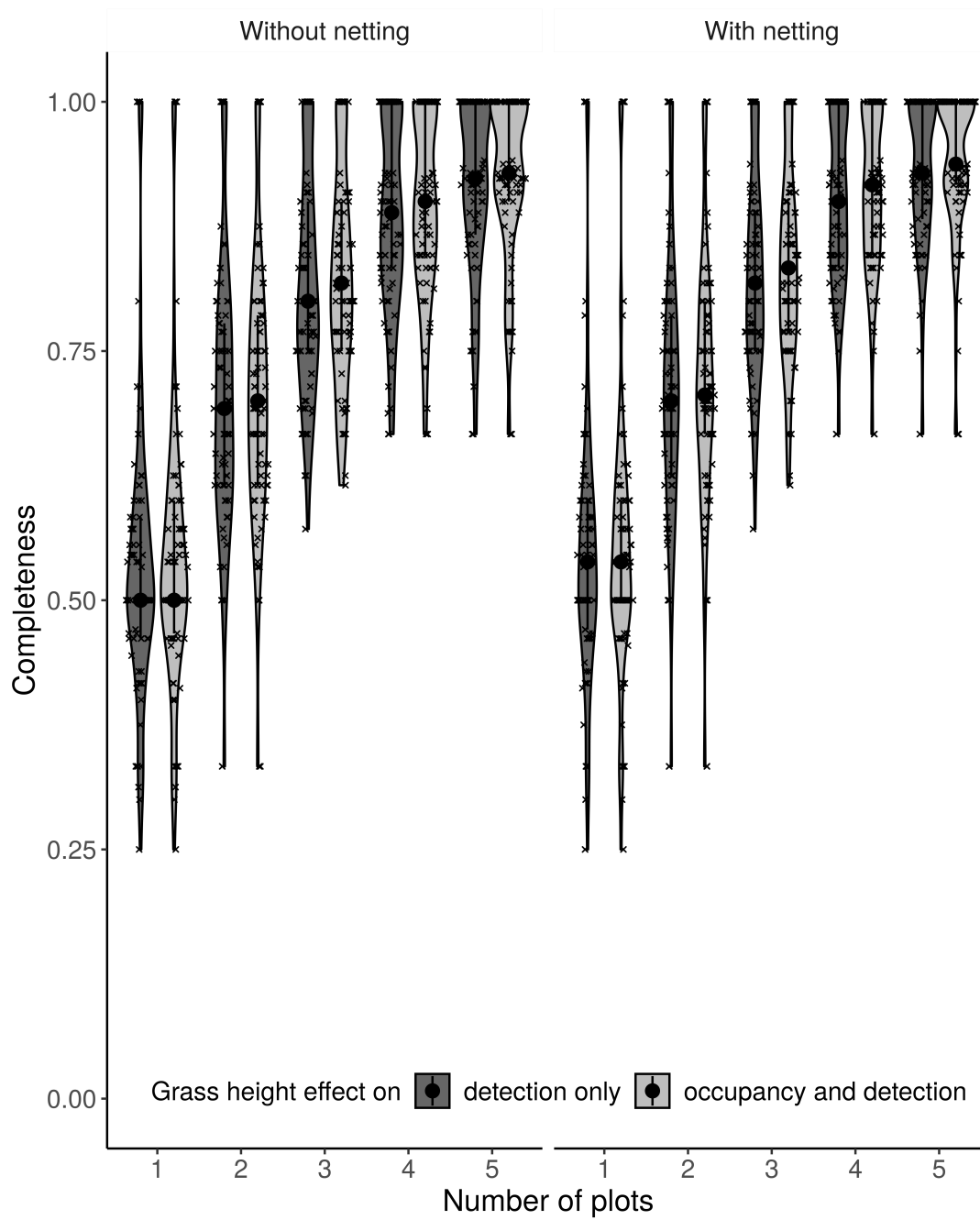


Figure S3: Inventory completeness at site level according to the model used, the number of plots sampled and the detection techniques used. Black dots represent the medians, and black segments represent the first and third quartiles.

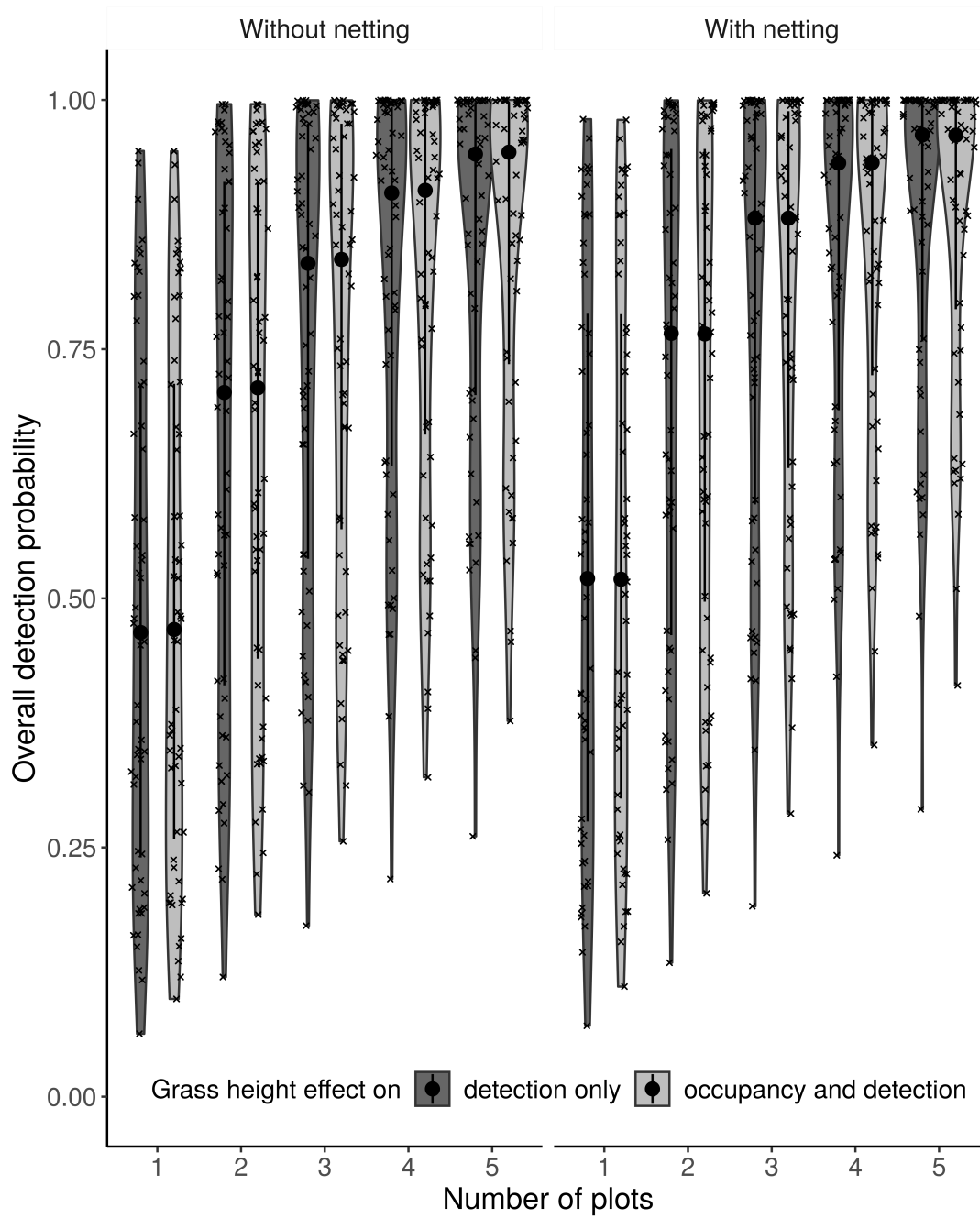


Figure S4: Distribution of the species-specific overall detection probabilities at the site level depending on the number of sampled plots, the model specification and the detection techniques used. Black dots represent the medians, and black segments represent the first and third quartiles.

Electronic supplementary material 2

Multi-species occupancy model JAGS code

```
model{

  ## Prior distributions for community-level model parameters ##

  omega ~ dunif(0,1)          #probability of species inclusion

  # Community means #

  a0.mean ~ dunif(0,1)

  mu.a0 <- log(a0.mean) - log(1-a0.mean)

  mu.a1 ~ dnorm(0, 0.001)

  mu.a2 ~ dnorm(0, 0.001)

  mu.b0 ~ dnorm(0, 0.001)

  mu.b1 ~ dnorm(0, 0.001)

  mu.b2 ~ dnorm(0, 0.001)

  mu.b3 ~ dnorm(0, 0.001)

  # Community precisions (inverse variances) #

  tau.a0 ~ dgamma(0.1,0.1)

  tau.a1 ~ dgamma(0.1,0.1)

  tau.a2 ~ dgamma(0.1,0.1)

  tau.b0 ~ dgamma(0.1,0.1)

  tau.b1 ~ dgamma(0.1,0.1)
```

```

tau.b2 ~ dgamma(0.1,0.1)

tau.b3 ~ dgamma(0.1,0.1)

for (i in 1:(n+nzeroes)) {

  ## Prior distributions for species level model parameters ##

  w[i] ~ dbern(omega) #latent variable: species i belongs (w[i]=1) or not (w[i]=0) to the
community

  # Parameters for occupancy model #

  a0[i] ~ dnorm(mu.a0, tau.a0) #random intercepts, occupancy at mean elevation
  a1[i] ~ dnorm(mu.a1, tau.a1) #random linear effects of altitude
  a2[i] ~ dnorm(mu.a2, tau.a2) #random quadratic effects of altitude

  # Parameters for occupancy model #

  b0[i] ~ dnorm(mu.b0, tau.b0) #random intercept, sighting detection
  b1[i] ~ dnorm(mu.b1, tau.b1) #random effect of listening
  b2[i] ~ dnorm(mu.b2, tau.b2) #random effect of sweep netting
  b3[i] ~ dnorm(mu.b3, tau.b3) #random linear effect of grass height on sighting

  ## Occupancy model specification ##

  for (j in 1:J) {

    logit(psi[j,i]) <- a0[i] + #species occupancy at mean elevation
      a1[i]*covSite1[j] + a2[i]*pow(covSite1[j],2) #altitude effects

    mu.psi[j,i] <- psi[j,i]*w[i] #site j could be occupied by species i only if it belongs to the
community
  }
}

```



```

#if species i belongs to the community: mu.psi=psi, else: mu.psi=0

Z[j,i] ~ dbern(mu.psi[j,i])    #latent variable: true occurrence status of species i at site j

## Detection model specification ##

for (k in 1:K[j]) {

logit(p[j,k,i]) <- b0[i] + #sighting intercept

                b1[i]*covDetection1[j,k] + #listening effect

                b2[i]*covDetection2[j,k] + #sweep netting effect

                b3[i]*covDetection0[j,k]*covDetection3[j,k] #grass height effect on sighting

mu.p[j,k,i] <- p[j,k,i]*Z[j,i]  #species i could be detected in site j during survey k only if site j
is occupied

#if j is occupied: mu.p=p, else: mu.p=0

X[j,k,i] ~ dbern(mu.p[j,k,i])  #binary indicator: observed occurrence status for species i in
site j on plot k

        ## Assess model adjustment ##

#Create simulated dataset to calculate the Bayesian p-value

Xnew[j,k,i] ~ dbern(mu.p[j,k,i])

#Pearson residuals

d[j,k,i]<- abs(X[j,k,i] - mu.p[j,k,i])

dnew[j,k,i]<- abs(Xnew[j,k,i]- mu.p[j,k,i])

d2[j,k,i]<- pow(d[j,k,i],2)

dnew2[j,k,i]<- pow(dnew[j,k,i],2)

```

```
}
```

```
dsum[j,i]<- sum(d2[j,1:K[j],i])
```

```
dnewsum[j,i]<- sum(dnew2[j,1:K[j],i])
```

```
}
```

```
}
```

```
#Calculate the discrepancy measure
```

```
p.fit<-sum(dsum[1:J,1:(n)])
```

```
p.fitnew<-sum(dnewsum[1:J,1:(n)])
```

```
## Species diversity estimates ##
```

```
#Sum all species observed (n) and unobserved species (n0) to estimate the
```

```
#true richness
```

```
n0 <- sum(w[(n+1):(n+nzeroes)])
```

```
N <- n + n0
```

```
#Create a loop to determine site level richness estimates for the
```

```
#whole community.
```

```
for(j in 1:J){
```

```
  Nsite[j]<- sum(Z[j,1:(n+nzeroes)])
```

```
}
```

```
}
```

Electronic supplementary material 3

Estimates of occupancy, detection, effects of altitude on occupancy and effects of grass height on sighting detection of the Orthoptera community at species- and community-levels. Estimates of occupancy, detection, effects of altitude on occupancy and effects of grass height on sighting detection of the Orthoptera community at species- and community-levels.

Table S1: Mean and 95% credible intervals for species-specific probabilities of occupancy at a medium altitude of 1869 m and, species-specific detection probabilities at the plot level for each technique and at the site level considering the three techniques combined or without the sweep netting. Also included is the naive occupancy, i.e. the ratio between the number of sites where the species was detected and the total number of sites (81 sampled sites).

Species	Naive occupancy	Plot-level detection probability									Site-level detection probability								
		Occupancy probability at 1869 m of elevation			Sighting			Listening			Netting			Complete sampling			Sampling without netting step		
		Mean	Q2.5	Q97.5	Mean	Q2.5	Q97.5	Mean	Q2.5	Q97.5	Mean	Q2.5	Q97.5	Mean	Q2.5	Q97.5	Mean	Q2.5	Q97.5
<i>Anonconotus baracunensis occidentalis</i>	0.02	0.02	0	0.11	0.18	0.03	0.51	0.01	0	0.11	0.02	0	0.16	0.67	0.19	0.99	0.63	0.17	0.98
<i>Anonconotus ghiliani</i>	0.25	0.29	0.15	0.47	0.3	0.2	0.41	0	0	0.02	0.05	0.02	0.09	0.86	0.72	0.95	0.83	0.67	0.93
<i>Anonconotus mercantouri</i>	0.1	0.06	0.02	0.17	0.43	0.24	0.64	0.01	0	0.05	0.07	0.02	0.18	0.95	0.8	1	0.93	0.76	0.99
<i>Antaxius pedestris</i>	0.1	0.21	0.08	0.49	0.15	0.06	0.3	0	0	0.03	0.03	0.01	0.09	0.63	0.3	0.88	0.57	0.26	0.84
<i>Arcyptera fusca</i>	0.36	0.48	0.34	0.62	0.4	0.32	0.48	0.25	0.19	0.33	0.06	0.03	0.1	0.99	0.97	0.99	0.98	0.96	0.99
<i>Bicolorana bicolor</i>	0.14	0.23	0.12	0.39	0.33	0.21	0.47	0.18	0.09	0.31	0.06	0.02	0.14	0.96	0.88	0.99	0.94	0.85	0.99
<i>Calliptamus italicus</i>	0.02	0.02	0	0.22	0.09	0.01	0.34	0	0	0.07	0.01	0	0.09	0.47	0.09	0.93	0.43	0.08	0.9
<i>Calliptamus siciliae</i>	0.22	0.1	0.04	0.21	0.54	0.43	0.65	0	0	0.03	0.16	0.09	0.25	0.99	0.97	1	0.98	0.94	1
<i>Chorthippus apricarius</i>	0.46	0.64	0.5	0.76	0.63	0.56	0.69	0.28	0.22	0.35	0.45	0.38	0.53	1	1	1	1	1	1
<i>Chorthippus biguttulus</i>	0.26	0.28	0.17	0.42	0.48	0.38	0.57	0.38	0.29	0.47	0.22	0.15	0.3	1	1	1	1	0.99	1
<i>Chorthippus brunneus brunneus</i>	0.01	0.02	0	0.21	0.12	0.01	0.49	0	0	0.14	0.02	0	0.17	0.58	0.09	0.99	0.54	0.08	0.98
<i>Chorthippus dorsatus</i>	0.14	0.12	0.05	0.24	0.88	0.79	0.94	0.51	0.37	0.64	0.71	0.59	0.82	1	1	1	1	1	1
<i>Chorthippus saulcyi daimei</i>	0.37	0.52	0.38	0.66	0.77	0.7	0.84	0.34	0.27	0.43	0.53	0.44	0.61	1	1	1	1	1	1
<i>Chorthippus vagans vagans</i>	0.09	0.01	0	0.06	0.71	0.53	0.85	0.4	0.23	0.6	0.23	0.11	0.4	1	1	1	1	1	1
<i>Decticus verrucivorus verrucivorus</i>	0.37	0.58	0.42	0.73	0.26	0.19	0.34	0.1	0.06	0.15	0.03	0.01	0.06	0.89	0.8	0.94	0.87	0.77	0.93
<i>Depressotetrix depressa</i>	0.02	0.04	0.01	0.19	0.15	0.03	0.41	0	0	0.08	0.06	0.01	0.24	0.69	0.22	0.98	0.58	0.16	0.94
<i>Ephippiger terrestris</i>	0.6	0.78	0.64	0.88	0.54	0.48	0.6	0.27	0.22	0.33	0.18	0.14	0.24	1	1	1	1	0.99	1
<i>Euchorthippus declivus</i>	0.51	0.64	0.47	0.78	0.78	0.71	0.83	0.26	0.2	0.32	0.59	0.52	0.66	1	1	1	1	1	1
<i>Euchorthippus elegantulus</i>	0.04	0.01	0	0.04	0.67	0.39	0.87	0.19	0.06	0.46	0.29	0.11	0.57	1	0.98	1	0.99	0.96	1
<i>Eupholidoptera chabrieri</i>	0.06	0.15	0.02	0.62	0.04	0.01	0.18	0	0	0.02	0	0	0.03	0.27	0.06	0.69	0.25	0.05	0.65
<i>Euthystira brachyptera</i>	0.16	0.22	0.11	0.38	0.29	0.17	0.45	0.01	0	0.04	0.04	0.01	0.09	0.86	0.67	0.97	0.82	0.62	0.95
<i>Gomphocerus sibiricus sibiricus</i>	0.25	0.05	0.01	0.14	0.6	0.45	0.73	0.4	0.26	0.55	0.25	0.15	0.39	1	1	1	1	0.99	1
<i>Gryllus campestris</i>	0.02	0.02	0	0.09	0.27	0.08	0.57	0.01	0	0.11	0.06	0.01	0.22	0.83	0.43	0.99	0.78	0.36	0.99
<i>Leptophyes punctatissima</i>	0.11	0.05	0.01	0.23	0.09	0.03	0.2	0	0	0.02	0.05	0.02	0.13	0.52	0.26	0.82	0.38	0.16	0.68
<i>Metrioptera saussuriana</i>	0.04	0.05	0.01	0.14	0.43	0.21	0.66	0.17	0.04	0.42	0.05	0.01	0.19	0.97	0.86	1	0.96	0.84	1
<i>Myrmeleotettix maculatus</i>	0.12	0.05	0.01	0.14	0.35	0.17	0.58	0.16	0.06	0.34	0.2	0.08	0.4	0.97	0.85	1	0.93	0.75	1
<i>Nemobius sylvestris</i>	0.17	0.15	0.06	0.32	0.05	0.02	0.12	0.24	0.13	0.38	0.01	0	0.03	0.8	0.59	0.95	0.79	0.58	0.94
<i>Oecanthus pellucens</i>	0.1	0.01	0	0.05	0.17	0.09	0.31	0	0	0.04	0.32	0.19	0.49	0.94	0.82	0.99	0.63	0.37	0.85

<i>Psophus stridulus</i>	0.1	0.14	0.06	0.27	0.34	0.21	0.5	0.06	0.02	0.15	0.05	0.02	0.13	0.93	0.8	0.99	0.9	0.75	0.98
<i>Roeseliana roeselii</i>	0.14	0.16	0.07	0.29	0.4	0.27	0.54	0.17	0.09	0.29	0.06	0.02	0.13	0.98	0.93	1	0.97	0.9	0.99
<i>Sepiana sepium</i>	0.01	0.01	0	0.07	0.31	0.06	0.76	0.01	0	0.23	0.05	0	0.35	0.83	0.29	1	0.8	0.26	1
<i>Stauroderus scalaris</i>	0.64	0.82	0.7	0.9	0.75	0.7	0.8	0.6	0.54	0.66	0.33	0.27	0.38	1	1	1	1	1	1
<i>Stenobothrus coticus</i>	0.1	0.01	0	0.06	0.52	0.3	0.73	0.38	0.19	0.62	0.2	0.08	0.4	1	0.97	1	0.99	0.95	1
<i>Stenobothrus lineatus</i>	0.6	0.73	0.6	0.84	0.61	0.55	0.67	0.27	0.22	0.33	0.18	0.14	0.23	1	1	1	1	1	1
<i>Stenobothrus nigromaculatus</i>	0.32	0.33	0.19	0.49	0.55	0.45	0.65	0.22	0.15	0.31	0.28	0.2	0.37	1	1	1	0.99	0.98	1
<i>Stenobothrus rubicundulus</i>	0.01	0.02	0	0.17	0.12	0.01	0.46	0	0	0.13	0.02	0	0.15	0.58	0.09	0.98	0.54	0.07	0.97
<i>Tessellana tessellata</i>	0.02	0	0	0.03	0.43	0.2	0.7	0.01	0	0.13	0.13	0.03	0.36	0.95	0.78	1	0.92	0.68	1
<i>Tettigonia cantans</i>	0.06	0.09	0.03	0.22	0.29	0.14	0.5	0	0	0.05	0.05	0.01	0.14	0.85	0.59	0.98	0.81	0.54	0.97
<i>Tettigonia viridissima</i>	0.06	0.06	0.01	0.19	0.19	0.08	0.38	0	0	0.04	0.03	0.01	0.09	0.71	0.39	0.94	0.67	0.34	0.92
<i>Tylopsis lilifolia</i>	0.04	0	0	0.02	0.81	0.59	0.93	0.35	0.13	0.66	0.32	0.13	0.6	1	1	1	1	1	1
<i>Yersinella beybienkoi</i>	0.2	0.12	0.04	0.24	0.55	0.42	0.67	0	0	0.02	0.09	0.04	0.17	0.99	0.95	1	0.98	0.93	1

Table S2: Mean and 95% credible intervals for the species-specific effects of the habitat covariates on the logit-scale, as estimated by the MSOM.

Species	Effect of altitude on occupancy probability						Effect of grass height on sighting detection probability		
	Linear			Quadratic			Mean	Q2.5	Q97.5
	Mean	Q2.5	Q97.5	Mean	Q2.5	Q97.5			
<i>Anonconotus baracunensis occidentalis</i>	2.26	0.15	4.81	-0.62	-1.46	0.26	0.04	-1.04	1.14
<i>Anonconotus ghiliani</i>	1.8	0.79	3.01	-0.96	-1.7	-0.3	-0.66	-1.21	-0.18
<i>Anonconotus mercantouri</i>	2.39	0.92	4.19	-0.92	-1.76	-0.17	-0.27	-1.2	0.61
<i>Antaxius pedestris</i>	-1.6	-3.31	-0.19	-1.19	-2.13	-0.44	-0.27	-0.9	0.32
<i>Arcyptera fusca</i>	-0.25	-0.88	0.33	-0.67	-1.18	-0.22	-0.04	-0.25	0.16
<i>Bicolorana bicolor</i>	0.13	-0.98	1.27	-1.29	-2.26	-0.56	0.19	-0.19	0.59
<i>Calliptamus italicus</i>	-2.92	-5.64	-0.67	-0.63	-1.46	0.23	-0.19	-1.22	0.8
<i>Calliptamus siciliae</i>	-2.91	-4.47	-1.58	-0.5	-1.16	0.18	0.04	-0.31	0.39
<i>Chorthippus apricarius</i>	0.2	-0.41	0.83	-1.02	-1.6	-0.53	0.16	-0.01	0.33
<i>Chorthippus biguttulus</i>	-0.72	-1.42	-0.09	-0.31	-0.76	0.13	0.14	-0.15	0.41
<i>Chorthippus brunneus brunneus</i>	0.77	-1.6	3.37	-0.86	-1.78	-0.01	-0.01	-1.09	1.05
<i>Chorthippus dorsatus</i>	-1.76	-3.07	-0.64	-0.7	-1.37	-0.1	0.92	0.4	1.49
<i>Chorthippus saulcyi daimeii</i>	0.09	-0.54	0.72	-0.89	-1.45	-0.39	-0.78	-1.01	-0.56
<i>Chorthippus vagans vagans</i>	-3.47	-5.58	-1.6	-0.66	-1.38	0.05	0.49	0.02	0.98
<i>Decticus verrucivorus verrucivorus</i>	-0.02	-0.7	0.64	-0.92	-1.5	-0.42	-0.15	-0.44	0.13
<i>Depressotetrix depressa</i>	-1.3	-3.63	0.72	-1	-1.96	-0.19	-0.45	-1.57	0.58
<i>Ephippiger terrestris</i>	-2.1	-3.01	-1.33	-0.96	-1.52	-0.46	0	-0.16	0.16
<i>Euchorthippus declivus</i>	-3.19	-4.51	-2.1	-1.23	-1.91	-0.63	-0.02	-0.2	0.16
<i>Euchorthippus elegantulus</i>	-3.14	-5.41	-1.11	-0.54	-1.27	0.26	0.67	-0.05	1.43
<i>Eupholidoptera chabrieri</i>	-2.78	-6.03	-0.58	-0.93	-1.9	-0.04	0.65	-0.05	1.35
<i>Euthystira brachyptera</i>	-1.38	-2.6	-0.31	-0.86	-1.53	-0.26	0.79	0.36	1.25
<i>Gomphocerus sibiricus sibiricus</i>	4.33	2.54	6.32	-0.56	-1.39	0.33	-0.48	-1.09	0.11
<i>Gryllus campestris</i>	-2.42	-4.71	-0.39	-0.86	-1.72	-0.08	-0.6	-1.67	0.4
<i>Leptophyes punctatissima</i>	-3.87	-6.67	-1.82	-0.8	-1.7	0.11	0.47	-0.04	0.98
<i>Metroptera saussuriana</i>	0.56	-1.22	2.44	-1.03	-1.97	-0.24	-0.36	-1.16	0.41
<i>Myrmeleotettix maculatus</i>	2.86	1.31	4.69	-0.75	-1.56	0	-0.82	-1.73	0.03
<i>Nemobius sylvestris</i>	-2.9	-4.65	-1.48	-1.06	-1.81	-0.38	0.38	-0.07	0.81
<i>Oecanthus pellucens</i>	-3.9	-6.18	-1.89	-0.55	-1.29	0.25	0.25	-0.3	0.8
<i>Oedipoda caerulescens caerulescens</i>	-3.59	-6.24	-1.56	-0.54	-1.3	0.26	-0.74	-1.5	-0.02
<i>Oedipoda germanica</i>	-0.92	-1.88	-0.07	-0.53	-1.09	0.02	-0.46	-0.94	0
<i>Omocestus haemorrhoidalis</i>	0.46	-0.08	1.05	-0.05	-0.52	0.38	-0.45	-0.88	-0.05
<i>Omocestus raymondi raymondi</i>	-1.23	-2.87	0.16	-0.53	-1.26	0.19	0.61	0.02	1.23
<i>Omocestus rufipes</i>	-2.97	-5.38	-0.76	-0.35	-1.08	0.52	0.33	-0.4	1.11
<i>Omocestus viridulus</i>	1.4	0.46	2.51	-0.95	-1.67	-0.33	-0.23	-0.52	0.05
<i>Pezotettix giornae</i>	-4.07	-6.23	-2.1	-0.49	-1.23	0.31	0.44	-0.01	0.91
<i>Pholidoptera aptera</i>	0.68	-1.28	2.83	-1.01	-1.99	-0.21	0.66	-0.04	1.35
<i>Pholidoptera fallax</i>	-2.23	-4.33	-0.43	-0.99	-1.87	-0.26	0.74	-0.02	1.51
<i>Pholidoptera griseoptera</i>	-2.2	-3.98	-0.73	-0.55	-1.29	0.26	0.38	-0.33	1.16
<i>Platycleis albopunctata</i>	-3.02	-4.75	-1.78	-0.72	-1.39	-0.01	-0.19	-0.48	0.11
<i>Podisma dechambrei</i>	1.64	0.74	2.7	-0.58	-1.24	0.04	-0.59	-1.01	-0.21
<i>Podisma pedestris</i>	0.41	-2.09	3.05	-0.91	-1.88	-0.08	0.08	-0.94	1.09
<i>Polysarcus denticauda</i>	0.06	-1.94	2.15	-1	-1.94	-0.21	0.2	-0.81	1.05
<i>Pseudochorthippus parallelus</i>	-0.78	-1.53	-0.07	-0.87	-1.46	-0.37	0.37	0.14	0.6
<i>Psophus stridulus</i>	-0.9	-2.22	0.25	-0.94	-1.75	-0.27	0.25	-0.33	0.78
<i>Roeseliana roeselii</i>	-1.32	-2.53	-0.3	-0.79	-1.47	-0.2	0.42	0.04	0.8
<i>Sepiana sepium</i>	-1.89	-4.36	0.36	-0.82	-1.71	-0.01	0.23	-0.66	1.14
<i>Stauroderus scalaris</i>	-0.47	-1.03	0.06	-0.94	-1.4	-0.53	0.08	-0.07	0.24
<i>Stenobothrus coticus</i>	3.65	1.73	5.87	-0.53	-1.36	0.32	-0.36	-1.21	0.51
<i>Stenobothrus lineatus</i>	-1.05	-1.66	-0.48	-0.63	-1.04	-0.23	-0.2	-0.37	-0.03
<i>Stenobothrus nigromaculatus</i>	2.16	1.16	3.33	-1.03	-1.8	-0.38	-0.83	-1.18	-0.5
<i>Stenobothrus rubicundulus</i>	-1.76	-4.47	0.71	-0.87	-1.77	-0.05	0	-1.12	1.14
<i>Tessellana tessellata</i>	-2.95	-5.33	-0.77	-0.34	-1.07	0.52	0.33	-0.47	1.17
<i>Tettigonia cantans</i>	-1.15	-2.8	0.36	-1.13	-2.11	-0.39	0.73	0.08	1.41
<i>Tettigonia viridissima</i>	-2.17	-4.06	-0.56	-0.98	-1.86	-0.27	0.74	-0.09	1.61
<i>Tylopsis lilifolia</i>	-3.3	-5.87	-1.01	-0.22	-0.99	0.66	-0.26	-1.07	0.56
<i>Yersinella beybienkoi</i>	-2.9	-4.41	-1.57	-0.87	-1.54	-0.23	0.69	0.29	1.11

Table S3: Community-level means and 95% credible intervals for the occupancy and detection parameters on the logit-scale.

	Hyper parameters	Mean	Q2.5	Q97.5
$\mu\alpha_0$	Occupancy at average altitude	-3.09	-4.42	-2.09
$\mu\alpha_1$	Linear effect of altitude	-0.87	-1.47	-0.31
$\mu\alpha_2$	Quadratic effect of altitude	-0.8	-1.04	-0.59
$\mu\beta_0$	Sighting detection at average grass height	-0.62	-1.08	-0.21
$\mu\beta_1$	Listening detection	-2.67	-3.47	-2.01
$\mu\beta_2$	Netting detection	-1.84	-2.18	-1.53
$\mu\beta_3$	Linear effect of grass height on sighting	0.06	-0.13	0.24

Electronic supplementary material 4

It is a link for a web application developed with shinyapps using the software R.

https://bastien-mourguiart.shinyapps.io/shiny_MSOM/

Description: An interactive and dynamic application to visualize species-by-species: the detection map, the altitudinal distribution, the detection probability for each sampling technique, and the effect of grass height on the sighting detection.

2.3 Conclusion

In this chapter, we highlighted the need to account for imperfect detection when studying insect communities. We showed that even relatively highly detectable species, such as Orthoptera, sampled with an efficient sampling design, could have detection probabilities lower than one. Thus, estimates of occupancy probabilities for these species could be biased if imperfect detection is not addressed (Lahoz-Monfort *et al.*, 2014). We then encourage entomologists to systematically use MSOMs when studying distribution of insect communities. To do so, they can use the model developed in the paper. The model can also help to evaluate sampling efficiency and potential optimisation scenarios. Besides, the sampling design presented in the paper could be used in other areas. Thus, this chapter provides a full methodology to study SER of Orthoptera communities in grassland habitats in the presence of partially observed response data.

3 Chapter 2: Modelling species-environment relationships using coarse scale and spatially misaligned environmental data

This work is in preparation for submission to *Journal of applied ecology*: B. Mourguiart, M. Chevalier, M. Marzloff, N. Caill-Milly, F. Ganthy, K. Mengersen, B. Liquet. Modelling fine-scale species-environment relationships using coarse-scale and spatially-misaligned environmental data.

This work also resulted in oral presentations:

- “Determining abundance-environment relationships of manila clam using misaligned environmental data”. Journées du GdR Ecologie Statistique (April 2022)
- “A comparison of hierarchical models to estimate species-environment relationships using spatially misaligned data”. Seminar for the statistics research team of the *Institut de recherche mathématique de Rennes* (May 2022)
- “Modelling abundance-environment relationships using spatially misaligned data” Seminar for the DYNECO research team of Ifremer (June 2022)

3.1 Synopsis

Describing the environment at a coarser resolution than the scale at which the species experiences the environment (i.e., the scale of effect) leads to bias in estimates of species-environment relationships (McInerny & Purves, 2011; Connor *et al.*, 2018). For instance, environmental data are often characterised by average conditions at a coarser resolution than the scale of effect (Potter *et al.*, 2013). In such situations, the coarse-scale covariate only partially describes the environment experienced by the species (McInerny & Purves, 2011). Such an error in covariate leads to bias in SER estimates and should be accounted for (McInerny & Purves, 2011; Foster *et al.*, 2012). Two methods have been proposed in the context of SER modelling: point-level spatial GLMs (Latimer *et al.*, 2006) and Berkson measurement Error Models (BEMs, McInerny & Purves, 2011). However, to our knowledge, the efficiency of these two models has not been thoroughly tested.

In this study, we investigated the effect of partially observed covariate when environmental data are averaged at a larger scale than the scale of effect on the performance of three models (a GLM, a spatial GLM and a BEM). In addition, we evaluated models’ performances regarding their explanatory and predictive power. This assessment was performed through simulations and application to a real case example of the commercially-harvested mollusc (i.e. the Manila clam *Ruditapes philippinarum*) in Arcachon Bay (southwestern France).

3.2 Publication

Modelling fine-scale species-environment relationships using coarse-scale and spatially-misaligned environmental data.

Abstract

- Species distribution models (SDMs) are extensively used in conservation ecology. Given that SDMs estimate species-environment relationships to predict species distribution across space and time, it is key to choose a relevant spatio-temporal scale to study species-environment relationships at the onset of the modelling process. However, this choice is usually constrained by data availability rather than driven by ecological knowledge. For instance, environmental descriptors are often derived from global climate models, which only resolve environmental variability at a given resolution (e.g. 1 km²) that does not necessarily correspond to the scale at which climate influences organisms (i.e. the scale of effect). Thus, spatial (and temporal) misalignment between ecological and environmental data is a common challenge in SDMs. Such a misalignment can bias estimates of species-environment relationships and hence jeopardise the robustness of predicted species distribution.
- We used simulation to study the effect of spatial misalignment on SDMs' performance using covariates averaged across scales that are coarser than the scale of the ecological process. We applied a classic GLM, a spatial GLM, and a GLM with a Berkson error structure in covariates to simulated data. Performances of models were evaluated on species-environment relationship estimates (i.e. response curves), abundance fit, and prediction performance both within and beyond the range of environmental data (i.e., interpolation

and extrapolation). We also fitted the three models to a real case example of the commercially-harvested mollusc (i.e. the Manila clam *Ruditapes philippinarum*) in Arcachon Bay (southwestern France).

- In the simulation, species response curves estimated by the classic and the spatial GLMs were flattened when using coarsened environmental data (i.e. averaged at larger scales than the scale of effect). While the GLM with a Berkson error structure improved estimates of species response curves, it did not improve predictions that were poor for all models. Results from the case study were similar to those obtained with simulations.
- Biases in estimated species-environment relationships and species distribution seem inevitable when environmental data are available at coarser scales than the scale at which ecological processes influence the study system. Thus, estimated species response curves should be interpreted with care if uncertainty exists on the scale at which a species experiences its environment. While Berkson error models performed better than other models in estimating species response curves, further investigations on their accuracy in more complex settings (e.g. more complex settings with additional covariates) are required before advocating for their use for broader application.
- Managers generally require a fine-scale understanding of species-environment relationships and species distribution to guide effective conservation actions. While SDM outputs can theoretically provide such worthy information, it is crucial that predictor variables are available at the scale at which the managed species experienced the environment before using model estimates to support fine-scale management actions. Using a Berkson error model can indicate mis-estimated variability in environmental descriptors and prevent the use of SDM estimates to guide management actions when it occurs.

Introduction

Since their emergence several decades ago (Guisan and Zimmermann 2000), species distribution models (SDMs) have been broadly used to guide management or

conservation actions (Guisan et al. 2013). By estimating species-environment relationships (i.e. response curves) SDMs have made it possible to characterise suitable environmental conditions for a species of interest and hence guide management plans to maintain favourable conditions over a given study area (Greenwood et al. 2016). SDMs have also been used to predict species distribution ranges both under current and future environmental conditions and therefore helped managers prioritise areas for conservation actions (Zurell et al. 2021). For such purposes, SDM species-environment relationships and associated species distribution need to be estimated at fine scales to scope ecologically-relevant management areas (McPherson, Jetz, and Rogers 2006).

So far, SDMs have mostly been fitted with coarse-scale environmental covariates (Austin and Van Niel 2011), with the underlying assumption that species distributions are mainly driven by global or regional climate (Pearson and Dawson 2003). However, recent evidence questioned this hypothesis (Rebaudo, Faye, and Dangles 2016) with a number of studies showing that species distribution can also vary locally, owing to micro-climate variability (Meineri and Hylander 2017; Lenoir, Hattab, and Pierre 2017). For instance, Ashcroft, Chisholm, and French (2009) highlighted the importance of fine-scale climate variability on the distribution of mountainous species. Similar studies conducted in forest systems revealed that canopy protection from sun radiation can buffer the effect of regional climate conditions by acting as a micro-refugia for some species (Zellweger et al. 2020; Stark and Fridley 2022). Hence, only relying on coarsely-resolved environmental variables as covariates for SDMs can produce a mismatch between the scale at which the environment influences a species (i.e., the scale of effect; Chandler and Hepinstall-Cymerman (2016)) and the scale at which the environment is described by covariates (Potter, Arthur Woods, and Pincebourde 2013). Such a mismatch might lead to mis-estimated response curves and lower predictive power of SDMs (Seo et al. 2009; McInerney and Purves 2011), particularly if the scale of effect is finer than the resolution of environmental covariates (Connor et al. 2018).

Ideally, the relevant scale for environmental covariates/predictors in SDM should be defined by ecological knowledge regarding the study species to avoid a mismatch between the scale at which the environment is considered in the model and the

scale of effect (Dormann et al. 2007). However, SDM studies commonly rely on existing species and environmental data, which largely limits the choice of scales. For example, many species data used in SDMs come from previous scientific surveys (Zipkin et al. 2010), citizen science (e.g., GBIF; Faurby and Araújo (2018)), or museum records (Marcer et al. 2012). Associated environmental covariates are often derived from free-to-use datasets, which usually come from climate models Liu et al. (2017) or remote sensing data (Pettorelli et al. 2014). Thus, combining existing data from heterogeneous datasets is likely to induce spatial misalignment (Gotway and Young 2002), which often includes a scale mismatch, between the response variable (i.e., species data) and the covariates (i.e., environmental data). Three types of spatial misalignments can occur in species distribution modelling: 1) point-to-point misalignment when response and explanatory variables are collected at similar scales but in different locations (e.g. using data from nearest weather stations as a proxy for the environmental conditions experienced by ecological communities in the vicinity; Foster, Shimadzu, and Darnell (2012)), 2) point-to-area misalignment when response variables are available at coarser scales than environmental data (e.g. when species data come from museum records available within coarsely-resolved spatial units, such as 5 km by 5 km grid cells; Marcer et al. (2012)), 3) area-to-point misalignment when covariates are observed at coarser scales than the response (e.g., species records collected at multiple georeferenced locations during a scientific survey are associated *a posteriori* to climate data; Latimer et al. (2006)).

Spatial misalignment is well recognized as a crucial issue in SDMs (Martínez-Minaya et al. 2018) and has been shown to induce bias in coefficient estimates and predicted species distribution if not accounted for (Foster, Shimadzu, and Darnell 2012; Stoklosa et al. 2015; Latimer et al. 2006). Several techniques exist to address spatial misalignment between species and environmental data. For instance, predicting environmental conditions in locations where species data is available using interpolation techniques (Foster, Shimadzu, and Darnell 2012), can help overcome point-to-point misalignment. For point-to-area and area-to-point misalignments, heterogeneous datasets can be spatially matched by either upscaling or downscaling the resolution of certain variables (i.e. response and/or covariates)

(Latimer et al. 2006; Keil et al. 2013). Upscaling methods inevitably lead to a loss of information and downscaling might be privileged when the analysis aims to predict fine-scale species distribution. While methods for downscaling species data exist (Keil et al. 2013; McPherson, Jetz, and Rogers 2006), downscaling of environmental data, which is a time-consuming task mostly performed by physicists, is beyond most ecologists' expertise (Hewitson and Crane 1996). Thus, area-to-point misalignment is sometimes overlooked when interest lies in describing fine-scale species distribution, which leads to association of spatially-discrete and distinct species records to identical environmental covariates within coarse environmental spatial cells (Latimer et al. 2006). This "naïve" downscaling may introduce errors in model covariates, which can compromise SDM accuracy (Latimer et al. 2006). For instance, environmental data are often characterised by average conditions at a scale larger than the scale of effect (Potter, Arthur Woods, and Pincebourde 2013). In such situations, the observed environment is a smoothed (less variable) version of the environment experienced by the species (Latimer et al. 2006; McNerny and Purves 2011). McNerny and Purves (2011) showed that such errors in covariate estimates could flatten unimodal species-environment relationships estimated by SDM with consequences on SDM's predictive performance. Two methods have been previously proposed to account for misestimated fine-scale variability in the coarse-scale environmental covariates that could arise from area-to-point misaligned data: point-level spatial GLMs (Latimer et al. 2006) and Berkson error models (BEMs, McNerny and Purves (2011); Martínez-Minaya et al. (2018)).

While BEMs and spatial GLMs can be used to account for area-to-point spatial misalignment, they do so in a different way. BEMs estimate species response curves by fitting the response variable to an unobserved (error-free) covariate that is assumed to be more variable than the available (error-prone) covariate (see, e.g., Muff et al. (2015) for more details). Alternatively, spatial GLMs include a random spatial effect that accounts for the fine-scale variability that is unexplained by coarse descriptors (Latimer et al. 2006). To our knowledge, the efficiency of these two models in accounting for area-to-point misalignment in species-environment relationship estimates has not been thoroughly tested. Indeed, although McNerny and Purves (2011) investigated the ability of BEM for estimating

species-environment relationships, they assumed that the variance error between the unobserved error-free and the observed error-prone environment was known. This information is however unavailable for most climatic datasets. Regarding spatial GLMs, Latimer et al. (2006) advocated for their use to account for area-to-point misalignment but without providing any external evidence of accuracy or predictive power of this method.

In this study, we investigated the effect of area-to-point misalignment when environmental data are averaged at a larger scale than the scale of effect on the performance of three alternative SDMs: a GLM (a frequently used SDM; Norberg et al. (2019)), a spatial GLM and a BEM (two potential solutions to address area-to-point misalignment; Latimer et al. (2006); Martínez-Minaya et al. (2018)). This assessment was performed both through simulations and application to a real case study.

Materials and methods

Data structure

Area-to-point misalignment describes the case where we observe a response variable (e.g., species counts) at n spatial point locations, while associated explanatory variables (e.g., environmental descriptors) are available at a coarser scale, typically across a grid of I cells that can contain multiple sampling points. Importantly, environmental variability is usually neglected within each grid-cell (i.e. all points are assumed to have the same value within a given cell), potentially leading to an area-to-point misalignment problem (i.e., a mismatch between the species point-scale and the covariate grid-scale). We denote $Y_{j(i)}$ the observed count at sampling point $j(i)$ which is included in grid-cell i , for $i = 1, \dots, I$ and $j(i) = 1, \dots, J_i$, with $\sum J_i = n$. \mathbf{W}_i represents a vector of misaligned explanatory environmental variables associated with grid-cell i . $X_{j(i)}$ is the value of the environmental variable at sampling point $j(i)$.

Modelling framework

Three models were considered: a Generalized Linear Model (GLM), a spatial Generalized Linear Model (sp-GLM), and a Berkson Error Model (BEM). All models assume that counts of the target species (e.g. clams in our case study) $Y_{j(i)}$ at sampling point $j(i)$ included in cell i relies on the expected abundance $\lambda_{j(i)}$: $Y_{j(i)} \sim f(\lambda_{j(i)}, \phi)$, where $f(\cdot)$ is the “Poisson” or the “Negative Binomial” probability distribution function with mean parameter $\lambda_{j(i)}$ and variance parameter ϕ (for the Negative Binomial). We describe the observation data as an outcome of a Poisson trial in the simulation. In the case study, we chose a Negative Binomial distribution to model overdispersed counts. We assumed that expected abundance is related to the environment through different formulations depending on the model. We fitted all models using scaled covariates (with mean 0 and standard deviation 1). In the following, we present the model formulations using a unique misaligned covariate for simplicity, but this can be extended to more covariates.

The GLM considers that the environment at the grid-scale, W_i , is the only driver of variation in expected abundance which is measured at the point-scale. It assumes that the grid is the environmental scale of effect on the species. Thus the species-environment relationship is modelled as follows:

$$\log(\lambda_{j(i)}) = \beta_0 + \beta_1 W_i + \beta_2 W_i^2$$

where β_0 is the expected abundance in log-scale in average environmental conditions (i.e., when the scaled covariate is null), β_1 and β_2 are the coefficients representing respectively the linear and the quadratic effects of the grid-scale environmental covariate W . Note that we used a log-link function as we modelled counts, but other link functions can be used depending on the type of species data (e.g. a logit-link function can be used for presence-absence data).

The spatial GLM is a mixed model that incorporates a spatial random effect, $\gamma_{j(i)}$, that allows variability between spatial points within a given grid-cell. It assumes that expected abundance not only depends on the grid-scale environment (as for the classic GLM) but also varies depending on a fine-scale latent spatial field. This spatial random effect allows capturing spatial signals not explained

by the predictors (Zurell et al. 2021). For example, unobserved spatial patterns could result from missing spatially coherent or biological predictors (e.g., dispersal ability). Here, the spatial random effect is supposed to capture variability within environmental grid cells. The model is written as follows:

$$\begin{aligned} \log(\lambda_{j(i)}) &= \beta_0 + \beta_1 W_i + \beta_2 W_i^2 + \gamma_{j(i)} \\ \boldsymbol{\gamma} &\sim MVN(\mathbf{0}, \boldsymbol{\Sigma}) \end{aligned}$$

where $\boldsymbol{\gamma}$ is a vector of dimension n , $\boldsymbol{\Sigma}$ is the spatial covariance structure whose generic element is $\Sigma_{u,v} = \sigma_\gamma^2 \times \text{Matern}(d_{u,v}, \kappa)$ where σ_γ^2 is the variance component, $\text{Matern}(\cdot, \cdot)$ is the Matérn function which describes how the correlation between two points (here u and v) decreases with the Euclidean distance separating them (denoted $d_{u,v}$), and where κ is a scaling parameter related to the spatial range r , i.e., the distance at which the spatial correlation becomes almost null. Generally, the range is defined as the distance at which the spatial correlation is close to 0.1. It could be derived from κ by: $r = \frac{\sqrt{8\nu}}{\kappa}$, with ν representing the degree of smoothness of the spatial process and usually fixed to one (Zuur, Ieno, and Saveliev (2017), p197).

The BEM jointly estimates the covariate and abundance at the point-scale. It considers that the observed environment at the grid-scale W_i is a smoothed version of the environment at the point-scale $X_{j(i)}$, which is assumed to be the “true” driver of abundance. This latent variable is modelled as follows:

$$X_{j(i)} \sim N(W_i, \sigma_X^2)$$

with σ_X^2 the variance parameter that describes the fine-scale variability lost by averaging the environment at grid-scale. The expected abundance modelled as a function of the latent variable X :

$$\log(\lambda_{j(i)}) = \beta_0 + \beta_1 X_{j(i)} + \beta_2 X_{j(i)}^2.$$

Parameter estimation

We fitted the GLM and the BEM within a Bayesian framework using MCMC sampling with the software JAGS (Plummer 2003) and the R package jagsUI Team (2018). We ran three chains for each analysis with a burn-in of 20,000 and an additional 20,000 iterations with a thinning rate of 20. For prior distributions of parameters β_1 and β_2 , we used normal distribution with a zero mean and a precision of 0.1. For the intercept, β_0 , we specified a prior distribution on a derived ecologically meaningful parameter, $\lambda^* = \exp(\beta_0)$, representing the expected abundance in average environmental conditions. We used a uniform prior bounded between 0 and 50 for this derived parameter, assuming that expected abundance in average environmental conditions should lay between those values. We used inverse gamma priors with shape 0.1 and rate 0.1 for precision variances. When we used the Negative Binomial distribution to model abundances for the case study, we set a uniform prior bounded between 0 and 50 for the dispersion parameter ϕ . We assessed convergence by examining the Gelman-Rubin statistic (\hat{R}) with a threshold fixed to 1.1 (Gelman et al. 2013). We fitted the spatial GLM using the INLA and SPDE approaches which were applied using the R package R-INLA (Rue, Martino, and Chopin 2009). We did not use MCMC sampling due to the known “big n problem” in spatial analysis that prevents the use of JAGS when sample size (i.e. number of locations) exceeds 100 (Kéry and Royle 2020). We kept the default priors specified by INLA for this model.

Simulation study

We tested the effect of spatial misalignment on models’ explanatory and predictive performances using a virtual species approach where abundance distribution was explained by one unique covariate. The simulation design, represented in Fig. 1, consisted of four steps: (1) simulate an ecological process with a known species-environment relationship at the point-scale, (2) simulate an observation process with a measured environment at the grid-scale representing averaged point-scale environmental conditions at different spatial resolutions, (3) fit three models to the different simulated datasets, and (4) assess the ability of models to explain, interpolate and extrapolate simulated abundances.

We simulated a grid of 2000 x 2000 points representing a virtual sampling area. For simplicity, we assumed that sampling points were separated by one meter but any type of distance can be considered. We simulated the virtual environment as a Gaussian spatial random field using a Matérn covariance matrix to represent spatial dependencies between points. We set the spatial range and variance parameters of the Matérn function to 40 meters and 1 respectively. From those point-scale environmental values we calculated virtual expected abundance at the point-scale using a quadratic linear relationship: $\lambda_j = \beta_0 + \beta_1 X_j + \beta_2 X_j^2$. We chose the coefficients ($\beta_0 = 3$, $\beta_1 = 1$, $\beta_2 = -1.5$) of the species-environment relationship to describe a sharp bell-shaped species response curve. We then simulated the observed counts, N_j , using random draws from a Poisson distribution with parameter λ_j .

In the virtual observation process, we averaged the point scale environment at four horizontal grid sizes (10 m, 20 m, 40 m, 80 m) to obtain the misaligned grid-scale environment. We defined the environmental spatial resolutions as multiplications (0.25, 0.5, 1, 2) of the environmental spatial range that we fixed at 40 m earlier in the simulation. By doing so, we assumed that the effect of misalignment depends on the ratio between the environmental spatial heterogeneity and the environmental spatial resolution rather than the absolute spatial resolution. We based this assumption on previous studies that found a positive relationship between the magnitude of spatial misalignment effect on regression models and spatial autocorrelation in covariates (Gotway and Young 2002; Naimi et al. 2014). We then virtually sampled the simulated survey area by randomly selecting 100, 300, or 500 points. We simulate no bias in species sampling (i.e., perfect detection). Observed counts were thus equal to the simulated abundance at the sampling points. For each sampling point, we allocated the environmental value of the grid containing the point to represent the observed coarse-scale environment. Finally, we replicated the virtual sampling 10 times, resulting in 120 simulated datasets (4 spatial resolutions x 3 sampling sizes x 10 replications).

Model performance assessment

To assess the accuracy of the three models, we fitted them to each of the 120 simulated datasets. We then compared the estimated species-environment relationships

and fitted counts to the true values (those simulated). We also assessed the models' predictive ability with regards to interpolation and extrapolation using a data splitting strategy (appendix 1 in Roberts et al. (2017)). First, we estimated model parameters for each of the 120 simulated datasets (train datasets). Then, we used two test datasets, one for interpolation (predicting in the same environment as the train dataset) and one for extrapolation (predicting in another environment than the train dataset). For the interpolation test dataset, we virtually sampled 100 new points for each replication in the same environment that the one used to fit the models. We then simulated virtual abundance at those 100 new points using the same species-environment relationship as for the train dataset. For the extrapolation test dataset, we applied the same scheme but environmental values at the 100 new points were increased by 1.5 relative to the interpolation test dataset to simulate a new environment. For both interpolation and extrapolation test datasets, we used the species-environment relationships estimated by the models using the train datasets to predict the abundance in the 100 new points using observed grid-scale environmental values at those new sampling points. We then compared predicted abundance ($\hat{\lambda}$) to the true simulated abundance (λ) by calculating root-mean-square errors ($RMSE = \sqrt{\frac{1}{R} \sum_{r=1}^R (\lambda - \hat{\lambda}_r)^2}$, with R the number of replications) and Spearman correlation coefficients. We chose the Spearman correlation coefficient to investigate the models' ability to accurately predict abundance patterns.

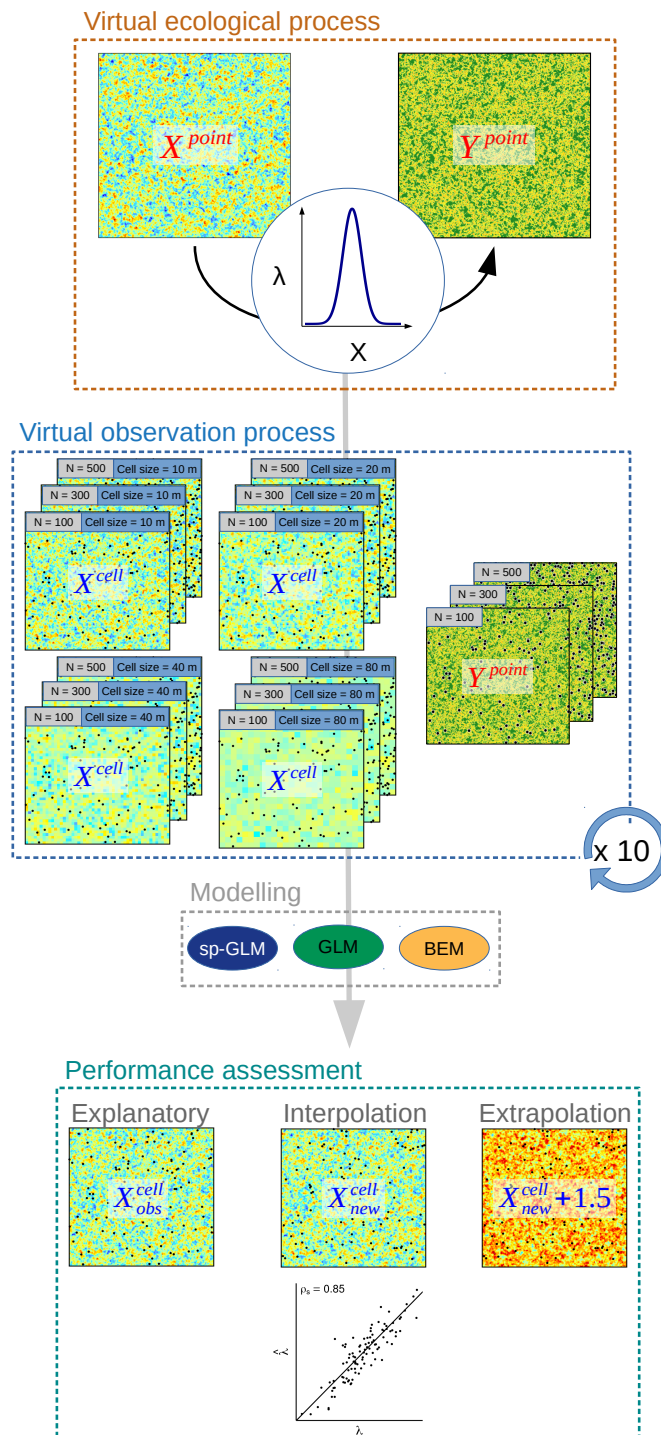


Figure 1: Representation of the simulation design used to investigate the effect of area-to-point spatial misalignment on models explanatory and predictive performances of a generalized linear model (GLM), a spatial GLM (sp-GLM) and a Berkson error model (BEM).

Case study

As a complement to the virtual simulation-based study, we applied the three models to a real case study. We estimated abundance-environment relationships of Manila clam, *Ruditapes philippinarum* (Adams & Reeve, 1980), in Arcachon Bay (SW of France). We used count data collected at 491 sampling points in June 2008 and 466 sampling points in June 2018 to fit and test the three models. We splitted the 2008 dataset into train and test datasets containing respectively 393 and 98 sampling points (75% and 25% of the sampling size). Models were fitted on the train dataset. We then used parameter estimates to predict abundance on the 2008 test dataset (to assess interpolation accuracy) and on the full 2018 dataset (to assess extrapolation accuracy).

Sampling points were spatially distributed within the study area following a generalized random tessellation stratified sampling design which is spatially balanced and particularly appropriate for patchy species such as clams (Kermorvant et al. 2019). Clams were captured by a Hamon grab, collecting a sediment core of 0.25 m² on a 0.2 m depth. We assumed no sampling bias as Manila clams bury no more than 0.12 m depth. Individuals captured were counted and measured. We only kept individuals longer than 30 mm in the counts, removing individuals that fisheries could target. We chose the threshold length to be smaller by 5 mm than the current regulatory catch size.

As the goal is to investigate the misalignment effects on models performance rather than select environmental drivers of Manila clam distribution, we chose to pre-select three covariates assumed to influence Manila clam's abundance (Tezuka et al. 2013; Yin et al. 2017; Bae et al. 2021): water salinity (psu), water temperature (°C) and immersion time (h.day⁻¹). We extracted those variables from the MARS3D numerical model (Lazure and Dumas 2008). This hydrodynamic model was implemented on Arcachon Bay by Plus et al. (2009) and improved by Kombiadou et al. (2014). The numerical model used an empirically-based tide (FES2012 solution, Carrere et al. 2013) to capture the natural variability of tidal forcing. Freshwater inputs were set as constants equal to yearly average values. We used a 2-dimensional version with a horizontal resolution of 235 m. A

fourteen-month simulation was performed before the survey, e.g., between April 2007 and May 2008 for the 2008 dataset, a period corresponding to the life span of larvae and juvenile phases of Manila clam (Caill-Milly 2012).

Results

Simulation study

Sample size only had a marginal effect on models' explanatory performances and almost no impact on models' predictive performance (Appendix 1). We only observed a slightly lower bias in species-environment relationship estimates with increasing sampling size, but that virtually did not modify the relative effect of area-to-point misalignment on alternative models' accuracy. Thus, we here only present the results for the intermediate scenario with 300 simulated sampling sites. Extended results with all sampling scenarios can be found in Appendix 1.

Estimation of species-environment relationships (SER)

All three models showed lower accuracy in Species-Environment Relationship (SER) estimates with increasing spatial scale (i.e. decreasing resolution) in environmental covariates, but responses varied across models (Fig. 2). SER estimated by the GLM were increasingly flattened with increasing spatial scale of covariates compared to the true curve, i.e. with lower maximums (see Appendix 2 for results on maximum estimates) and larger widths (Fig. 2 second row and Appendix 2). Similar effects of area-to-point misalignment were observed for the spatial-GLM (Fig. 2 third row). In the most severe scenario, where the observed scale (i.e., 80 m) is twice the spatial range of the simulated environmental variable (i.e., 40 m), both GLMs (classic and spatial) failed to retrieve the correct shape of species response in some replications, estimating U-shaped or exponential curves (Fig. 2, fourth column of the second and third rows). The BEM had much better accuracy, with almost no difference between estimated and simulated curves in the first two scale scenarios (Fig. 2, first row of the first and second columns) and slight departures from the simulated optimum in the most two severe scenarios (i.e., optimum positions biased toward zero; see Appendix 2 for results on bias in optimum estimates).

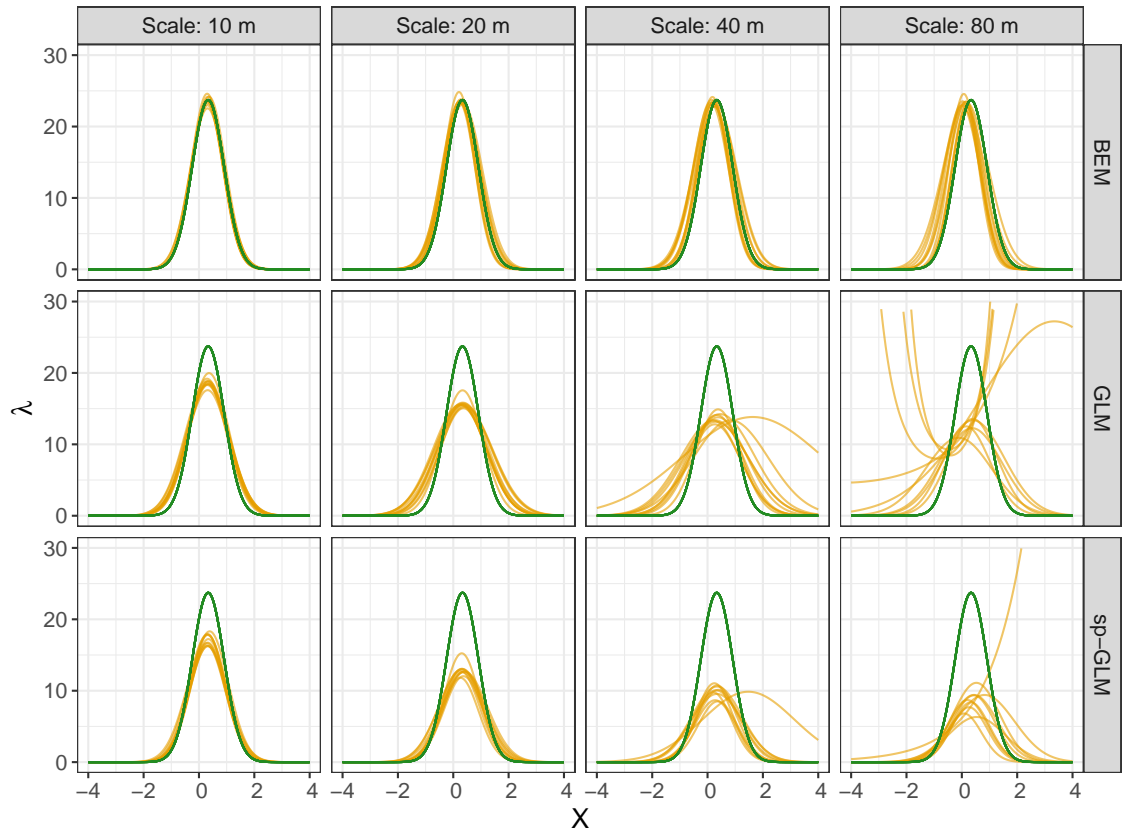


Figure 2: Species response curves estimated by the three models for the ten replications (orange curves) compared to simulated species response curves (green curves) in four scenarios with different spatial scales at which environmental covariates were used.

Abundance predictions

GLM explanatory power was lower than the two other models for all the simulated spatial scales (Table 1). Even under the highest 10 m grid-size resolution (i.e., the lowest misalignment), GLM had much lower explanatory performance than the two other models with, for instance, an average RMSE 2.5 and 1.9 times higher than for the BEM and the spatial GLM, respectively. GLM fit accuracy decreased rapidly with increasing grid-size (i.e., lowering spatial resolution of covariates) from 10 m to 80 m. Predicted abundance became almost uncorrelated to the observed counts for the third scale scenario with an average correlation coefficient of 0.33 (SD = 0.07). In contrast, the explanatory power of the BEM and the spatial GLM was only marginally affected by coarsening spatial resolution of descriptors. Their average RMSE remained almost constant across all scale scenarios, ranging from 1.79 to 1.93 for BEM and from 2.30 to 2.48 for spatial GLM, while the correlation

coefficients with observed abundance constantly being around 0.94.

The three models were very sensitive to the increase of spatial scales regarding the accuracy of their predictions for both interpolation and extrapolation (Table 1). Surprisingly, predictive performances were slightly better but still low for extrapolation datasets, except in the coarsest 80 m scale scenario for the GLM and the spatial GLM. This last result is related to the estimation of incorrect shapes of SER under the GLM and the spatial GLM in some replications scenario (Fig. 2). Indeed, both models estimated SER predicting unrealistic high counts at both extremes of the environmental gradient.

Table 1: RMSE and Spearman correlation coefficients between predicted and simulated counts for the three models on the three types of data used for prediction depending on the scale of observed environment.

Performance metric	Data type	Model	Scale size			
			10 m	20 m	40 m	80 m
RMSE	Train	BEM	1.79 (0.1)	1.9 (0.14)	1.92 (0.14)	1.93 (0.14)
		GLM	4.49 (0.22)	6.09 (0.29)	7.08 (0.28)	7.51 (0.26)
		sp-GLM	2.3 (0.09)	2.43 (0.15)	2.47 (0.16)	2.48 (0.16)
	Test (interpolation)	BEM	4.4 (0.43)	6.52 (0.55)	8.37 (0.51)	9.63 (0.69)
		GLM	4.63 (0.27)	6.26 (0.32)	7.26 (0.21)	7.54 (0.2)
		sp-GLM	4.49 (0.43)	6.58 (0.48)	8.08 (0.67)	8.33 (0.68)
	Test (extrapolation)	BEM	3.18 (0.5)	4.75 (0.45)	5.86 (0.66)	6.21 (0.83)
		GLM	3.43 (0.51)	5.01 (0.51)	6.73 (1.22)	28.4 (34.36)
		sp-GLM	3.46 (0.48)	5.14 (0.65)	8.09 (4.83)	28.88 (50.58)
Spearman correlation	Train	BEM	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)
		GLM	0.76 (0.03)	0.55 (0.05)	0.33 (0.07)	0.18 (0.06)
		sp-GLM	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)
	Test (interpolation)	BEM	0.74 (0.05)	0.51 (0.08)	0.26 (0.09)	0.11 (0.11)
		GLM	0.74 (0.05)	0.51 (0.08)	0.25 (0.08)	0.16 (0.08)
		sp-GLM	0.73 (0.05)	0.5 (0.07)	0.24 (0.08)	0.16 (0.1)
	Test (extrapolation)	BEM	0.86 (0.03)	0.71 (0.03)	0.5 (0.04)	0.31 (0.08)
		GLM	0.86 (0.03)	0.71 (0.03)	0.43 (0.17)	0.02 (0.33)
		sp-GLM	0.86 (0.02)	0.7 (0.05)	0.38 (0.31)	0.09 (0.28)

Case study

We observed similar patterns as in species-environment relationships estimated by the three models on Manila clam counts as those described above in the simulation context as both GLMs estimated flatter SER relative the BEM (Fig. 3). Still, all models predicted that Manila clam abundance is unimodally related to the three covariates with an optimum abundance estimated around an immersion time of 12 hours, a mean salinity of 31 psu and a mean temperature of 14.4 °C.

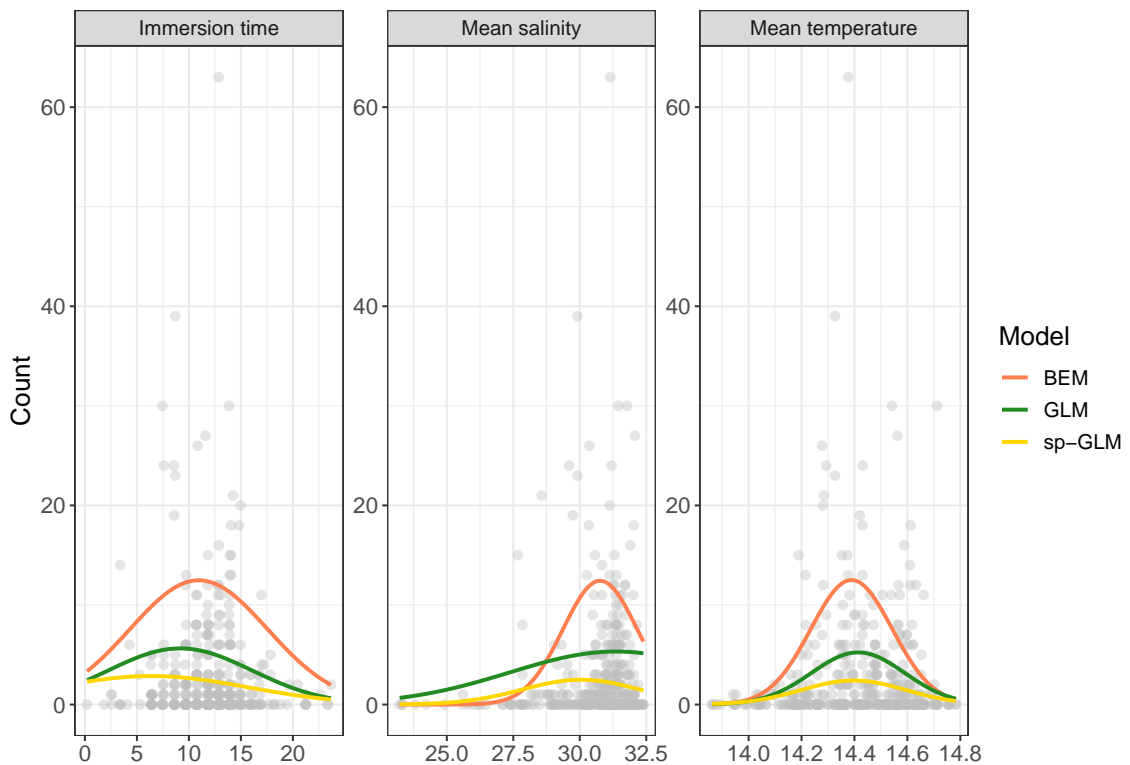


Figure 3: Estimates of abundance-environment relationships of the three models used (GLM, spatial GLM and Berkson error model) to fit Manila clam counts collected at 437 sampling sites during spring 2008 with three covariates obtained from hydrodynamical model MARS3D at an horizontal resolution of 235 m in Arcachon Bay.

Case study results of explanatory or predictive (interpolation and extrapolation) performances of the three models aligned with those obtained in the simulation study (Table 2). The BEM outperformed the two GLMs in explanatory power, and in particular the classic GLM that showed the lowest power. However, the three models had very low predictive power, especially regarding extrapolation given that, predicted counts were almost unrelated to observed counts (Table 2).

Table 2: RMSE and Spearman correlation coefficients between predicted and observed counts of Manila clam for the three alternative models, namely GLM, spatial GLM (sp-GLM) and Berkson Error Model (BEM), as estimated using the three different datasets, namely the training 2008 dataset, the test 2008 dataset (spatial interpolation) and the 2018 dataset (space and time extrapolation).

Data type	Model	RMSE	Spearman correlation
Train (2008)	BEM	1.84	0.91
	GLM	3.47	0.25
	sp-GLM	2.40	0.68
Test (2008)	BEM	5.00	0.30
	GLM	3.07	0.26
	sp-GLM	2.87	0.39
Test (2018)	BEM	7.69	0.22
	GLM	6.90	0.19
	sp-GLM	6.76	0.33

Discussion

Summary of major results

Species response curves (i.e., estimates of Species-Environment Relationships; or SER hereafter) are flattened in GLM and spatial GLM when environmental data is available at a coarser scale than the scale of effect. This bias increases with decreasing environmental data resolution (i.e. when spatial misalignment increases). Relative to the GLM and the spatial GLM, the Berkson Error Model (BEM) more accurately estimates SERs even in the worst case scenarios, despite a slight bias optimum conditions estimating (i.e. the response curve is slightly shifted) when the spatial resolution of environmental covariates is too coarse. All three models' predictive performances (in terms of both interpolation and extrapolation) rapidly decrease with coarsening of environmental data spatial resolution, and become poor in the worst case scenario.

Regression dilution precludes the interpretation of model coefficients and variable importances in GLMs

As expected, in our simulations, the area-to-point misalignment led to a diminution of heterogeneity in covariates, thus to an erroneously smoothed environment (Appendix 3). This error is known as the regression dilution problem and leads

to flattened unimodal SER estimated by GLMs (McInerny and Purves 2011). We highlighted in the simulations that adding point-level spatial random effects that allow for fine-scale heterogeneity in species response does not solve but can actually worsen the regression dilution problem. Our results may provide some insights on previously-observed differences between GLMs and spatial GLMs fitted to area-to-point misaligned data (Latimer et al. 2006). Latimer et al. (2006) correlated fine-scale plant species distribution in the Cape Floristic Region of South Africa to coarse-scale environmental factors in a GLM and a point-level spatial GLM. They obtained a better fit relying on fewer significant variables with a spatial GLM than with a classic GLM, which they interpreted as a better and larger estimation of coefficient variance in spatial GLMs (Dormann et al. 2007). However, our results suggest that this could also result from the attenuation effect of area-to-point misalignment which is stronger in spatial GLMs than in GLMs. For instance, they found differences between GLM and spatial GLM in the significance of model coefficients for edaphic variables (i.e. soil properties of samples). Those variables probably act at fine scales and thus may be subject to area-to-point misalignment inducing a more significant attenuation effect in spatial GLM than in GLM, which can decrease the effect size of the edaphic variables. Comparison of relative covariate effects based on a unique model can be misleading when exposed to area-to-point misalignment in variables. For instance, a (spatial) GLM including two covariates A and B observed at the same scale, which corresponds to the scale of effect of covariate A but is coarser than the scale of effect of covariate B, might underestimate the effect size of the covariate B. Thus, when the ecological process varies at fine scales, the consideration of coarse environmental covariates to model fine-scale species data should be avoided not only in GLMs, as previously highlighted (McInerny and Purves 2011), but also in spatial GLMs, a result that contradicts previous findings (Latimer et al. 2006).

Berkson error model as a potential solution to account for area-to-point misalignment when estimating fine-scale SER.

The BEM appeared as an interesting candidate model to overcome bias in SER estimates induced by area-to-point misalignment in GLMs. When applied to both

simulated data and a real-world case study, the BEM outperformed the GLM and the spatial GLM in terms of fit. Simulations also highlighted that the BEM, conversely to the two GLMs, accurately estimated fine-scale SER shape using coarse environmental data. The ability of the BEM to accurately estimate fine-scale SER using smoothed error-prone covariate has already been highlighted in particular settings where error variance was known (McInerny and Purves 2011). We showed in our simulation study that it is not necessary to specify the error variance in BEM to estimate SER accurately. This result could significantly improve the range of applications of BEM. Note, however, that in the BEM we developed, we assumed that error variance within a given cell is constant across grid cells and depends on the variability between coarse-grid cells, with the underlying assumption of spatial stationarity in covariates (Dormann et al. 2007). This assumption may not always be reasonable, especially in broad-scale studies. Further research is needed to investigate the effect of departure from this assumption. Potential solutions around this may include incorporating local environmental values to describe cell-specific variance error between coarse and fine-scale environments.

In the Manila clam case study, differences observed between estimates of GLMs and BEM could be related to area-to-point misalignment, as observed in the simulation study, but they may also flow from other sources of variability unexplained by the covariates. For instance, in our application, we found better explanatory performance for BEM than spatial GLM, which contrasts with simulation-based results that revealed equivalent goodness-of-fit for both. Hence, we could have expected comparable explanatory performances in BEM and spatial GLM for the Manila clam application if area-to-point misalignment was the only unobserved source of uncertainty. However, other sources of heterogeneity induced by non-spatial patterns and not described by the covariates can affect Manila clam distribution and thus explain the better performance of BEM. Temporal misalignment can, for instance, lead to unobserved variability in covariates. We chose to use summary statistics of environmental conditions over a period of fourteen months before the survey according to the life span of larvae and juvenile phases potentially affected by the environment (Tezuka et al. 2013). However, as juveniles and adults have a relatively high environmental tolerance (Dang 2009), we could have assumed

environment filtering on the larval phase only. Thus, using summarised environmental conditions over a longer period may have induced temporal misalignment. Further research should therefore assess how BEM responds to multiple sources of variability unexplained by covariates. In addition, further work is required to investigate if model selection between GLMs and different formulations of BEMs can be used to discriminate environmental variables that have a finer scale of effect than the observed scale. We showed that BEM better fitted the data if one variable was assumed to explain fine-scale species distribution and was observed at a coarser scale than its scale of effect. Does this result stand when considering more than one variable? For instance, if two effective descriptors of species distribution are observed at the same scale, which corresponds to the effect scale of one covariate but is coarser than the scale of effect of the other; then, does a BEM modelling one error-free and one error-prone covariate, better fit such data than a misaligned GLM or a BEM assuming two error-prone covariates? Will the misspecified BEM with two error-prone covariates accurately estimate a zero variance error for the error-free covariate, or will it share the variance error induced by the error-prone covariate between the two covariates? Answers to those questions are required before interpreting the case study results. If error variance is shared between covariates, even if there is only one error-prone covariate, our results could be misleading. Hence, further research on BEM behavior, especially regarding various sources of unexplained variability and multiple scales of effects in the descriptors, is needed before advocating for the broad use of BEMs in species distribution modelling.

Low predictive power may limit the applicability of BEM in SDM

Another salient problem limiting BEM's use in SDM is its low predictive power. While improving estimates of species-environment relationships is interesting, planning conservation actions also needs accurate predictions of species distribution (Zurell et al. 2021). Unfortunately, neither of the two candidate models (BEM and spatial GLM) fulfilled this task under area-to-point misalignment, and GLM performed even worse. While it is not surprising that a GLM trained with coarse data failed to predict fine-scale species distribution, we expected the spatial GLM

to perform better, at least for interpolation in simulation-based examples. However, the spatial GLM over-fitted the data and failed to estimate the environmental spatial pattern. Further investigations will be performed before submission of this chapter to a peer-reviewed journal to understand the behaviour of the spatial GLM. It is easier to understand why BEM and non-spatial GLM produced inaccurate predictions. Inaccurate predictions by the BEM and the non-spatial GLM are likely due to averaging of environment values at the grid-scale. This produced a shrinkage in the observed environmental gradient, where grid-scale observed values become concentrated around the mean value of the gradient. Hence, the range of predictions made by GLM and BEM with coarse-scale environmental data is also shrunk around the estimated count at the environmental mean (i.e. $exp(\beta_0)$, see Methods and Appendix). Solutions may be developed to improve predictions in BEM. For instance, instead of using the environment available at the grid-scale, we could have used the posterior distribution of the spatial point environmental variable with the closest coarse environmental value in the training dataset.

Conclusions - Perspectives

Managers may require fine-scale information on species-environment relationships and species distribution (Guisan et al. 2013). While, due to data availability constraints, it is tempting to allocate coarse environmental variables to describe fine-scale distributional data, this widely used practice should be avoided. Allocating coarse-scale environmental variables might lead to bias in estimates of fine-scale species-environment relationships and predictions of species distribution. At the moment, SDMs should only be used to inform conservation and management actions if evidence demonstrates that descriptors were observed at their scale of effect. Further research is needed to define those scales of effect. New advances in environmental data collection techniques, e.g. remote sensing, may help by in this regard by increasing the availability of fine-scale environmental data (Lembrechts, Nijs, and Lenoir 2019). When exposed to a scale mismatch between the observed environment and its scale of effect, the Berkson error model offers a better alternative to GLMs and deserves further investigation to study the extent to which scale mismatches can affect SDMs inferences.

Acknowledgments

B.M. was supported by the UPPA-E2S international chair of Kerrie Mengersen. Collaboration between B.M., M.M. and M.C. has been possible thanks to the Collaborative Hub for the Extension of Coastal Knowledge (CHECK) of the Ifremer's unit DYNECO.

References

- Ashcroft, Michael B., Laurie A. Chisholm, and Kristine O. French. 2009. "Climate Change at the Landscape Scale: Predicting Fine-Grained Spatial Heterogeneity in Warming and Potential Refugia for Vegetation." *Global Change Biology* 15 (3): 656–67. <https://doi.org/10.1111/j.1365-2486.2008.01762.x>.
- Austin, Mike P., and Kimberly P. Van Niel. 2011. "Improving Species Distribution Models for Climate Change Studies: Variable Selection and Scale." *Journal of Biogeography* 38 (1): 1–8. <https://doi.org/10.1111/j.1365-2699.2010.02416.x>.
- Bae, Hyeonmi, Jibin Im, Soobin Joo, Boongho Cho, and Taewon Kim. 2021. "The Effects of Temperature and Salinity Stressors on the Survival, Condition and Valve Closure of the Manila Clam, *Venerupis Philippinarum* in a Holding Facility." *Journal of Marine Science and Engineering* 9 (7). <https://doi.org/10.3390/jmse9070754>.
- Caill-Milly, Nathalie. 2012. "Relations Entre l'état d'une Ressource Et Son Exploitation via La Compréhension Et La Formalisation Des Interactions de Socio-écosystèmes. Application à La Palourde Japonaise (*Venerupis Philippinarum*) Du Bassin d'arcachon." PhD thesis, Université de Pau et des Pays de l'Adour.
- Carrere, Loren, Florent Lyard, M Cancet, A Guillot, and Laurent Roblou. 2013. "FES 2012: A New Global Tidal Model Taking Advantage of Nearly 20 Years of Altimetry." *20 Years of Progress in Radar Altimetry* 710: 13.
- Chandler, Richard, and Jeffrey Hepinstall-Cymerman. 2016. "Estimating the Spatial Scales of Landscape Effects on Abundance." *Landscape Ecology* 31 (6): 1383–94. <https://doi.org/10.1007/s10980-016-0380-z>.
- Connor, Thomas, Vanessa Hull, Andrés Vina, Ashton Shortridge, Ying Tang,

- Jindong Zhang, Fang Wang, and Jianguo Liu. 2018. "Effects of Grain Size and Niche Breadth on Species Distribution Modeling." *Ecography* 41: 1270–82. <https://doi.org/10.1111/ecog.03416>.
- Dang, Cécile. 2009. "Dynamique Des Populations de Palourdes Japonaises (Ruditapes Philippinarum) Dans Le Bassin d'arcachon: Conséquences Sur La Gestion Des Populations Exploitées." PhD thesis, Bordeaux 1.
- Dormann, Carsten F., Jana M. McPherson, Miguel B. Araújo, Roger Bivand, Janine Bolliger, Gudrun Carl, Richard G. Davies, et al. 2007. "Methods to Account for Spatial Autocorrelation in the Analysis of Species Distributional Data: A Review." *Ecography* 30 (5): 609–28. <https://doi.org/10.1111/j.2007.0906-7590.05171.x>.
- Faurby, Søren, and Miguel B. Araújo. 2018. "Anthropogenic Range Contractions Bias Species Climate Change Forecast." *Nature Climate Change* 8: 252–56. <https://doi.org/10.1038/s41558-018-0089-x>.
- Foster, Scott D., Hideyasu Shimadzu, and Ross Darnell. 2012. "Uncertainty in Spatially Predicted Covariates: Is It Ignorable?" *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 61 (4): 637–52. <https://www.jstor.org/stable/23251164>.
- Gelman, Andrew, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. 2013. *Bayesian Data Analysis*. Third edition. New York: Chapman; Hall.
- Gotway, Carol A., and Linda J. Young. 2002. "Combining Incompatible Spatial Data." *Journal of the American Statistical Association*, no. 97:458: 632–48. <https://doi.org/10.1198/016214502760047140>.
- Greenwood, Owen, Hannah L. Mossman, Andrew J. Suggitt, Robin J. Curtis, and Ilya M. D. Maclean. 2016. "Using in Situ Management to Conserve Biodiversity Under Climate Change." *Journal of Applied Ecology* 53 (3): 885–94. <https://doi.org/10.1111/1365-2664.12602>.
- Guisan, Antoine, Reid Tingley, John B. Baumgartner, Ilona Naujokaitis-Lewis, Patricia R. Sutcliffe, Ayesha I. T. Tulloch, Tracey J. Regan, et al. 2013. "Predicting Species Distributions for Conservation Decisions." *Ecology Letters* 16 (12): 1424–35. <https://doi.org/https://doi.org/10.1111/ele.12189>.

- Guisan, Antoine, and Niklaus E Zimmermann. 2000. "Predictive Habitat Distribution Models in Ecology." *Ecological Modelling* 135 (2-3): 147–86.
- Hewitson, B. C., and R. G. Crane. 1996. "Climate Downscaling: Techniques and Application." *Climate Research* 07 (2): 85–95. <https://doi.org/10.3354/cr007085>.
- Hijmans, Robert J., Susan E. Cameron, Juan L. Parra, Peter G. Jones, and Andy Jarvis. 2005. "Very High Resolution Interpolated Climate Surfaces for Global Land Areas." *International Journal of Climatology* 25 (15): 1965–78. <https://doi.org/10.1002/joc.1276>.
- Keil, Petr, Jonathan Belmaker, Adam M. Wilson, Philip Unitt, and Walter Jetz. 2013. "Downscaling of Species Distribution Models: A Hierarchical Approach." *Methods in Ecology and Evolution* 4 (1): 82–94. <https://doi.org/10.1111/j.2041-210x.2012.00264.x>.
- Kellner, Ken, and Mike Meredith. 2021. "jagsUI: A Wrapper Around 'Rjags' to Streamline 'JAGS' Analyses." <https://CRAN.R-project.org/package=jagsUI>.
- Kermorvant, Claire, Nathalie Caill-Milly, Noëlle Bru, and Frank D'Amico. 2019. "Optimizing Cost-Efficiency of Long Term Monitoring Programs by Using Spatially Balanced Sampling Designs: The Case of Manila Clams in Arcachon Bay." *Ecological Informatics* 49 (January): 32–39. <https://doi.org/10.1016/j.ecoinf.2018.11.005>.
- Kéry, Marc, and J Andrew Royle. 2020. *Applied Hierarchical Modeling in Ecology: Analysis of Distribution, Abundance and Species Richness in r and BUGS: Volume 2: Dynamic and Advanced Models*. Academic Press.
- Kombiadou, Katerina, Florian Ganthy, Romaric Verney, Martin Plus, and Aldo Sottolichio. 2014. "Modelling the Effects of *Zostera Noltei* Meadows on Sediment Dynamics: Application to the Arcachon Lagoon." *Ocean Dynamics* 64 (10): 1499–1516.
- Latimer, Andrew M., Shanshan Wu, Alan E. Gelfand, and John A. Silander. 2006. "Building Statistical Models To Analyze Species Distributions." *Ecological Applications* 16 (1): 33–50. <https://doi.org/10.1890/04-0609>.
- Lazure, Pascal, and Franck Dumas. 2008. "An External–Internal Mode Coupling for a 3d Hydrodynamical Model for Applications at Regional Scale (MARS)."

Advances in Water Resources 31 (2): 233–50.

- Lembrechts, Jonas J., Ivan Nijs, and Jonathan Lenoir. 2019. “Incorporating Microclimate into Species Distribution Models.” *Ecography* 42 (7): 1267–79. <https://doi.org/10.1111/ecog.03947>.
- Lenoir, Jonathan, Tarek Hattab, and Guillaume Pierre. 2017. “Climatic Microrefugia Under Anthropogenic Climate Change: Implications for Species Redistribution.” *Ecography* 40 (2): 253–66. <https://doi.org/https://doi.org/10.1111/ecog.02788>.
- Liu, Gang, William J Skirving, Erick F Geiger, Jacqueline L De La Cour, Ben L Marsh, Scott F Heron, Kyle V Tirak, Alan E Strong, and C Mark Eakin. 2017. “NOAA Coral Reef Watch’s 5km Satellite Coral Bleaching Heat Stress Monitoring Product Suite Version 3 and Four-Month Outlook Version 4.” *Reef Encounter* 32 (1): 39–45.
- Marcer, Arnald, Joan Pino, Xavier Pons, and Lluís Brotons. 2012. “Modelling Invasive Alien Species Distributions from Digital Biodiversity Atlases. Model Upscaling as a Means of Reconciling Data at Different Scales.” *Diversity and Distributions* 18 (12): 1177–89. <https://doi.org/10.1111/j.1472-4642.2012.00911.x>.
- Martínez-Minaya, Joaquín, Michela Cameletti, David Conesa, and Maria Grazia Pennino. 2018. “Species Distribution Modeling: A Statistical Review with Focus in Spatio-Temporal Issues.” *Stochastic Environmental Research and Risk Assessment* 32 (11): 3227–44. <https://doi.org/10.1007/s00477-018-1548-7>.
- McInerney, Greg J., and Drew W. Purves. 2011. “Fine-Scale Environmental Variation in Species Distribution Modelling: Regression Dilution, Latent Variables and Neighbourly Advice.” *Methods in Ecology and Evolution* 2 (3): 248–57. <https://doi.org/10.1111/j.2041-210X.2010.00077.x>.
- McPherson, Jana M., Walter Jetz, and David J. Rogers. 2006. “Using Coarse-Grained Occurrence Data to Predict Species Distributions at Finer Spatial Resolutions—Possibilities and Limitations.” *Ecological Modelling* 192 (3): 499–522. <https://doi.org/10.1016/j.ecolmodel.2005.08.007>.
- Meineri, Eric, and Kristoffer Hylander. 2017. “Fine-Grain, Large-Domain Climate Models Based on Climate Station and Comprehensive Topographic Information

- Improve Microrefugia Detection.” *Ecography* 40 (8): 1003–13. <https://doi.org/10.1111/ecog.02494>.
- Muff, Stefanie, Andrea Riebler, Leonhard Held, Håvard Rue, and Philippe Saner. 2015. “Bayesian Analysis of Measurement Error Models Using Integrated Nested Laplace Approximations.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 64 (2): 231–52. <https://doi.org/10.1111/rssc.12069>.
- Naimi, Babak, Nicholas A. S. Hamm, Thomas A. Groen, Andrew K. Skidmore, and Albertus G. Toxopeus. 2014. “Where Is Positional Uncertainty a Problem for Species Distribution Modelling?” *Ecography* 37 (2): 191–203. <https://doi.org/10.1111/j.1600-0587.2013.00205.x>.
- Norberg, Anna, Nerea Abrego, F. Guillaume Blanchet, Frederick R. Adler, Barbara J. Anderson, Jani Anttila, Miguel B. Araújo, et al. 2019. “A Comprehensive Evaluation of Predictive Performance of 33 Species Distribution Models at Species and Community Levels.” *Ecological Monographs* 89 (3): e01370. <https://doi.org/10.1002/ecm.1370>.
- Pearson, Richard G., and Terence P. Dawson. 2003. “Predicting the Impacts of Climate Change on the Distribution of Species: Are Bioclimate Envelope Models Useful?” *Global Ecology and Biogeography* 12 (5): 361–71. <https://doi.org/10.1046/j.1466-822X.2003.00042.x>.
- Pettorelli, Nathalie, William F. Laurance, Timothy G. O’Brien, Martin Wegmann, Harini Nagendra, and Woody Turner. 2014. “Satellite Remote Sensing for Applied Ecologists: Opportunities and Challenges.” *Journal of Applied Ecology* 51 (4): 839–48. <https://doi.org/10.1111/1365-2664.12261>.
- Plummer, Martyn. 2003. “JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling.” *Working Papers*, 8.
- Plus, Martin, Franck Dumas, Jean-Yves Stanisiere, and Daniele Maurer. 2009. “Hydrodynamic Characterization of the Arcachon Bay, Using Model-Derived Descriptors.” *Continental Shelf Research* 29 (8): 1008–13. <https://doi.org/10.1016/j.csr.2008.12.016>.
- Potter, Kristen A., H. Arthur Woods, and Sylvain Pincebourde. 2013. “Microclimatic Challenges in Global Change Biology.” *Global Change Biology* 19 (10): 2932–39. <https://doi.org/10.1111/gcb.12257>.

- Rebaudo, François, Emile Faye, and Olivier Dangles. 2016. “Microclimate Data Improve Predictions of Insect Abundance Models Based on Calibrated Spatiotemporal Temperatures.” *Frontiers in Physiology* 7. <https://www.frontiersin.org/article/10.3389/fphys.2016.00139>.
- Roberts, David R., Volker Bahn, Simone Ciuti, Mark S. Boyce, Jane Elith, Gutzeta Guillera-Arroita, Severin Hauenstein, et al. 2017. “Cross-Validation Strategies for Data with Temporal, Spatial, Hierarchical, or Phylogenetic Structure.” *Ecography* 40 (8): 913–29. <https://doi.org/10.1111/ecog.02881>.
- Rue, Håvard, Sara Martino, and Nicholas Chopin. 2009. “Approximate Bayesian Inference for Latent Gaussian Models Using Integrated Nested Laplace Approximations (with Discussion).” *Journal of the Royal Statistical Society B* 71: 319–92.
- Seo, Changwan, James H Thorne, Lee Hannah, and Wilfried Thuiller. 2009. “Scale Effects in Species Distribution Models: Implications for Conservation Planning Under Climate Change.” *Biology Letters* 5 (1): 39–43. <https://doi.org/10.1098/rsbl.2008.0476>.
- Stark, Jordan R., and Jason D. Fridley. 2022. “Microclimate-Based Species Distribution Models in Complex Forested Terrain Indicate Widespread Cryptic Refugia Under Climate Change.” *Global Ecology and Biogeography* 31 (3): 562–75. <https://doi.org/10.1111/geb.13447>.
- Stoklosa, Jakub, Christopher Daly, Scott D. Foster, Michael B. Ashcroft, and David I. Warton. 2015. “A Climate of Uncertainty: Accounting for Error in Climate Variables for Species Distribution Models.” *Methods in Ecology and Evolution* 6 (4): 412–23. <https://doi.org/10.1111/2041-210X.12217>.
- Team, R Core. 2018. “R: A Language and Environment for Statistical Computing.” <https://www.R-project.org/>.
- Tezuka, Naoaki, Masaei Kanematsu, Kimio Asami, Kazutaka Sakiyama, Masami Hamaguchi, and Hironori Usuki. 2013. “Effect of Salinity and Substrate Grain Size on Larval Settlement of the Asari Clam (Manila Clam, *Ruditapes philippinarum*).” *Journal of Experimental Marine Biology and Ecology* 439: 108–12. <https://doi.org/https://doi.org/10.1016/j.jembe.2012.10.020>.
- Yin, Xuwang, Peng Chen, Hai Chen, Wen Jin, and Xiwu Yan. 2017. “Physiological

- Performance of the Intertidal Manila Clam (*Ruditapes Philippinarum*) to Long-Term Daily Rhythms of Air Exposure.” *Scientific Reports* 7 (1): 1–12.
- Zellweger, Florian, Pieter De Frenne, Jonathan Lenoir, Pieter Vangansbeke, Kris Verheyen, Markus Bernhardt-Römermann, Lander Baeten, et al. 2020. “Forest Microclimate Dynamics Drive Plant Responses to Warming.” *Science* 368 (6492): 772–75. <https://doi.org/10.1126/science.aba6880>.
- Zipkin, Elise F., J. Andrew Royle, Deanna K. Dawson, and Scott Bates. 2010. “Multi-Species Occurrence Models to Evaluate the Effects of Conservation and Management Actions.” *Biological Conservation* 143 (2): 479–84. <https://doi.org/10.1016/j.biocon.2009.11.016>.
- Zurell, Damaris, Christian König, Anne-Kathleen Malchow, Simon Kapitza, Greta Bocedi, Justin Travis, and Guillermo Fandos. 2021. “Spatially Explicit Models for Decision-Making in Animal Conservation and Restoration.” *Ecography* 44: 1–16. <https://doi.org/10.1111/ecog.05787>.
- Zuur, Alain F., Elena N Ieno, and Anatoly A Saveliev. 2017. *Beginner’s Guide to Spatial, Temporal and Spatial-Temporal Ecological Data Analysis with R-INLA – Volume1: Using GLM and GLMM*. Highland Statistics Ltd. Vol. 1.

Supplementary materials

Appendix 1: Extended results with all sampling sizes

Sampling size = 100

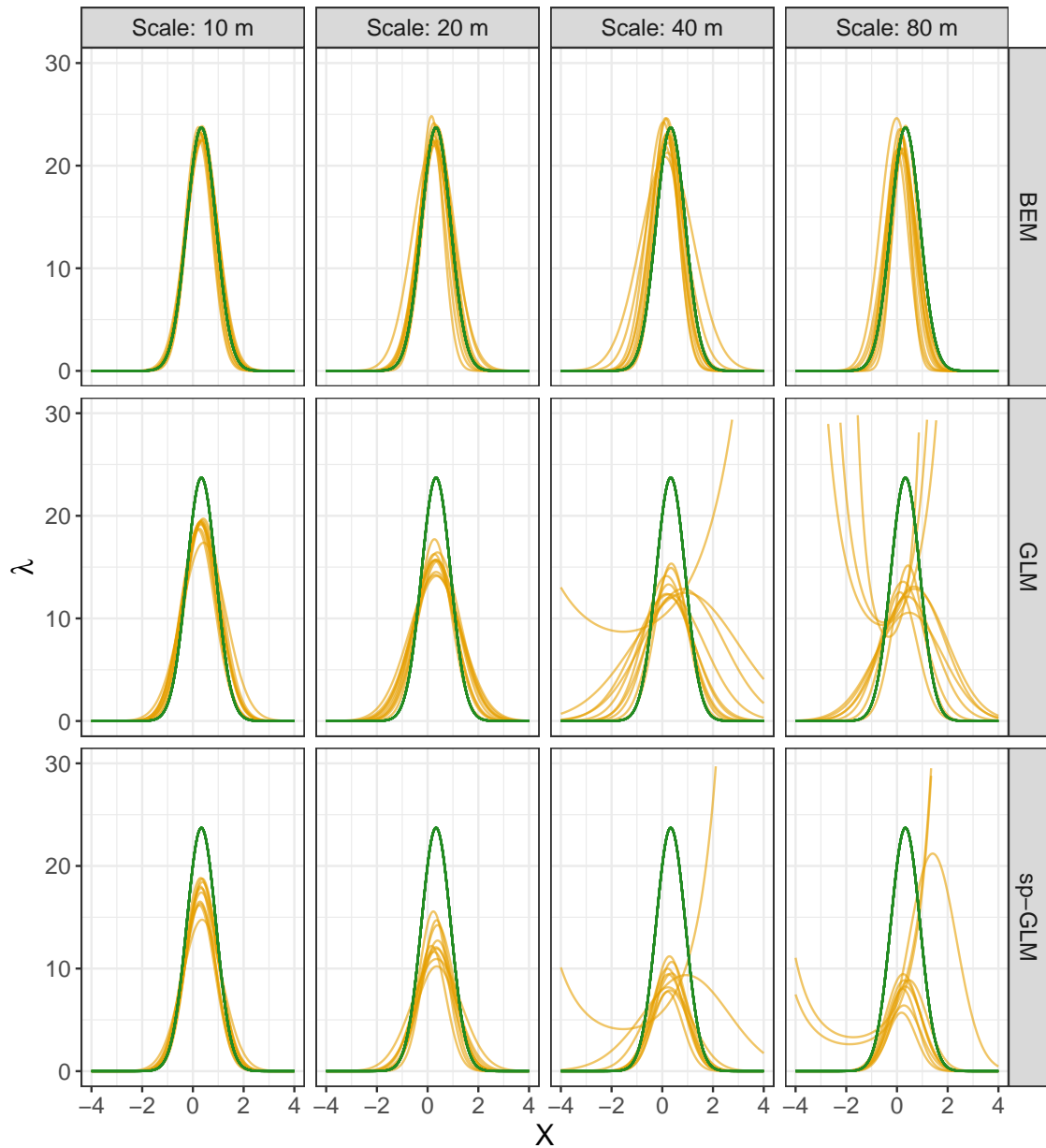


Figure 4: Species response curves estimated by the three models (orange curves) compared to simulated species response curves (green curves) in four scenarios according to different spatial scales at which environmental covariates were used. Each orange curve corresponds to one of the ten replications ran in the simulation analysis. (Sampling size = 100 sites)

Sampling size = 300

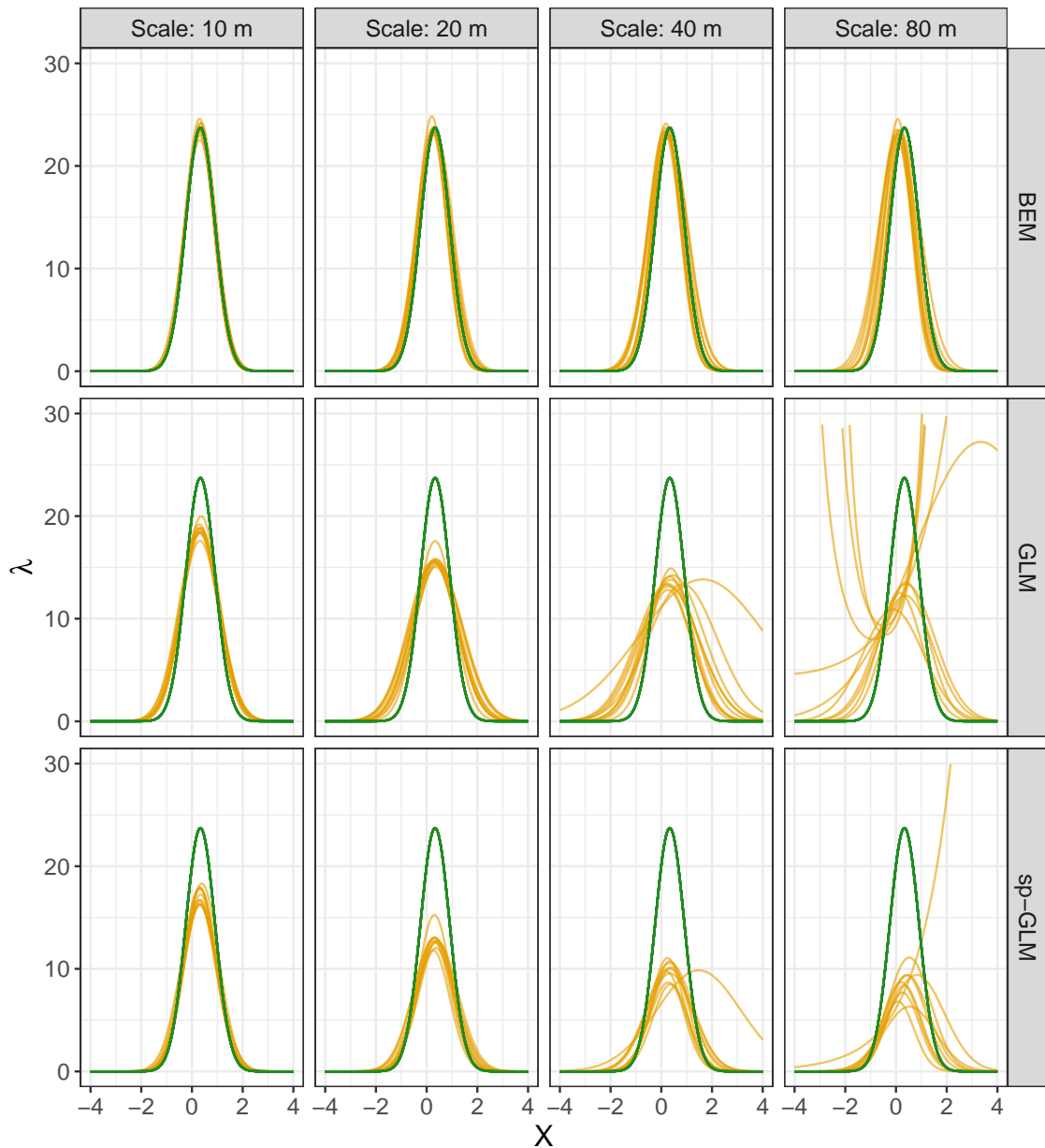


Figure 5: Species response curves estimated by the three models (orange curves) compared to simulated species response curves (green curves) in four scenarios according to different spatial scales at which environmental covariates were used. Each orange curve corresponds to one of the ten replications ran in the simulation analysis. (Sampling size = 300 sites)

Sampling size = 500

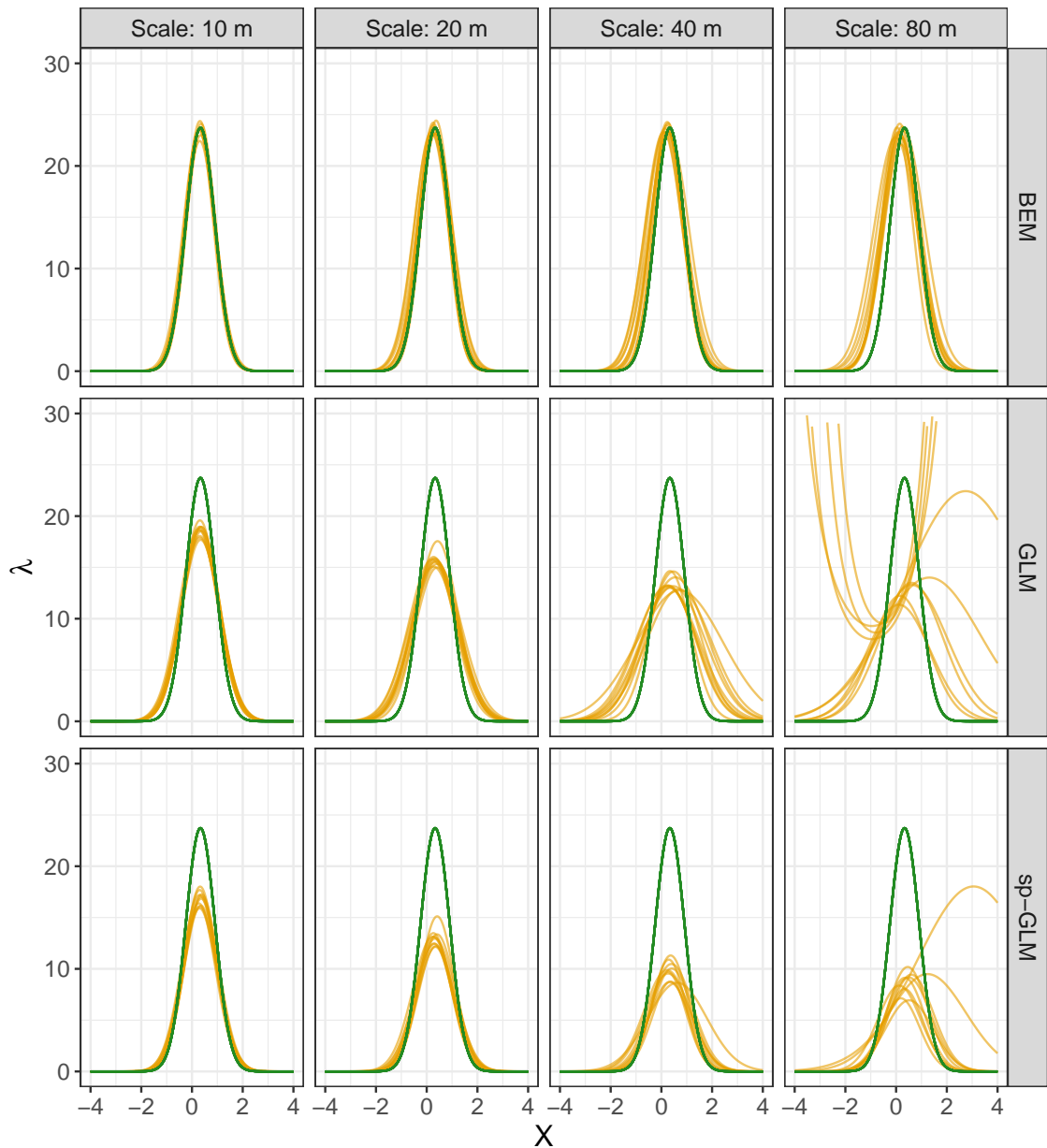


Figure 6: Species response curves estimated by the three models (orange curves) compared to simulated species response curves (green curves) in four scenarios according to different spatial scales at which environmental covariates were used. Each orange curve corresponds to one of the ten replications ran in the simulation analysis. (Sampling size = 500 sites)

Table 3: RMSE between predicted and simulated counts for the three models on the three types of data used for prediction depending on the scale of observed environment.

Sampling size	Data type	Model	Scale size			
			10 m	20 m	40 m	80 m
100	Train	BEM	1.79 (0.13)	1.85 (0.19)	1.86 (0.2)	1.86 (0.2)
		GLM	4.32 (0.31)	6.09 (0.4)	6.96 (0.54)	7.34 (0.54)
		sp-GLM	2.21 (0.19)	2.33 (0.23)	2.36 (0.25)	2.39 (0.25)
	Test (interpolation)	BEM	4.37 (0.23)	5.97 (0.61)	7.86 (0.88)	8.92 (0.64)
		GLM	4.52 (0.21)	6.03 (0.45)	7.23 (0.46)	7.68 (0.34)
		sp-GLM	4.51 (0.32)	6.4 (1.01)	8.3 (2.16)	9.27 (1.88)
	Test (extrapolation)	BEM	3.38 (0.47)	4.67 (0.83)	5.77 (0.82)	6.24 (0.9)
		GLM	3.44 (0.4)	4.68 (0.55)	6.85 (2.4)	40.05 (87.95)
		sp-GLM	3.53 (0.41)	5.19 (0.79)	17.03 (30.57)	47556147.47 (127843173.2)
	Train	BEM	1.79 (0.1)	1.9 (0.14)	1.92 (0.14)	1.93 (0.14)
		GLM	4.49 (0.22)	6.09 (0.29)	7.08 (0.28)	7.51 (0.26)
		sp-GLM	2.3 (0.09)	2.43 (0.15)	2.47 (0.16)	2.48 (0.16)
BEM		4.4 (0.43)	6.52 (0.55)	8.37 (0.51)	9.63 (0.69)	

300	Test (interpolation)	GLM	4.63 (0.27)	6.26 (0.32)	7.26 (0.21)	7.54 (0.2)
		sp-GLM	4.49 (0.43)	6.58 (0.48)	8.08 (0.67)	8.33 (0.68)
	Test (extrapolation)	BEM	3.18 (0.5)	4.75 (0.45)	5.86 (0.66)	6.21 (0.83)
		GLM	3.43 (0.51)	5.01 (0.51)	6.73 (1.22)	28.4 (34.36)
		sp-GLM	3.46 (0.48)	5.14 (0.65)	8.09 (4.83)	28.88 (50.58)
		BEM	1.81 (0.07)	1.93 (0.12)	1.95 (0.12)	1.95 (0.12)
Train	GLM	4.54 (0.19)	6.11 (0.21)	7.15 (0.19)	7.55 (0.17)	
	sp-GLM	2.3 (0.05)	2.45 (0.1)	2.49 (0.1)	2.5 (0.1)	
	BEM	4.34 (0.4)	6.68 (0.57)	8.34 (0.79)	10.01 (0.92)	
500	Test (interpolation)	GLM	4.43 (0.22)	6.12 (0.36)	7.06 (0.31)	7.55 (0.38)
		sp-GLM	4.5 (0.3)	6.54 (0.59)	7.68 (0.68)	8.13 (0.72)
		BEM	3.03 (0.28)	4.64 (0.42)	5.56 (0.48)	6.31 (0.45)
	Test (extrapolation)	GLM	3.18 (0.29)	4.84 (0.36)	6.48 (0.73)	19.59 (16.49)
		sp-GLM	3.14 (0.28)	4.87 (0.38)	6.49 (1.31)	14.06 (10.83)
		BEM	3.14 (0.28)	4.87 (0.38)	6.49 (1.31)	14.06 (10.83)

Table 4: Spearman correlation coefficients between predicted and simulated counts for the three models on the three types of data used for prediction depending on the scale of observed environment.

Sampling size	Data type	Model	Scale size				
			10 m	20 m	40 m	80 m	
100	Train	BEM	0.95 (0.01)	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	
		GLM	0.77 (0.05)	0.54 (0.07)	0.32 (0.1)	0.24 (0.09)	
		sp-GLM	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	
	Test (interpolation)	BEM	0.74 (0.02)	0.55 (0.06)	0.29 (0.14)	0.14 (0.14)	
		GLM	0.75 (0.03)	0.56 (0.07)	0.31 (0.14)	0.15 (0.1)	
		sp-GLM	0.74 (0.04)	0.54 (0.08)	0.27 (0.14)	0.12 (0.11)	
	Test (extrapolation)	BEM	0.86 (0.02)	0.72 (0.06)	0.54 (0.05)	0.35 (0.1)	
		GLM	0.86 (0.01)	0.72 (0.06)	0.42 (0.34)	0.14 (0.35)	
		sp-GLM	0.86 (0.02)	0.71 (0.06)	0.31 (0.39)	-0.31 (0.12)	
	1000	Train	BEM	0.95 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)
			GLM	0.76 (0.03)	0.55 (0.05)	0.33 (0.07)	0.18 (0.06)
			sp-GLM	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)
		BEM	0.74 (0.05)	0.51 (0.08)	0.26 (0.09)	0.11 (0.11)	

300	Test (interpolation)	GLM	0.74 (0.05)	0.51 (0.08)	0.25 (0.08)	0.16 (0.08)
		sp-GLM	0.73 (0.05)	0.5 (0.07)	0.24 (0.08)	0.16 (0.1)
	Test (extrapolation)	BEM	0.86 (0.03)	0.71 (0.03)	0.5 (0.04)	0.31 (0.08)
		GLM	0.86 (0.03)	0.71 (0.03)	0.43 (0.17)	0.02 (0.33)
		sp-GLM	0.86 (0.02)	0.7 (0.05)	0.38 (0.31)	0.09 (0.28)
		BEM	0.95 (0)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)
Train	GLM	0.75 (0.02)	0.55 (0.04)	0.31 (0.06)	0.18 (0.05)	
	sp-GLM	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	0.94 (0.01)	
	BEM	0.73 (0.07)	0.49 (0.07)	0.28 (0.09)	0.09 (0.11)	
500	Test (interpolation)	GLM	0.73 (0.07)	0.5 (0.05)	0.27 (0.1)	0.12 (0.12)
		sp-GLM	0.71 (0.07)	0.48 (0.07)	0.27 (0.07)	0.14 (0.06)
		BEM	0.87 (0.03)	0.7 (0.06)	0.5 (0.09)	0.32 (0.12)
	Test (extrapolation)	GLM	0.87 (0.03)	0.7 (0.06)	0.5 (0.09)	-0.06 (0.34)
		sp-GLM	0.86 (0.04)	0.69 (0.06)	0.48 (0.08)	0.09 (0.23)
		BEM	0.87 (0.03)	0.7 (0.06)	0.5 (0.09)	0.32 (0.12)

Appendix 2: Bias in ecological parameters

We also calculated the bias for three ecological meaningful parameters derived from model parameters: optimum $\theta = -\frac{\beta_1}{2\beta_2}$; maximum of abundance $\lambda_{max} = \exp(\beta_0 - \frac{\beta_1^2}{4\beta_2})$; ecological tolerance $\tau = \sqrt{-\frac{1}{2\beta_2}}$.

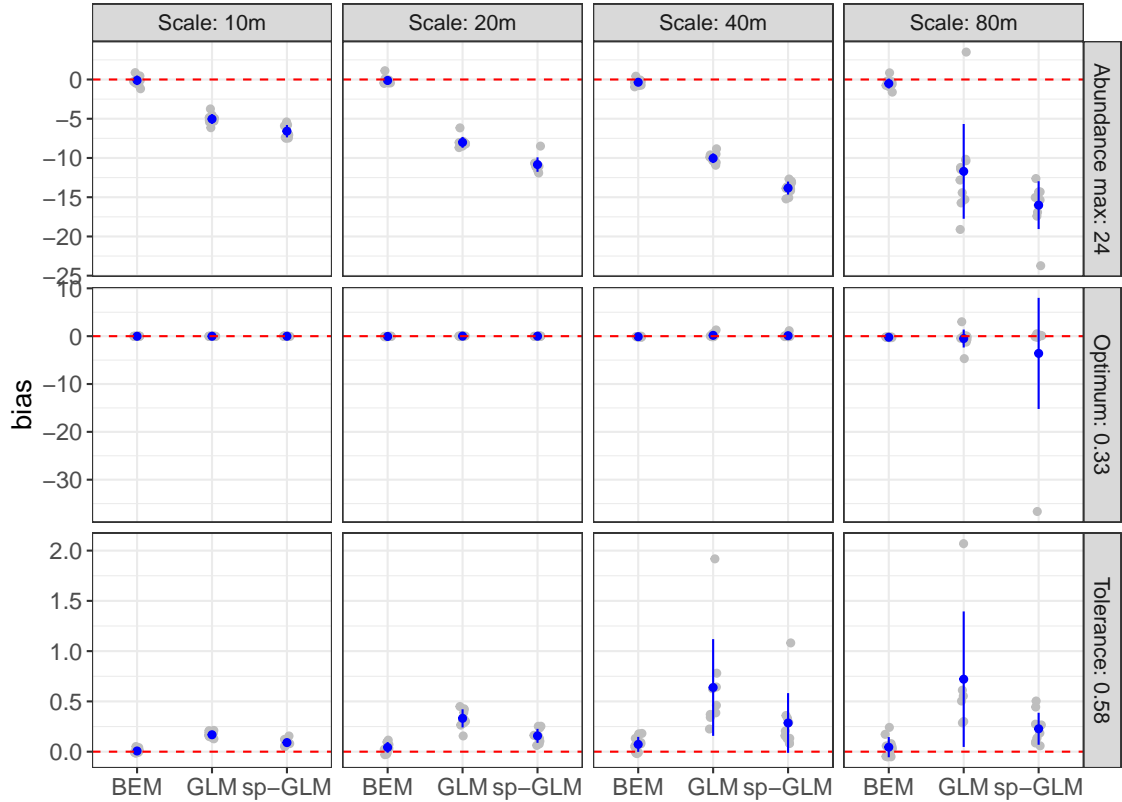


Figure 7: Bias of ecological parameters (optimum, tolerance, maximum of abundance) derived from model coefficients (the betas) estimated by the three models fitted at four different spatial scales at which environmental data was averaged prior to analysis. Grey points indicate bias at the replication level. Blue points represent average bias on ten replications. Blue errorbars represent departure from the mean by one standard-deviation.

Table 5: Average bias and RMSE of the three models for the ecological parameters derived from coefficient estimates in four scale scenarios after ten simulated replications. Numbers in brackets indicate standard-deviations. Sampling size in each replication was 300 sites.

scale	param	model	True	Estimate	Bias	RMSE
10 m	max	BEM	23.73	23.62 (0.61)	-0.11 (0.61)	0.58
		GLM	23.73	18.69 (0.62)	-5.04 (0.62)	5.07
		sp-GLM	23.73	17.14 (0.79)	-6.58 (0.79)	6.63
	opt	BEM	0.33	0.32 (0.02)	-0.02 (0.02)	0.03
		GLM	0.33	0.32 (0.03)	-0.01 (0.03)	0.03
		sp-GLM	0.33	0.32 (0.03)	-0.01 (0.03)	0.03
	tol	BEM	0.58	0.58 (0.02)	0.01 (0.02)	0.02
		GLM	0.58	0.75 (0.03)	0.17 (0.03)	0.17
		sp-GLM	0.58	0.67 (0.03)	0.09 (0.03)	0.10
20 m	max	BEM	23.73	23.6 (0.48)	-0.12 (0.48)	0.47
		GLM	23.73	15.72 (0.69)	-8.01 (0.69)	8.04
		sp-GLM	23.73	12.89 (0.92)	-10.84 (0.92)	10.88
	opt	BEM	0.33	0.27 (0.05)	-0.06 (0.05)	0.08
		GLM	0.33	0.34 (0.05)	0.01 (0.05)	0.05
		sp-GLM	0.33	0.33 (0.05)	0 (0.05)	0.05
	tol	BEM	0.58	0.62 (0.05)	0.04 (0.05)	0.07
		GLM	0.58	0.91 (0.09)	0.33 (0.09)	0.34
		sp-GLM	0.58	0.73 (0.07)	0.16 (0.07)	0.17
40 m	max	BEM	23.73	23.38 (0.41)	-0.35 (0.41)	0.52
		GLM	23.73	13.69 (0.61)	-10.04 (0.61)	10.05
		sp-GLM	23.73	9.89 (0.83)	-13.84 (0.83)	13.86
	opt	BEM	0.33	0.22 (0.07)	-0.11 (0.07)	0.13
		GLM	0.33	0.5 (0.44)	0.17 (0.44)	0.45
		sp-GLM	0.33	0.43 (0.37)	0.1 (0.37)	0.37
	tol	BEM	0.58	0.65 (0.07)	0.07 (0.07)	0.10
		GLM	0.58	1.22 (0.48)	0.64 (0.48)	0.78
		sp-GLM	0.58	0.86 (0.3)	0.29 (0.3)	0.40
80 m	max	BEM	23.73	23.22 (0.63)	-0.51 (0.63)	0.79
		GLM	23.73	12.02 (6.04)	-11.71 (6.04)	13.03
		sp-GLM	23.73	7.72 (3.05)	-16.01 (3.05)	16.27
	opt	BEM	0.33	0.11 (0.06)	-0.22 (0.06)	0.23
		GLM	0.33	-0.16 (1.88)	-0.49 (1.88)	1.85
		sp-GLM	0.33	-3.27 (11.61)	-3.6 (11.61)	11.59
	tol	BEM	0.58	0.62 (0.1)	0.04 (0.1)	0.11
		GLM	0.58	1.3 (0.67)	0.72 (0.67)	0.95
		sp-GLM	0.58	0.81 (0.16)	0.23 (0.16)	0.27

Appendix 3: Visualize misalignment effect

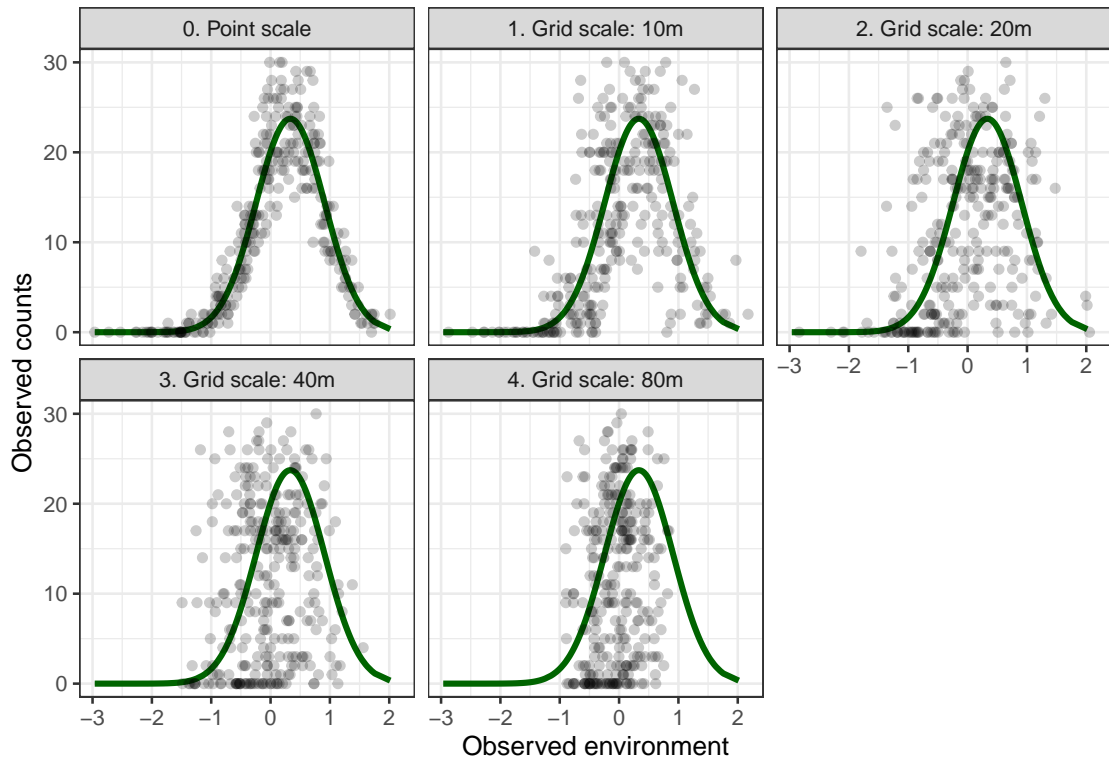


Figure 8: Observed relationships between point-scale counts and environmental variable collected at different spatial scales. Green curves represent the simulated fine-scale relationship.

Appendix 4: Bias in coefficients

GLM and sp-GLM had in average all their coefficients biased towards zero with increasing following spatial scales. BEM had almost no bias for β_0 , to note however a slight increase, but had minor bias in β_2 and major bias in β_1 .

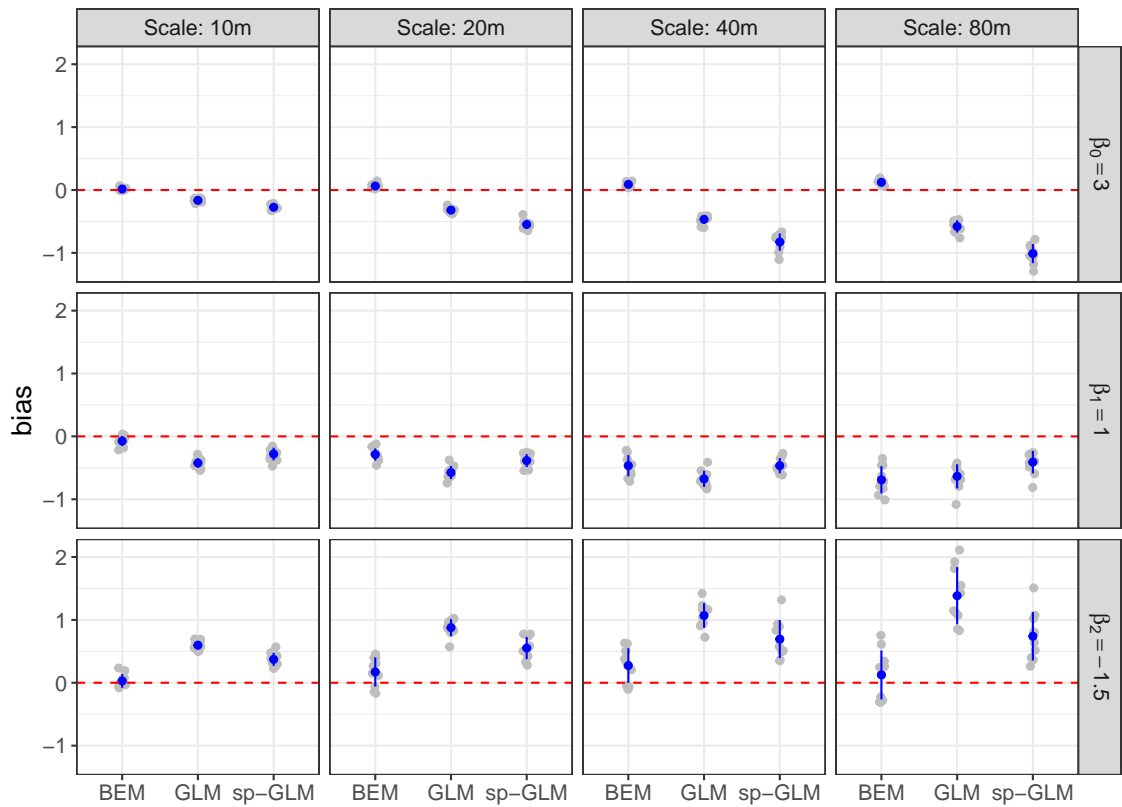


Figure 9: Bias of coefficients estimated by the three models in four scenarios of different spatial scales at which environmental data was averaged prior to analysis. Grey points indicate bias at the replication level. Blue points represent average bias on ten replications. Blue errorbars represent departure from the mean by one standard-deviation.

Table 6: Average bias and RMSE of the three models for the coefficient estimates in four scale scenarios after ten simulated replications. Numbers in brackets indicate standard-deviations. Sampling size in each replication was 300 sites.

scale	param	model	True	Estimate	Bias	RMSE
10 m	β_0	BEM	3.0	3.02 (0.02)	0.02 (0.02)	0.03
		GLM	3.0	2.83 (0.03)	-0.17 (0.03)	0.17
		sp-GLM	3.0	2.73 (0.04)	-0.27 (0.04)	0.28
	β_1	BEM	1.0	0.93 (0.08)	-0.07 (0.08)	0.11
		GLM	1.0	0.58 (0.08)	-0.42 (0.08)	0.43
		sp-GLM	1.0	0.72 (0.1)	-0.28 (0.1)	0.30
	β_2	BEM	-1.5	-1.47 (0.11)	0.03 (0.11)	0.11
		GLM	-1.5	-0.9 (0.07)	0.6 (0.07)	0.60
		sp-GLM	-1.5	-1.13 (0.1)	0.37 (0.1)	0.38
20 m	β_0	BEM	3.0	3.06 (0.04)	0.06 (0.04)	0.07
		GLM	3.0	2.68 (0.04)	-0.32 (0.04)	0.32
		sp-GLM	3.0	2.45 (0.07)	-0.55 (0.07)	0.55
	β_1	BEM	1.0	0.71 (0.1)	-0.29 (0.1)	0.30
		GLM	1.0	0.43 (0.1)	-0.57 (0.1)	0.58
		sp-GLM	1.0	0.62 (0.1)	-0.38 (0.1)	0.40
	β_2	BEM	-1.5	-1.33 (0.23)	0.17 (0.23)	0.28
		GLM	-1.5	-0.62 (0.13)	0.88 (0.13)	0.89
		sp-GLM	-1.5	-0.95 (0.18)	0.55 (0.18)	0.58
40 m	β_0	BEM	3.0	3.09 (0.04)	0.09 (0.04)	0.10
		GLM	3.0	2.53 (0.07)	-0.47 (0.07)	0.47
		sp-GLM	3.0	2.17 (0.14)	-0.83 (0.14)	0.84
	β_1	BEM	1.0	0.53 (0.17)	-0.47 (0.17)	0.49
		GLM	1.0	0.32 (0.13)	-0.68 (0.13)	0.69
		sp-GLM	1.0	0.53 (0.12)	-0.47 (0.12)	0.48
	β_2	BEM	-1.5	-1.23 (0.28)	0.27 (0.28)	0.38
		GLM	-1.5	-0.43 (0.2)	1.07 (0.2)	1.09
		sp-GLM	-1.5	-0.8 (0.3)	0.7 (0.3)	0.75
80 m	β_0	BEM	3.0	3.12 (0.04)	0.12 (0.04)	0.13
		GLM	3.0	2.42 (0.1)	-0.58 (0.1)	0.59
		sp-GLM	3.0	1.99 (0.15)	-1.01 (0.15)	1.02
	β_1	BEM	1.0	0.31 (0.22)	-0.69 (0.22)	0.72
		GLM	1.0	0.36 (0.19)	-0.64 (0.19)	0.66
		sp-GLM	1.0	0.59 (0.18)	-0.41 (0.18)	0.44
	β_2	BEM	-1.5	-1.38 (0.39)	0.12 (0.39)	0.39
		GLM	-1.5	-0.12 (0.45)	1.38 (0.45)	1.45
		sp-GLM	-1.5	-0.76 (0.38)	0.74 (0.38)	0.82

Appendix 5: Graphs of prediction results - simulation study

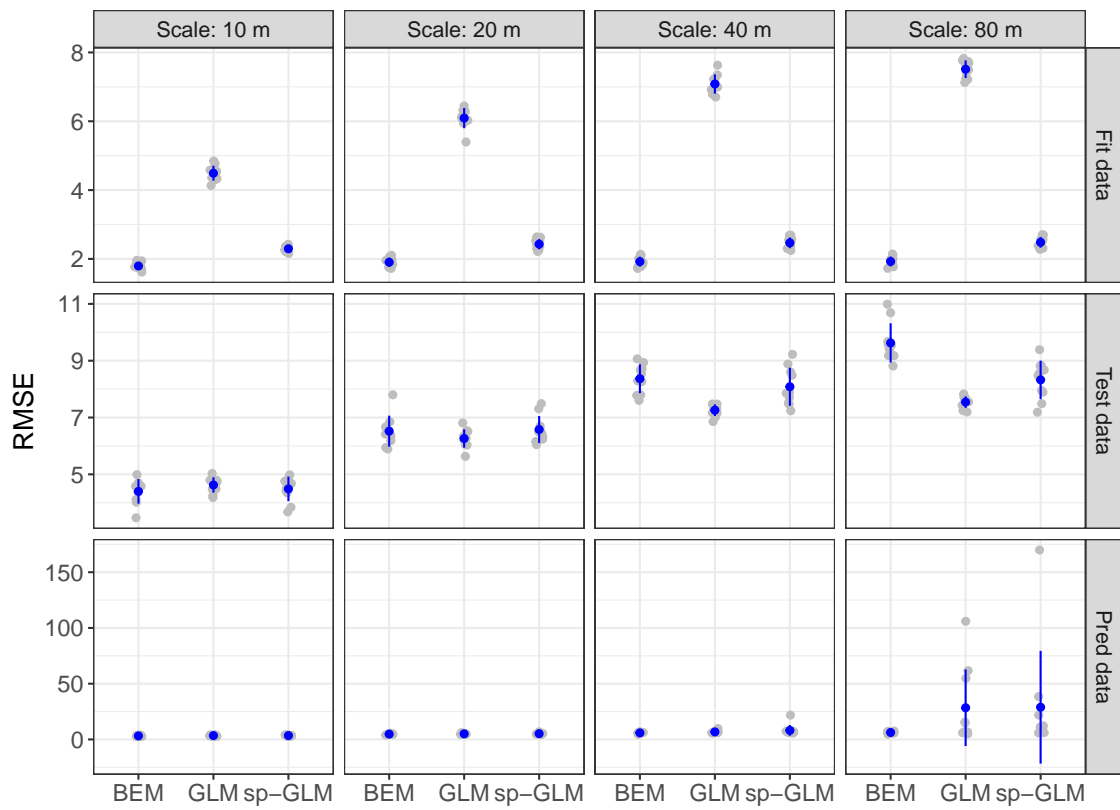


Figure 10: Root-mean-squared-errors (RMSEs) of models predicted counts compared to the simulated ones for each type of prediction and spatial scale at which environmental covariate was used for model fitting. Grey points indicate RMSE at the replication level. Blue points represent average RMSE on the ten replications. Blue errorbars represent departure from the mean by one standard-deviation.

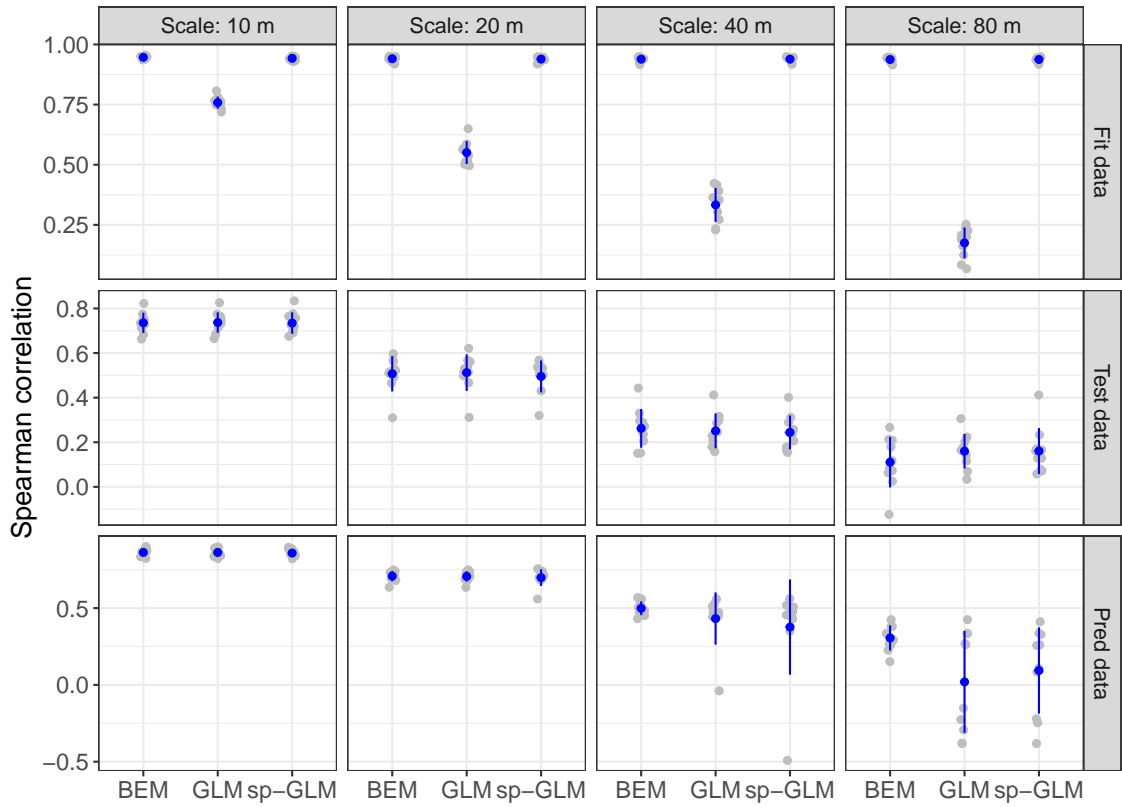


Figure 11: Spearman correlation of models predicted counts compared to the simulated ones for each type of prediction and spatial scale at which environmental covariate was used for model fitting. Grey points indicate correlation at the replication level. Blue points represent average correlation on the ten replications. Blue errorbars represent departure from the mean by one standard-deviation.

Appendix 6: Graph of predictions - Case study

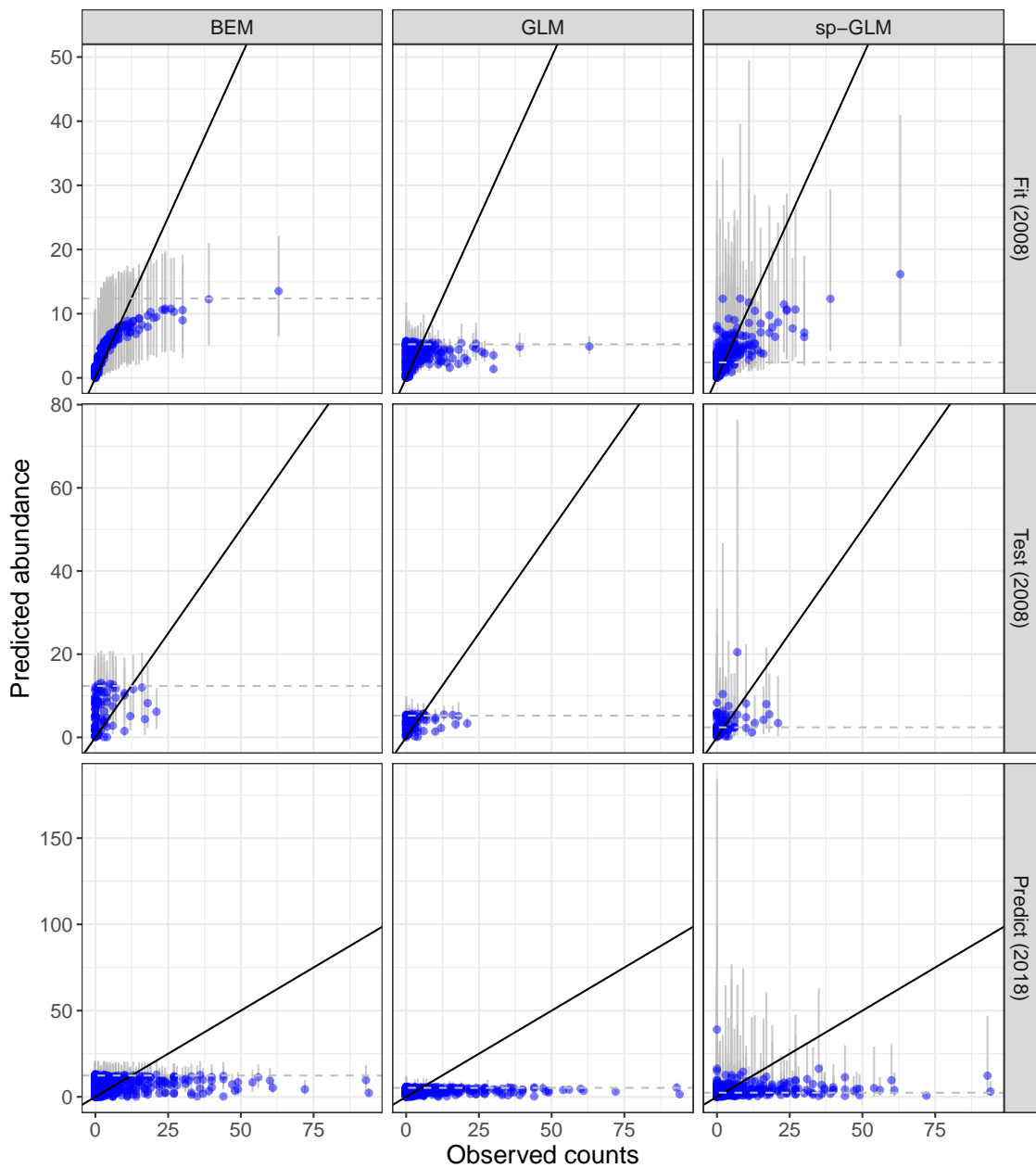


Figure 12: Relationship between predicted counts by the three models and observed counts of Manila clam at the same locations in Arcachon Bay. Rows correspond to different type of environmental data on which predictions had been done (see Materials and Methods part for further informations). The grey dashed lines represented estimated abundance at average environmental conditions.

3.3 Conclusion

In this study, we showed that partially observing the covariate, when environmental data are averaged at a larger scale than the scale of effect, leads to low accuracy in SER estimates and predictions produced by GLMs and spatial GLMs. In addition, we found that a hierarchical GLM with a Berkson error structure improved estimates of SERs. However, the Berkson error model (BEM) did not improve predictions and thus needs further investigations for broad applications. In conclusion, we advised ecologists not to interpret GLMs estimates of SERs when environmental data may describe the environment at a coarser scale than the scale of effect. In addition, we discussed the potential use of BEM to indicate misestimated variability in environmental descriptors and prevent the use of SER estimates to guide management actions when it occurs.

4 Chapter 3: A new method to explicitly estimate the shift of optimum along gradients in multispecies studies

The core of this chapter is a paper under revision in *Journal of Biogeography*: B. Mourguiart, B. Liquet, K. Mengersen, T. Couturier, J. Mansons, Y. Braud, A. Besnard. A new method to explicitly estimate the shift of optimum along gradients in multispecies studies. In revision in *Journal of Biogeography*. I also made an online oral presentation, “Explicitly estimating shifts in optimum positions of multiple species along environmental gradient”, for the *52èmes Journées de Statistique de la Société Française de Statistique* (7 June 2021).

4.1 Synopsis

Optimum shifts in species-environment relationships are intensively studied in a wide range of ecological topics, including climate change, species invasion, or theoretical ecology. In our case, we were motivated by a case study in the context of climate change. We wanted to estimate species-specific optimum shifts in Orthoptera communities between two surveys conducted 30 years apart along an elevation gradient in the French Alps.

Different modelling frameworks were developed to estimate optimum shifts between two surveys but, to our knowledge, none explicitly. Maybe the two most widely used methods are based on the mean comparison (such as Student t-test) (Chen *et al.*, 2009; Menéndez *et al.*, 2014; Freeman *et al.*, 2018) and regression (such as GLMs or GAMs) methods (Lembrechts *et al.*, 2017). The mean comparison method compares the mean values of occupied sites on the gradient between the two samples. Such optimum estimates could be biased when sampling effort is uneven along the gradient (ter Braak & Looman, 1986; Shoo *et al.*, 2006), which is common in ecology (Rumpf *et al.*, 2018; Veen *et al.*, 2021). Regression methods fit two GLMs (or GAMs), one for each survey, and use coefficient estimates to calculate optima and derive shifts (Coudun & Gégout, 2005). Such a separate analysis makes it difficult to estimate uncertainty around the shift, so this uncertainty is often omitted (e.g. Lembrechts *et al.* (2017); Urli *et al.* (2014)). This approach also raises issues for species with an optimum close to the edge of the sampled gradient. For these species, the regression can not accurately estimate the optimum (ter Braak & Looman, 1986; Coudun & Gégout, 2006), and edge species are usually withdrawn from the analysis (e.g. Lenoir *et al.*, 2008). Hence, none of the two methods appeared satisfying to us.

We formulated a new model to explicitly estimate optimum shifts of multiple species, which we called the Explicit Hierarchical Model of Optimum Shifts (EHMOS). We compared the accuracy of optimum shifts, and precision estimates of EHMOS to the mean comparison mean and a

generalised linear mixed model through simulations. Specifically, we investigated the effects of sampling design, species ecological specialisation and marginality (i.e., partially observed relationship) on the accuracy of the three methods. We also fitted the three methods to the motivating case of study.

4.2 Publication

A new method to explicitly estimate the shift of optimum along gradients in multispecies studies

Running title: Modelling of gradient optimum shifts

Bastien Mourguiart^{1*}, Benoit Lique^{1,2}, Kerrie Mengersen^{1,3}, Thibaut Couturier⁴, Jérôme Mansons⁵, Yoan Braud⁶, Aurélien Besnard⁴

¹ Laboratoire de Mathématiques et de leurs Applications, Université de Pau et des Pays de l'Adour, E2S UPPA, CNRS, Anglet, France

² Department of Mathematics and Statistics, Macquarie University, Sydney, NSW 2109, Australia

³ ARC Centre of Excellence for Mathematical and Statistical Frontiers, School of Mathematical Science, Queensland University of Technology, Brisbane, QLD 4072, Australia

⁴ CEFÉ, Univ Montpellier, CNRS, EPHE-PSL University, IRD, Univ Paul Valéry Montpellier 3, Montpellier, France

⁵ Parc National du Mercantour, Nice, France

⁶ Bureau d'études Entomia, Vaumeilh, France

Abstract

Aim

Optimum shifts in species–environment relationships are intensively studied in a wide range of ecological topics, including climate change and species invasion. Numerous statistical methods are used to study optimum shifts but, to our knowledge, none explicitly estimate it. We extended an existing model to explicitly estimate optimum shifts for multiple species. We called this new Bayesian hierarchical model the Explicit Hierarchical Model of Optimum Shifts (EHMOS).

Location

All locations

Taxon

All taxa

Methods

In a simulation study, we compared the accuracy of EHMOS to a mean comparison method and a Bayesian generalized linear mixed model (GLMM). Specifically, we tested if the accuracy of the methods was sensitive to (1) sampling design, (2) species optimum position, and (3) species ecological specialization. In addition, we compared the three methods using a real dataset of investigated optimum shifts in 24 Orthopteran species between two time periods along an elevation gradient.

Results

Of all the simulated scenarios, EHMOS was the most accurate method. GLMM was the most sensitive method to species optimum position, providing unreliable estimates in the presence of marginal species, i.e. species with an optimum close to a sampling boundary. The mean comparison method was also sensitive to species optimum position and ecological specialization, especially in an unbalanced sampling design, with high negative bias and low interval coverage compared to EHMOS. The case study results obtained with EHMOS were consistent with what

is expected considering ongoing climate change, with mostly upward shifts, which further improved confidence in the accuracy of the EHMOS method.

Main conclusions

Our findings indicate that EHMOS could be used for a wide range of topics and extended to produce new insights, especially in climate change studies. Explicit estimation of optimum shifts notably allows investigation of ecological assumptions that could explain interspecific variability of these shifts.

Keywords: Bayesian hierarchical modelling, Generalized Linear Mixed Model (GLMM), mean comparison method, optimum shift, sampling design, species distribution, species ecological specialization, species response curve

Introduction

Understanding changes in relationships between species and environmental, geographical or temporal gradients is of central interest in biogeography. For instance, assessing changes in species distribution along latitudinal (Thorson et al., 2016), elevational (Maggini et al., 2011) or climatic (Tayleur et al., 2015) gradients in response to ongoing climate change is crucial for developing relevant biodiversity change scenarios (Taheri et al., 2021) and predicting its impact on ecosystem health and human well-being (Pecl et al., 2017). Modelling these changes is also central in the study of invasive species to predict the area of potential colonization (Elith et al., 2010). The shifts in species phenology, i.e. changes along a temporal gradient, in response to climate change are also well-studied to obtain a mechanistic understanding of climate change impacts on species distributions (Moussus et al., 2010; Strebel et al., 2014; Vitasse et al., 2021).

All this research aims at modelling the relationship between species characteristics – such as occurrence probability (Maggini et al., 2011), abundance (Thorson et al., 2016) or singing activity (Strebel et al., 2014) – and linear covariates describing the studied gradient. These modelled relationships, often named “species response curves” (Tayleur et al., 2015), are usually unimodal curves (Oksanen & Minchin, 2002). They are characterized by a width, a maximum, an optimum (i.e. the value

on the gradient where the maximum abundance or occurrence probability of the species is reached), and a lower and upper limit. Changes in the relationship could then occur in each of these parameters, providing different and complementary information on the way species are affected by change over the gradient. However, these changes may be difficult to estimate. Changes in width or range limits of the species response curve may for instance be more difficult to estimate than optimum shifts due to sampling design or detectability issues linked to lower abundance near range limits (Shoo et al., 2006). Hence, studying optimum shifts is often favoured over shifts at range limits to disentangle ecological changes from potential sampling artifacts (Shoo et al., 2006; Lenoir et al., 2008; Moussus et al., 2010).

Different modelling frameworks have been developed to study optimum shifts. The simplest, and probably the most widely used, is the mean comparison method (e.g. Shoo et al., 2006; Freeman et al., 2018), also known as the weighted average method (ter Braak & Looman, 1986; Thorson et al., 2016), in which mean values of occupied sites on the gradient are compared between two or more samples. However, such optimum estimates could be biased when sampling effort is uneven along the gradient (ter Braak & Looman, 1986; Shoo et al., 2006). Such uneven sampling is common in ecology (Shoo et al., 2006; Rumpf et al., 2018), especially in cases of sampling difficult-to-access terrains such as mountains. Designing a random sample of a species distribution requires knowing the species distribution prior to sampling it, which is almost never the case. Regression methods, such as generalized linear models (GLMs), are also commonly used to estimate optimum shifts (Coudun & Gégout, 2005; Lenoir et al., 2008; Urli et al., 2014) and are known to be more robust to unbalanced sampling design (ter Braak & Looman, 1986). In this approach, two regressions are conducted to estimate the species optimum for each sample separately (Coudun & Gégout, 2005). Shift is then calculated as the difference between the two estimated optima (as in Lenoir et al., 2008). Such a separate analysis makes it difficult to estimate uncertainty around the shift, so this uncertainty is often omitted (Urli et al., 2014; e.g. Lembrechts et al., 2017). This absence of uncertainty is particularly prejudicial when post-hoc tests are performed on shift point estimates (Felde et al., 2012). This approach also raises issues for species with an optimum close to the edge of the sampled gradient,

hereafter referred to as edge species, also called marginal species (Hernandez et al., 2006). In these cases, the regression follows a logistic relationship and does not provide an optimum (Citores et al., 2020). As a consequence, such edge species are usually withdrawn from analysis (Felde et al., 2012; e.g. Rumpf et al., 2018); however, they may be particularly sensitive to changes as they are limited by certain constraints (physiological, environmental, etc.) and subject to extirpation phenomenon, i.e. local extinctions caused by a lack of remaining suitable areas (Freeman et al., 2018). Moreover, estimates of the mean shift at the community level could be biased by the removal of edge species, as the optimum shift could depend on the position along the gradient (Rumpf et al., 2018). A model robust to unbalanced sampling design, providing estimates of uncertainties and including all species whatever the location of their optimum would be a substantial improvement for modelling shift in optima.

We developed an extension of the Gaussian logistic model (ter Braak & Looman, 1986; Jamil et al., 2014) that explicitly estimates shifts in optima for several species simultaneously, with their uncertainties. We also fitted a Bayesian Generalized Linear Mixed Model to simultaneously analyse both sampling occasions and all species and use the posterior distributions of the estimated coefficients to compute posterior distributions of species optimum shifts. The performance of these models was evaluated and compared to the estimates provided by a simple mean comparison method. The comparison was undertaken through a simulation study involving different sampling designs and ecological scenarios. The primary goal of the simulation study was to determine if the tested methods accurately estimate the optimum shifts of multiple species along a virtual environmental gradient. Specifically, we tested if the accuracy of the methods was sensitive to (1) sampling design, i.e. how sampling sites are distributed along the gradient, (2) species marginality, i.e. the location of the species optimum on the gradient, and (3) species ecological specialization, i.e. the width of the species–environment relationship. We finally applied the three methods to 24 Orthoptera species sampled twice 30 years apart in the French Alps along an elevation gradient.

Materials and methods

Shift modelling

We used a mean comparison method based on the Student t-test (hereafter, t-test) that compares species' average positions along a gradient between two sampling occasions, e.g. two periods or two regions. In this approach, species' optima are estimated as the mean of environmental values at which the species have been observed. We performed unpaired two-sample Student t-tests between surveys for each species separately. This provided species-specific shift estimates and associated confidence intervals.

We developed and fitted two hierarchical models: a 'classic' Generalized Linear Mixed Model (hereafter, cGLMM) and a new Explicit Hierarchical Model of Optimum Shifts (EHMOS). Both models assumed that the occurrence state (present or absent) observed for species i at site j during sampling occasion k , $Y_{i,j,k}$, relies on the probability $\psi_{i,j,k}$ that species i can occupy site j during sampling occasion k . In both models, we described the observation data as an outcome of a Bernoulli trial:

$$Y_{i,j,k} \sim \text{Bernoulli}(\psi_{i,j,k}).$$

We assumed that species occupancy probability is related to the gradient through a symmetric bell-shaped curve, with an optimum where the maximum occupancy probability is reached. However, the two models differ in the formulation of the species–environment relationship. cGLMM allows five shapes of species–environment relationships (flat, monotone, 'plateau', U-shaped or bell-shaped unimodal symmetric response curves) some of which could be ecologically meaningless (e.g. U-shaped curves). EHMOS by its ecological formulation models only symmetric unimodal response curves. However, note that if EHMOS estimates optimum outside the sampling range, then the estimated SRC will look like a monotone or a 'plateau' relationship. EHMOS ecological formulation also allows direct interpretation of its model coefficients, especially about optimum and shift values. Note that both models assumed that unimodal species response curves are symmetrical, i.e. occu-

pancy probabilities decrease at the same rate on both sides of the optimum. We discuss this potential limitation further in the manuscript.

Classic Generalized Linear Mixed Model (cGLMM) – The occupancy probability was modelled on the logit scale as a regression of linear and quadratic effects of environment. To allow for species-specific occurrence relationships with the gradient, we added a random intercept and random slope effects for the species ID. Changes in species distribution along the gradient between sampling occasions were assessed including a binary covariate, S_k , taking the value of 0 for one sampling occasion and 1 for the other, and its interaction with the linear and quadratic effects of the environment. As the same sites were included in the two samples, we added a random intercept term for the site ID, $\gamma_j^{site} \sim N(0, \sigma_\gamma^2)$, to account for data dependencies. Hence, in our cGLMM the occupancy model was formulated as follows:

$$\text{logit}(\psi_{i,j,k}) = \beta_{0i} + \beta_{1i} \times X_j + \beta_{2i} \times X_j^2 + \beta_{3i} \times S_k + \beta_{4i} \times S_k \times X_j + \beta_{5i} \times S_k \times X_j^2 + \gamma_j^{site} \quad (1)$$

where β_i represents the coefficients related to the species-specific effects mentioned above and X_j is the value of the environmental variable at site j . Hence, β_{0i} , β_{1i} and β_{2i} describe the relationship between species occurrence and elevation during the first sampling for species i , while β_{3i} , β_{4i} and β_{5i} denote the changes in the species-specific occurrence–gradient relationship between the two sampling occasions. We calculated the species-specific optimum positions of each sampling occasion, denoted by $\theta_{i,k}$, from these regression coefficients (see Equation 2 in ter Braak & Looman, 1986). We then calculated species shifts by subtracting optimum values of each sampling occasion as follows:

$$\begin{aligned} \delta_i &= \theta_{i,2} - \theta_{i,1} \\ &= \frac{-(\beta_{1i} + \beta_{4i})}{2 \times (\beta_{2i} + \beta_{5i})} - \frac{-\beta_{1i}}{2 \times \beta_{2i}} \end{aligned} \quad (2)$$

where δ_i is the optimum shift of species i .

Explicit Hierarchical Model of Optimum Shifts (EHMOS) – We extended the multi-species Gaussian logistic model developed to describe species unique unimodal response curves (Jamil et al., 2014) to explicitly estimate shift in optimum positions between two sampling occasions. This model, we called Explicit Hierarchical Model of Optimum Shifts (EHMOS), describes multiple species–environment relationships simultaneously for two samples:

$$\text{logit}(\psi_{i,j,k}) = \alpha_{i,k} - \frac{[X_j - (\theta_i + (S_k \times \delta_i))]^2}{2 \times \tau_{i,k}^2} + \gamma_j^{\text{site}} \quad (3)$$

where $\alpha_{i,k}$ represents the maximum occupancy probability on the logit scale reached by species i in sampling occasion k ; θ_i is the environmental optimum of species i for the first sampling occasion (as the binary variable S_k is coded as $S_1 = 0$ and $S_2 = 1$); δ_i is the shift between the two optima for species i ; $\tau_{i,k}$ is the species' ecological tolerance, a measure of the gradient portion length on which the species could occur (see Appendix S1 in Supporting Information).

In both models (Eq. (1) and (3)), the gradient covariate was standardized to have a mean equal to zero and variance equal to one. Each species-specific parameter ($\Theta_i \in \{\beta_i, \theta_i, \delta_i\}$) was modelled as a random species effect drawn from a normal distribution: $\Theta_i \sim N(\mu_\Theta, \sigma_\Theta^2)$ described by the community mean (μ_Θ) and the variance between species (σ_Θ^2). Those community parameters are shared between the two sampling occasions for all the GLMM parameters (β_i 's), and for the optimum (θ_i) and the shift (δ_i) in the EHMOS model. For the species-specific and occasion-specific parameters $\alpha_{i,k}$ and $\tau_{i,k}$, the community means could vary between sampling occasions (e.g. $\tau_{i,k} \sim N(\mu_{\tau_k}, \sigma_\tau^2)$).

We implemented the two models in a Bayesian context. We then obtained posterior distributions for each species-specific parameter. Assuming a lack of prior knowledge of a parameter's true value, hyper-parameters were afforded default priors. We used wide normal priors (with mean 0 and variance 1000) for the means of community-level effects (the μ 's) and inverse-gamma priors (Inv-gamma(0.1, 0.1)) for the community-level variances (the σ^2 's). In the cGLMM, we computed posterior

distributions of the shifts from MCMC samples of the regression coefficients (Eq. (2)). Hence, for both models we had posterior distributions for the species-specific shifts in optimum position. We used medians as point estimates and symmetric 95% credible intervals as corresponding measures of uncertainty. Symmetric 95% credible intervals were computed based on 0.025 and 0.975 quantiles of the MCMC posterior distributions.

We coded the two hierarchical Bayesian models in BUGS language and ran them in JAGS (Plummer, 2003), using the *jagsUI* package (Kellner & Meredith, 2021) in R software (R Core Team, 2018). The code is available on GitHub. We ran the analysis with the function `autojags()`, which updates the number of iterations within the burn-in phase until all parameters have a Gelman-Rubin statistic (\hat{R}) less than 1.1, suggesting satisfactory convergence (Gelman et al., 2013), or when the total number of iterations reached a maximum set at 250,000 iterations. We ran three MCMC chains with a thinning rate of 10 and with an initial burn-in phase of 15,000 iterations that was updated by 15,000 iterations until the specified convergence level was met or the maximum number of iterations was reached.

Simulation study

We simulated the sampling of a community of 20 species during two sampling occasions at 300 sampling sites distributed along an elevational gradient ranging from 1000 m to 3000 m. We used two sampling design sub-scenarios: the distribution along the virtual elevation gradient of the sampling sites was either uniform, in sampling sub-scenario A1, or unbalanced, in sampling sub-scenario A2 (Fig. 1a). In the unbalanced design, the simulated sampling sites followed a truncated normal distribution of mean 1970 m, standard-deviation 335 m restricted between 1000 m and 3000 m. We chose the mean and standard deviation according to those observed in the motivating example (see below).

We then positioned on the virtual gradient the optima of each species for the two virtual sampling occasions. We defined two optimum sub-scenarios (coded B1 and B2). In optimum sub-scenario B1, 100% of species had both their optima in the middle of the sampling range, hereafter ‘middle species’ (green points in Fig. 1b), while in B2 30% of species (six species) were considered as ‘edge species’, i.e. had

their optima close to the boundaries of the sampled gradient (red points in Fig. 1b). We defined the middle of the sampling range as the part of the gradient that contains 80% of the sampling sites (green areas in Fig. 1ab), the upper and lower parts of the remaining sampling range were considered the sampling edges (orange areas in Fig. 1ab). We simulated half of the optima of the edge species close to the upper edge, and the other half close to the lower edge. Simultaneously with the optimum of the first sampling occasion (full points in Fig. 1b), we simulated the optimum shift, thus obtaining the second optimum (circles in Fig. 1b), to respect percentages of edge and middle species in both sampling occasions. For all optimum sub-scenarios, we computed 60% of upward shifts, 20% of downward shifts and 20% of no shift to test all kinds of response while keeping the majority of upward shifts as expected in conditions of ongoing climate change. The values of upward and downward shifts were sampled from a uniform distribution bounded between 80 m and 250 m, and -250 m and -80 m respectively.

We used two ecological specialization sub-scenarios (C1 and C2). In sub-scenario C1, we simulated only specialist species, i.e. species having a narrow ecological niche and high maximum occurrence probability (Fig. 1c). In sub-scenario C2, we simulated 50% of specialist and 50% of generalist species having a wider ecological niche and lower maximum occurrence probability than specialist species (Fig. 1c). We allowed for changes in shape parameters between sampling occasions in both sub-scenarios, but not in ecological specialization type: i.e. a generalist species remains a generalist species, but could have slightly different species response curves in the two sampling occasions. Niche width and maximum occurrence probability were sampled from uniform distributions, independently for the two sampling occasions, with the distribution parameters depending on the ecological specialization type. Specialist species had widths ranging from 400 m and 600 m, while generalists' width ranged between 1200 m and 1400 m. Maximum probability was set between 0.9 and 0.99 for specialist species, and between 0.75 and 0.85 for generalist species. These distribution parameter values were chosen to clearly separate the two types of species and match the values observed in the motivating example.

We ran all possible combinations of the sub-scenarios, resulting in eight simulated

scenarios. In each, we considered each species to have a symmetric unimodal response curve in both sampling occasions (Fig. 1d). We produced site-specific probability of occupancy for each species in each scenario following:

$$\text{logit}(\psi_{i,j,k}^s) = \alpha_{i,k}^s - \frac{[X_j^s - (\theta_i^s + (S_k \times \delta_i^s))]^2}{2 \times \tau_{i,k}^s{}^2}$$

with $\psi_{i,j,k}^s$ the occurrence probability of species i at site j during sampling occasion k simulated in scenario s , X_j^s is the value on the elevation gradient of site j under scenario s , θ_i^s and δ_i^s are respectively the optimum and shift of species i in scenario s , and $\alpha_{i,k}^s$ and $\tau_{i,k}^s$ are related to niche width and maximum probability of occupancy (see Appendix S1). Observed presence or absence of species i at site j during sampling occasion k under scenario s , $Y_{i,j,k}^s$, was then drawn from a Bernoulli distribution: $Y_{i,j,k}^s \sim \text{Bernoulli}(\psi_{i,j,k}^s)$. We replicated the Bernoulli trials 30 times for each of the 8 scenarios, resulting in 240 simulated datasets.

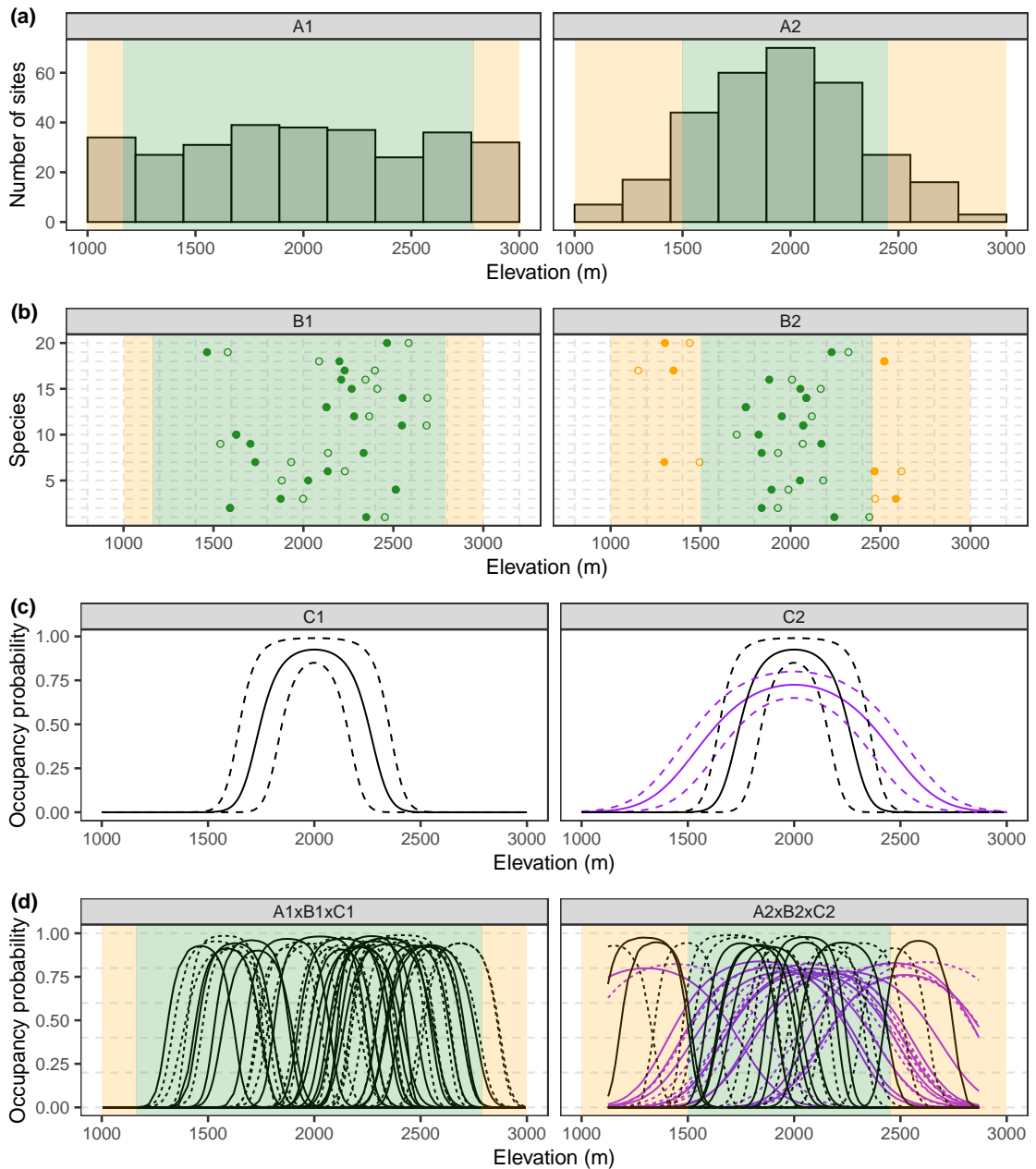


Figure 1: Example of the simulation process for two scenarios each broken down in three categories of sub-scenarios: (a) the sampling design either uniform (A1) or unbalanced (A2) (b) the optimum positions for the first (full points) and second (circles) sampling occasions of the 20 species, which is either in the middle (green area) or close to the edges (orange areas) of the sampling range in variable proportions depending on sub-scenario category (B1: 100% middle, B2: 70% middle) (c) the ecological specialization of species, which can be specialist (dark curves) or generalist species (purple curves) in variable proportions depending on the sub-scenario category (100% in C1 or 50% in C2 of specialist species). The bottom panel represents (d) the species response curves for the first and second sampling occasion (solid and dashed lines) of the two scenarios that results in the combination of the three sub-scenarios above.

Model performance assessment

We assessed estimation accuracy and precision through four metrics (i) bias, the difference between true and estimated absolute optimum shifts, (ii) root mean squared error, which combines bias and imprecision, (iii) interval coverage, i.e. the proportion of estimated credible/confidence intervals that contain the true optimum shifts, and (iv) interval score, which combines information on width of credible/confidence intervals and on interval coverage. Details on the calculation of the four metrics are provided in Appendix S2. We also compared the computing performance of the Bayesian models. We counted the number of times when at least one parameter failed to converge after the maximum of iterations was reached. Due to the number of runs and parameters, we just looked at the Gelman-Rubin statistic threshold of 1.1 (Gelman et al., 2013) and did not make visual diagnostics. We also computed the average running time for each model.

Motivating example

Our goal when starting this study was to find an appropriate modelling approach that produces estimates, and associated uncertainties, of shifts in elevation optima of an Orthoptera community studied during two periods separated by 30 years. This study was conducted in 134 sampling sites distributed along an elevation gradient ranging from 928 m to 2614 m (mean = 1869 m) above sea level in Mercantour National Park in the southern French Alps. The first survey was carried out during the summers between 1983 and 1988 (Gueguen 1990). At each site, the density of Orthoptera species was assessed based on multiple counts made using box quadrats that trapped individuals. The second survey was conducted during the summers of 2018 and 2019. Each sampling site consisted of a circle with a 70-m radius where five spatial replicates (hereafter, plots) 30 m² in area were surveyed following three successive steps: (1) one minute of listening to species stridulating in the plot by standing close to its edge, (2) six minutes of sighting species by walking across the entire plot, and (3) two 45-second sweep netting sessions across the entire plot (see Mourguiart et al., 2021 for a full description of the sampling design). To minimize the effect of varying sampling effort and detection methods between surveys, we reduced the Orthoptera records to the presence/absence data at the

site level and only kept species detected in at least 5% of sampling sites. We pooled some species' data in a species complex when species confusion was possible. Thus, we aggregated data for species in the genera *Yersinella*, *Calliptamus*, *Podisma* and in the complex *Chorthippus biggutulus-brunneus-mollis-daimeii*. We fitted the three methods to the 24 taxa kept after data processing.

Results

Before analysing the results, we checked for convergence issues and running times of cGLMM and EHMOS. In the \hat{R} statistics check, we observed that 50% of runs failed to converge (at least one parameter had \hat{R} higher than 1.1) for cGLMM (see in Appendix S3 as Table S3.1). Due to the running time, 26 hours (SD = 15 hours) on average for cGLMM (see in Appendix S3 as Table S3.2), we chose not to re-run more iterations of models that did not converge. All models converged for EHMOS and took on average 5 hours (SD = 1 hour) for this process (see in Appendix S3 as Table S3.2). In the motivating example, both models converged: EHMOS ran in 82 minutes and cGLMM in 26 minutes.

Simulation

Averaging performance metrics for all species in the eight scenarios, we found that EHMOS performed better than cGLMM and t-test whatever the performance metrics. EHMOS had the smallest average RMSE (mean = 23, SD = 1), bias (mean = 0, SD = 1), interval score (mean = 115, SD = 7), and the largest interval coverage (mean = 0.96, SD = 0). cGLMM had similar interval coverage to EHMOS (mean = 0.95, SD = 0), but relatively high positive bias (mean = 12, SD = 3), the largest RMSE of all three methods (mean = 37, SD = 4), and interval score (mean = 1402, SD = 478). The t-test had an intermediate RMSE (mean = 28, SD = 2) and interval score (mean = 207, SD = 24) between cGLMM and EHMOS, with bias as high as cGLMM, but negative (mean = -12, SD = 2), and the smallest interval coverage of the three methods (mean = 0.85, SD = 0.02).

At the scenario level, EHMOS was still the most accurate method. It performed as well as cGLMM in scenarios with only middle species, i.e. scenarios including sub-scenario B1 (Fig. 2 upper panels), and was much more accurate than cGLMM

in scenarios with edge species, i.e. scenarios including sub-scenario B2 (Fig. 2 lower panels). Averaging performance metrics for only edge species, we obtained an average RMSE almost three times larger for cGLMM than for EHMOS. cGLMM estimates of the quadratic effect of environment on species occurrence could be close to zero for edge species, yielding an unreliably high ratio of linear and quadratic coefficients, thus unreliable optimum and optimum shift estimates (Eq. 2). Hence, on average for edge species, cGLMM had larger positive bias (79 m) and interval score (8820) compared to EHMOS (mean bias = -6 m; mean interval score = 268).

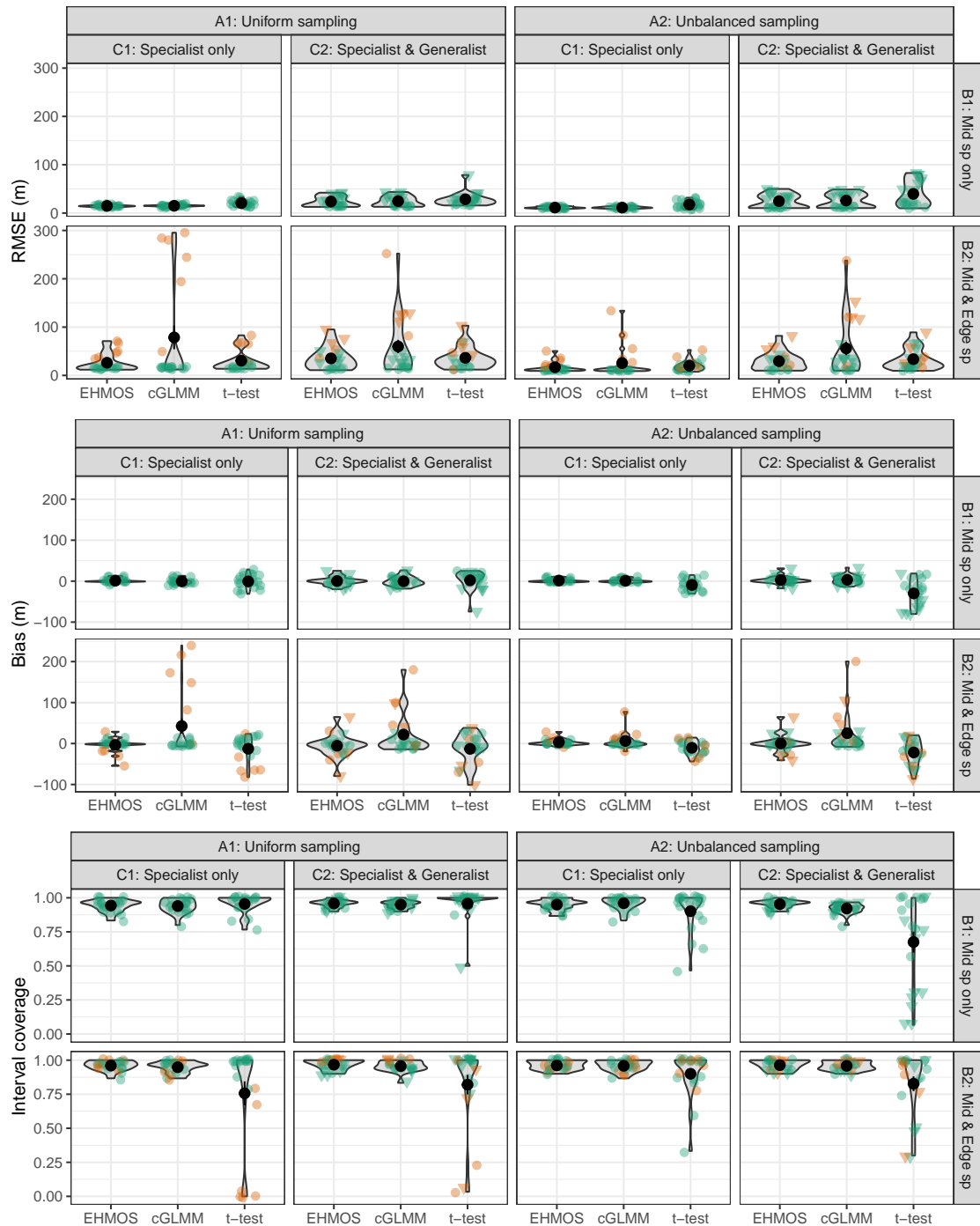


Figure 2: Distribution of species specific performance metrics for the eight simulated scenarios after 30 replications. The colour of the points represents species optimum position along the simulated gradient (in the edge, orange points, or in the middle, green points). Shape of points corresponds to species ecological specialization type relative to the width of their ecological niche (generalist species with broad ecological niche, represented by inverted triangles, or specialist species with narrower ecological niche, represented by circles).

The t-test method had comparable results to EHMOS in scenarios A1xB1xC1 and A1xB1xC2, with only slightly larger average RMSEs and interval score (Fig. 2,

Table 1). For scenarios including edge species, i.e. scenarios including sub-scenario B2, there was only a slight difference in RMSEs between EHMOS and the t-test method (Fig. 2). However, the t-test method, in scenarios with edge species, had on average negative bias that was at least two times larger than EHMOS bias (Fig. 2). The t-test method also had poorer interval metrics than EHMOS, with interval coverage never reaching the expected threshold of 95%, in contrast to EHMOS, which always reached this threshold (Fig. 2), and with interval scores 30% to 150% larger than EHMOS in scenarios with edge species (Table 1). The t-test method also resulted in poorer performance than EHMOS in scenarios with unbalanced sampling design (sub-scenario A2), with negative bias, smaller interval coverage and larger interval score (Fig. 2, Table 1). Differences between the EHMOS and t-test method performances in scenarios with unbalanced sampling design were exacerbated by the presence of generalist species (scenario A2xB1xC2 and A2xB2xC2 in Fig. 2 and Table 1). The t-test had a high average negative bias of -30 m and -22 m, small average interval coverage of 68% and 83%, and large interval scores of 360 and 217 in scenarios A2xB1xC2 and A2xB2xC2. In comparison, EHMOS is relatively unbiased, with an average bias of 3 m and 0 m, had average interval coverage superior to 95% for both scenarios, and had smaller interval scores: 117 and 150 for scenarios A2xB1xC2 and A2xB2xC2 respectively.

Table 1: Average species specific interval scores for the eight simulated scenarios after 30 replications. Numbers in brackets represent standard deviations.

Scenario	EHMOS	cGLMM	t-test
Uniform sampling			
A1xB1xC1	69.48 (3.5)	72.32 (3.58)	100.63 (4.01)
A1xB1xC2	113.39 (12.56)	116.17 (11.72)	148.46 (22.83)
A1xB2xC1	136.17 (24.66)	4558.69 (1929.05)	346 (120.27)
A1xB2xC2	192.37 (31.53)	5054.91 (3201.19)	281.66 (90.63)
Unbalanced sampling			
A2xB1xC1	53.57 (2.61)	53.28 (2.62)	90.11 (10.58)
A2xB1xC2	117.07 (14.94)	125.35 (16.29)	359.86 (90.49)
A2xB2xC1	85.46 (14.41)	288 (179.88)	109.22 (13.4)
A2xB2xC2	150.14 (24.76)	945.45 (507.25)	216.9 (55.42)

Application

On average, species shifted upslope (Fig. 3), with the mean shift ranging from 124 m (SD = 95 m) and 173 m (SD = 127 m) to 183 m (SD = 80 m) based on the estimates from the t-test, cGLMM and EHMOS respectively. All significant shifts were upslope for the three methods, with 11, 8 and 10 species having significant shifts based respectively on EHMOS, cGLMM and t-test credible or confidence intervals. However, we observed heterogeneity between species, with estimated shifts ranging from 19 m (*Stauroderus scalaris*) to 344 m (*Bicolorana bicolor*) based on EHMOS, -124 m (*Platycleis albopunctata*) to 419 m (*Omocestus haemorrhoidalis*) based on cGLMM estimates and, -30 m (*Platycleis albopunctata*) to 356 m (*Bicolorana bicolor*) based on t-test results.

The three methods provided heterogeneous shift estimates (Fig. 3). T-test estimates seemed to be smaller than those estimated by cGLMM or EHMOS. To verify this observation, we performed a paired t-test on species point estimates between methods. We found significant differences between the t-test and cGLMM estimates (mean difference = -49, $t(23) = -2.62$, $p = 0.02$) and EHMOS estimates (mean difference = -60, $t(23) = -5.51$, $p < 0.01$), confirming our observation on the graph. Methods also differed in the precision of estimates. The t-test had on average narrower confidence intervals with a mean confidence interval width of 271 m, compared to EHMOS (371 m) and cGLMM (4701 m). This high value for cGLMM is due to unreliably wide credible intervals for nine species (Fig. 3). Those unrealistic estimates could be explained by quadratic coefficients, the denominator in optimum formulae (Eq. 2), estimated close to zero and generating unreliable optimum estimates and thus shift estimates. This seemed to occur for taxa having their optimum close to the sampling edge (before the 10th or beyond the 90th percentile of the sampling site elevations, depicted as orange rectangles in Fig. 3), for example, *Yersinella sp.* and *Gomphocerus sibiricus sibiricus*, or species having low maximum occupancy probability, e.g. *Omocestus haemorrhoidalis*.

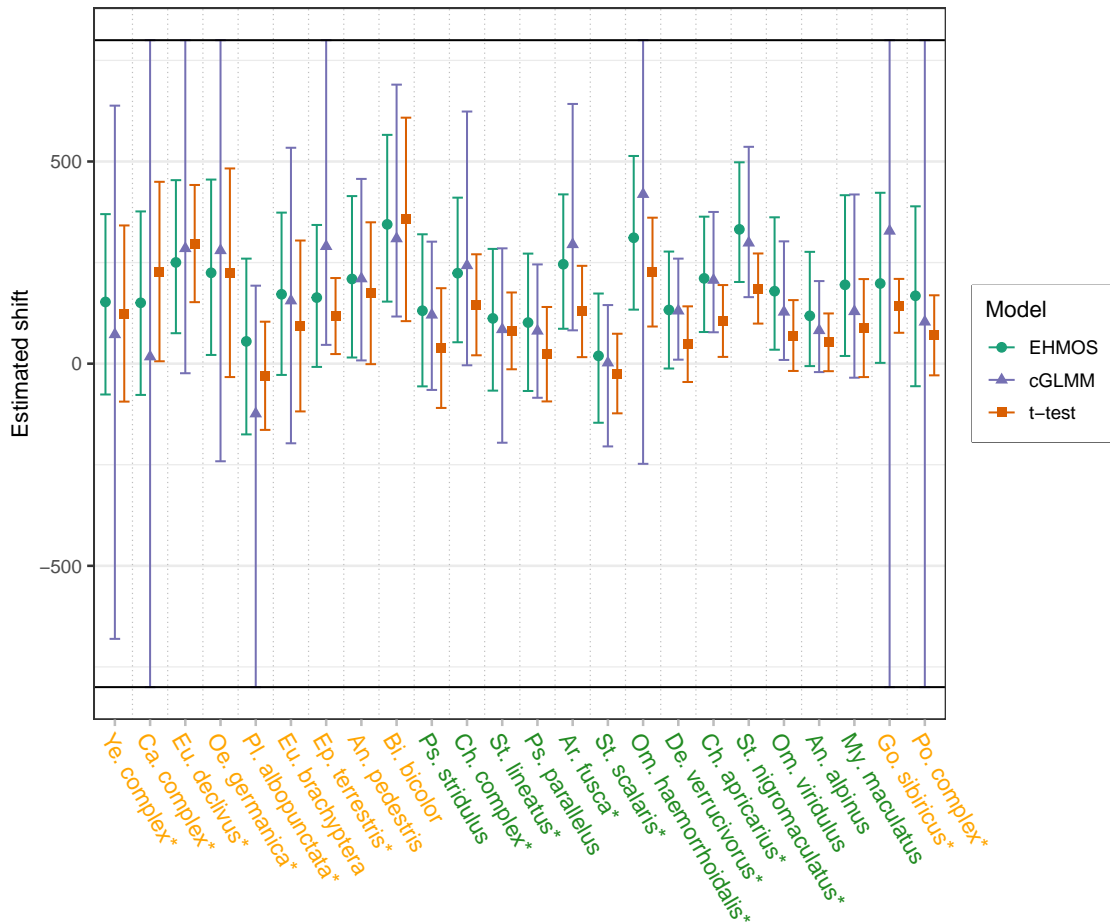


Figure 3: Shift estimates following three statistical methods on 24 Orthoptera species optimum positions between 1980–1988 and 2018–2019 along an elevational gradient in Mercantour National Park (French Alps). The points correspond to the mean estimates and the segments to the associated 95% credible or confidence intervals. The black horizontal lines represent graphical range, segments reaching those lines indicate intervals wider than the graphical range. Species are ordered by increasing optimum estimates of the first period. Species names are coloured depending on the optimum position types with edge species coloured in orange and middle species in green. Presence (or absence) of asterisk indicates species ecological type as defined in the simulation study (presence: generalist, absence: specialist).

Discussion

In this study, we performed simulations to test the accuracy of optimum shift estimates of a new model, the Explicit Hierarchical Model of Optimum Shifts (EHMOS), compared to a more standard GLMM and method based on the comparison of the mean. The simulation study provides evidence that EHMOS generally performed better than the two other methods. cGLMM performed identically to EHMOS when species had their optimum in the middle of the sampling gradient,

but was inaccurate and very imprecise when some species had their optimum close to the sampling limits. We also found that the mean comparison method tended to be negatively biased, with low interval coverage in scenarios involving edge species or unbalanced sampling design, especially when generalist species were included. The case study confirmed that the three methods provide different results. cGLMM provided unreliable estimates for nine of the 24 Orthoptera species, and the mean comparison method had on average smaller shift estimates than the other two methods. On average, we found that species had shifted upslope by 124 m (SD = 95 m) and 173 m (SD = 127 m) to 183 m (SD = 80 m) based on estimates from the t-test, cGLMM and EHMOS respectively.

EHMOS always performed better than the t-test-based method. Differences between these two methods were particularly large under unbalanced sampling design, and especially in the presence of generalist species. It is predictable that mean estimates are biased by unbalanced sampling designs, in particular in scenarios where a part of the gradient is preferentially sampled leading to bias in optimum estimates towards this part (ter Braak & Looman, 1986). We simulated one example where the middle part of the gradient was preferentially sampled, but similar results are expected in other preferential sampling design, e.g. preferential sampling towards low elevations due to the reduction in surface area with elevation (Lenoir et al., 2008). In such situations, both optima estimates could be biased towards the preferentially sampled part of the gradient and shift estimates could be biased towards zero. This explains the negative bias produced by t-test method in our simulated scenarios involving unbalanced sampling design. It is also not surprising that bias due to unbalanced sampling design is larger for generalist species than for specialist species. Bias is more influenced by the distribution of sampling sites occupied by the species rather than by the entire distribution of sampling sites. Occupied sites of a specialist species with a narrow range could be uniformly distributed even if the entire distribution of sampling sites is uneven. In contrast, the larger the species range, the more the occupied site distribution on the gradient will match the entire sampling site distribution, leading to biased shift estimates. A slightly larger bias with the t-test than with EHMOS was also observed in scenarios involving edge species, even with balanced sampling. This could be

explained by the sensitivity of t-tests to species response curve truncation (ter Braak & Looman, 1986). In this case, the species–environment relationship is only partially observed, and the optimum is underestimated as the species range is not entirely covered (ter Braak & Looman, 1986). Bias will thus increase with the magnitude of truncation. Hence, edge species that shifted towards the sampling range margins will have their second optimum more severely biased than the first, inducing underestimation in the shift estimates. Such a bias that depends on the species type could lead to misleading conclusions when comparing magnitudes of shifts between species types. Previous studies found that edge species shifted less than middle species (Rumpf et al., 2018). While this finding is ecologically consistent, it could have been exacerbated by the bias associated to the use of the mean comparison method.

Generalized linear models have been extensively used in species distribution modelling (Elith et al., 2010), and also in optimum shift modelling (e.g. Coudun & Gégout, 2005; Lenoir et al., 2008). Usually, multiple GLMs are conducted separately for each species and each survey that is to be compared (but see Lembrechts et al., 2017 using GLMM). This approach has the inconvenient of losing information and thus precision in estimation by splitting the data into multiple data sets. Estimates for data-poor species may for instance be improved with a GLMM, by borrowing information from data-rich species (Ovaskainen & Soininen, 2011). Hence, one could assume that GLM would have, at best equal and probably worse performances than cGLMM. We thus chose to compare EHMOS only with cGLMM to show that EHMOS has advantages against classical GLM-based approaches. We worked with a Bayesian framework that allows us to derive uncertainty measures associated with shift estimates. Note that it would have been possible to compute the asymptotic variance of optima or shift estimates using a frequentist approach, e.g. the delta method (Urli et al., 2014). The Bayesian GLMM we developed here performed as well as the EHMOS and better than the mean comparison method in scenarios where all species have optima in the middle of the sampling range. However, as expected, for edge species, the standard GLMM failed to estimate bell shape curves, fitting sigmoid curves instead and making derived optima and shifts unreliable (ter Braak & Looman, 1986; Citores et al., 2020). In contrast,

EHMOS, thanks to its ecological formulation and its explicit modelling of optimum shifts, allowed estimates of optimum shifts for edge species. Hence, edge species could be kept if one is using EHMOS, while they should be removed in t-test or GLM-based analysis. Ecological formulation of EHMOS, especially explicit modelling of shift, will also allow easier choice of prior in a Bayesian setting. cGLMM parameters should be transformed to be ecologically interpretable (see Appendix S1). Such a transformation may change an uninformative prior on the parameter scale into an informative prior on the ecological scale, making the choice of prior more complicated (see Lemoine, 2019 for an example in logistic regression). Hence, the standard GLMM appeared to be less flexible than EHMOS.

In addition to having more robust point estimate accuracy whatever the sampling design, ecological marginality and specialization than a t-test or cGLMM, EHMOS also seemed to produce more precise estimates. cGLMM was very sensitive to ecological marginality in terms of precision, providing very wide and uninformative credible intervals for edge species. The interval coverage of the t-test was also very sensitive to ecological marginality, sampling design and ecological specialization. Having informative precision estimates and reporting them is important when studying species shifts (Bates et al., 2015; Taheri et al., 2021). Some authors use the overlap of confidence intervals as a surrogate for a test of significant difference (Lenoir et al., 2008; Urli et al., 2014), so low interval coverage could lead to misleading conclusions, and the t-test method should be avoided.

In our case study, we studied optimum shifts between two surveys conducted 30 years apart of 24 Orthoptera species along an elevational gradient in the French Alps. All methods suggested that on average Orthoptera species shifted their optimum towards higher elevations, and that the magnitude of the shift varies between species. These results are consistent with previous findings on other insect taxa, which on average have shifted upslope, but at different paces depending on the species (Vitasse et al., 2021). Even if the magnitude of shifts varied between methods, our estimates are in line with the estimated warming rate of $0.36 \text{ }^\circ\text{C.decade}^{-1}$ in the Alps between 1970 and 2019, which roughly corresponds to a mean shift of about $+62\text{--}71 \text{ m decade}^{-1}$ (Vitasse et al. 2021), leading to an expected mean shift of around $+186\text{--}213 \text{ m}$ between the two surveys. Thus, the average shift estimated

by EHMOS (183 m) is closer to what is expected than the t-test, suggesting it may provide a more reliable estimate of optimum shift. Yet this observation has to be taken carefully, as the warming rate could vary between regions, and species might respond at a slower rate than expected (Vitasse et al., 2021). However, as the case study is close to the scenario A2xB2xC2 of our simulations, with the presence of both generalist and edge species sampled with an unbalanced sampling design, we could suspect an underestimation of optimum shifts by the t-test as observed in these simulations. It was also not surprising to obtain unreliable wide credible intervals for most of the edge species with the GLMM, as the simulations provide such evidence.

Notwithstanding its demonstrated comparative advantages, EHMOS could have some limitations. It might produce optimum estimates outside the sampling range and thus, depending on the sampling design, outside the environmental range. This property could be an advantage in the case of niche truncation. However, optimum estimates should be limited to the environmental range: for example, by imposing an appropriate prior. One potential other limitation, also valid for cGLMM, is the assumption that species response curves are symmetric. While this assumption could stand for numerous species (Rydgren et al., 2003), some species could have asymmetrical relationships with environment gradients (Oksanen & Minchin, 2002). We found through a simulation study that the three methods could be biased in presence of asymmetrical species response curves (see Appendix S4). More precisely, the three methods provided biased optimum shift estimates when species response curves changed of shapes between the sampling occasions. For instance, the three methods produced positively biased optimum shift estimates for species having a symmetrical response curve in the first sampling occasion and a right-skewed species response curve in the second sampling occasion, situations that could arise due to range margin contraction or expansion for instance. In contrast, optimum shift should be accurately estimated by the three methods for species having the same degree of skewness in both sampling occasions. Hence, shape of species response curves should be inspected carefully prior to analysis. In case where response curves are asymmetric, some alternative modelling methods have been proposed, such as generalized additive models (GAM, e.g. Heegaard,

2002) or Huisman-Olff-Fresco approach (HOF, e.g. Oksanen & Minchin, 2002), and applied in optimum shift modelling (Maggini et al., 2011; Urli et al., 2014). However, previous simulation studies have shown that both methods could fail to retrieve the true shape of species response (Jansen & Oksanen, 2013; Michaelis & Diekmann, 2017). Both methods appeared to be sensitive to gradient truncation and sampling intensity (Jansen & Oksanen, 2013; Citores et al., 2020), limiting their use in unbalanced sampling design and for edge species. Further research is thus needed to find more robust methods to detect and deal with asymmetrical response curves. Finally, we present the EHMOS within a hierarchical model that did not incorporate observation process. It is well known that imperfect detection could bias occupancy estimates, and that this should be addressed in modelling whenever it is possible. While our example does not allow the use of multi-species occupancy models, hierarchical linear models estimating both species occupancy and detection (Dorazio & Royle, 2005), due to the absence of sampling replication in historical data, it is straightforward to extend our framework in a MSOM by adding the observation process in the hierarchy of the model.

EHMOS presents some advantages and potential extensions that we did not explore through our simulation analysis that may be interesting for future research. Its explicit formulation of optimum shifts could allow direct testing of ecological assumptions about inter-specific variability: for example, assumptions about species traits that may affect the magnitude of shifts. Species trait effects could be added directly in the model through the random effect specification of shift (Jamil et al., 2014). The potential effects of species ecological marginality or ecological specialization on magnitude of shift could also be tested by adding the correlation between optimum and shift, or tolerance and shift parameters respectively. Finally, while we focused on the particular case of sampling the same sites at two sampling occasions, we assumed that EHMOS would perform equally in the case of different sites sampled at each sampling occasion. It could also be easily extended to deal with more than two sampling occasions, and test for changes in rates of optimum shifts between multiple sampling occasions. EHMOS thus appears a highly flexible method that could be used in various fields.

References

- Bates A.E., Bird T.J., Stuart-Smith R.D., Wernberg T., Sunday J.M., Barrett N.S., Edgar G.J., Frusher S., Hobday A.J., Pecl G.T., Smale D.A., & McCarthy M. (2015) Distinguishing geographical range shifts from artefacts of detectability and sampling effort. *Diversity and Distributions*, **21**, 13–22.
- Citores L., Ibaibarriaga L., Lee D.-J., Brewer M.J., Santos M., & Chust G. (2020) Modelling species presence–absence in the ecological niche theory framework using shape-constrained generalized additive models. *Ecological Modelling*, **418**, 108926.
- Conn P.B., Johnson D.S., Williams P.J., Melin S.R., & Hooten M.B. (2018) A guide to bayesian model checking for ecologists. *Ecological Monographs*, **88**, 526–542.
- Coudun C. & Gégout J.-C. (2005) Ecological behaviour of herbaceous forest species along a pH gradient: A comparison between oceanic and semicontinental regions in northern France. *Global Ecology and Biogeography*, **14**, 263–270.
- Dorazio R.M. & Royle J.A. (2005) Estimating Size and Composition of Biological Communities by Modeling the Occurrence of Species. *Journal of the American Statistical Association*, **100**, 389–398.
- Elith J., Kearney M., & Phillips S. (2010) The art of modelling range-shifting species. *Methods in Ecology and Evolution*, **1**, 330–342.
- Felde V.A., Kapfer J., & Grytnes J.-A. (2012) Upward shift in elevational plant species ranges in Sikkildalen, central Norway. *Ecography*, **35**, 922–932.
- Freeman B.G., Scholer M.N., Ruiz-Gutierrez V., & Fitzpatrick J.W. (2018) Climate change causes upslope shifts and mountaintop extirpations in a tropical bird community. *Proceedings of the National Academy of Sciences*, **115**, 11982–11987.
- Gelman A., Carlin J.B., Stern H.S., Dunson D.B., Vehtari A., & Rubin D.B. (2013) *Bayesian Data Analysis*. Chapman; Hall, New York.
- Heegaard E. (2002) The outer border and central border for species–environmental relationships estimated by non-parametric generalised additive models. *Ecological Modelling*, **157**, 131–139.
- 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.
- Huisman J., Olff H., & Fresco L. (1993) A hierarchical set of models for species response analysis. *Journal of Vegetation Science*, **4**, 37–46.
- Jamil T., Kruk C., & ter Braak C.J.F. (2014) A Unimodal Species Response Model Relating Traits to Environment with Application to Phytoplankton Communities. *PLoS ONE*, **9**, e97583.
- Jansen F. & Oksanen J. (2013) How to model species responses along ecological gradients—huisman–olff–fresco models revisited. *Journal of Vegetation Science*, **24**, 1108–1117.
- Kellner K. & Meredith M. (2021) jagsUI: A Wrapper Around 'rjags' to Streamline 'JAGS' analyses.
- Lembrechts J.J., Alexander J.M., Cavieres L.A., Haider S., Lenoir J., Kueffer C., McDougall K., Naylor B.J., Nuñez M.A., Pauchard A., Rew L.J., Nijs I., & Milbau A. (2017) Mountain roads shift native and non-native plant species' ranges. *Ecography*, **40**, 353–364.
- Lemoine N.P. (2019) Moving beyond noninformative priors: Why and how to choose weakly informative priors in bayesian analyses. *Oikos*, **128**, 912–928.
- Lenoir J., Gégout J.-C., Marquet P.A., Ruffray P. de, & Brisse H. (2008) A significant upward shift in plant species optimum elevation during the 20th century. *Science*, **320**, 1768–1771.
- Maggini R., Lehmann A., Kéry M., Schmid H., Beniston M., Jenni L., & Zbinden N. (2011) Are Swiss birds tracking climate change?: Detecting elevational shifts using response curve shapes. *Ecological Modelling*, **222**, 21–32.
- Michaelis J. & Diekmann M.R. (2017) Biased niches – Species response curves and niche attributes from Huisman-Olff-Fresco models change with differing species prevalence and frequency. *PLOS ONE*, **12**, e0183152.
- Mourguiart B., Couturier T., Braud Y., Mansons J., Combrisson D., & Besnard A. (2021) Multi-species occupancy models: An effective and flexible framework for studies of insect communities. *Ecological Entomology*, **46**, 163–174.
- Moussus J.-P., Julliard R., & Jiguet F. (2010) Featuring 10 phenological estimators using simulated data. *Methods in Ecology and Evolution*, **1**, 140–150.

- Oksanen J. & Minchin P.R. (2002) Continuum theory revisited: What shape are species responses along ecological gradients? *Ecological Modelling*, **157**, 119–129.
- Ovaskainen O. & Soininen J. (2011) Making more out of sparse data: Hierarchical modeling of species communities. *Ecology*, **92**, 289–295.
- Pecl G.T., Araújo M.B., Bell J.D., Blanchard J., Bonebrake T.C., Chen I.-C., Clark T.D., Colwell R.K., Danielsen F., Evengård B., & others (2017) Biodiversity redistribution under climate change: Impacts on ecosystems and human well-being. *Science*, **355**, eaai9214.
- Plummer M. (2003) JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. *Working Papers*, 8.
- R Core Team (2018) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rumpf S.B., Hülber K., Klöner G., Moser D., Schütz M., Wessely J., Willner W., Zimmermann N.E., & Dullinger S. (2018) Range dynamics of mountain plants decrease with elevation. *Proceedings of the National Academy of Sciences*, **115**, 1848–1853.
- Rydgren K., Økland R.H., & Økland T. (2003) Species response curves along environmental gradients. A case study from SE norwegian swamp forests. *Journal of Vegetation Science*, **14**, 869–880.
- Shoo L.P., Williams S.E., & Hero J.-M. (2006) Detecting climate change induced range shifts: Where and how should we be looking? *Austral Ecology*, **31**, 22–29.
- Strebel N., Kéry M., Schaub M., & Schmid H. (2014) Studying phenology by flexible modelling of seasonal detectability peaks. *Methods in Ecology and Evolution*, **5**, 483–490.
- Taheri S., Naimi B., Rahbek C., & Araújo M.B. (2021) Improvements in reports of species redistribution under climate change are required. *Science Advances*, **7**, eabe1110.
- Taylor C., Caplat P., Massimino D., Johnston A., Jonzén N., Smith H.G., & Lindström Å. (2015) Swedish birds are tracking temperature but not rainfall: Evidence from a decade of abundance changes. *Global Ecology and Biogeography*, **24**, 859–872.

- ter Braak C.J.F. & Looman C.W.N. (1986) Weighted averaging, logistic regression and the Gaussian response model. *Vegetatio*, **65**, 3–11.
- Thorson J.T., Pinsky M.L., & Ward E.J. (2016) Model-based inference for estimating shifts in species distribution, area occupied and centre of gravity. *Methods in Ecology and Evolution*, **7**, 990–1002.
- Urli M., Delzon S., Eyermann A., Couallier V., García-Valdés R., Zavala M.A., & Porté A.J. (2014) Inferring shifts in tree species distribution using asymmetric distribution curves: A case study in the Iberian mountains. *Journal of Vegetation Science*, **25**, 147–159.
- Vitasse Y., Ursenbacher S., Klein G., Bohnenstengel T., Chittaro Y., Delestrade A., Monnerat C., Rebetez M., Rixen C., Strebel N., Schmidt B.R., Wipf S., Wohlgemuth T., Yoccoz N.G., & Lenoir J. (2021) Phenological and elevational shifts of plants, animals and fungi under climate change in the European Alps. *Biological Reviews*, **96**, 1816–1835.
- Wilson R.J., Gutiérrez D., Gutiérrez J., Martínez D., Agudo R., & Monserrat V.J. (2005) Changes to the elevational limits and extent of species ranges associated with climate change. *Ecology Letters*, **8**, 1138–1146.
- Zipkin E.F., DeWan A., & Royle J.A. (2009) Impacts of forest fragmentation on species richness: A hierarchical approach to community modelling. *Journal of Applied Ecology*, **46**, 815–822.

Appendix S1: Description of EHMOS' parameters

In the occupancy model of EHMOS (see Eq. (3) in the main text), the four parameters (α , θ , δ and τ) represent ecological descriptors of the species-environment relationship (Fig. S1):

- α corresponds to the maximum probability of occurrence, ψ_{max} , on the logit scale. Indeed, the maximum probability of occurrence is reached when gradient value is equal to the optimum, i.e. we have $\psi_j = \psi_{max}$ when $X_j = \theta$. Thus, Eq. (3) in the main text becomes:

$$\log\left(\frac{\psi_{max}}{1 - \psi_{max}}\right) = \alpha \iff \psi_{max} = \frac{1}{1 + \exp(-\alpha)}; \quad (4)$$

- θ is the species optimum, the environmental value at which species reached its maximum probability of occupancy;
- δ is the species shift between two optima;
- τ represents the environmental tolerance of a species, i.e. the gradient range that a species can occupy. It can be related to the width, ω , of the species response curve at a specified occupancy probability threshold, p_ω , by:

$$\omega = 2 \times \tau \sqrt{2 \times (\alpha - \text{logit}(p_\omega))} \iff \tau = 0.5 \times \frac{\omega}{\sqrt{2 \times (\alpha - \text{logit}(p_\omega))}} \quad (5)$$

with a p_ω usually specified at 0.05 to estimate the suitable range of a species (Michaelis & Diekmann, 2017).

Those four EHMOS parameters can also be derived from cGLMM coefficient estimates (Jamil et al., 2014):

$$\alpha = \beta_0 - \frac{\beta_1^2}{4\beta_2}$$

$$\theta = -\frac{\beta_1}{2\beta_2}$$

$$\delta = -\frac{\beta_1}{2\beta_2} - \left(-\frac{\beta_1 + \beta_4}{2(\beta_2 + \beta_5)}\right)$$

$$\tau = \sqrt{-\frac{1}{2\beta_2}}$$

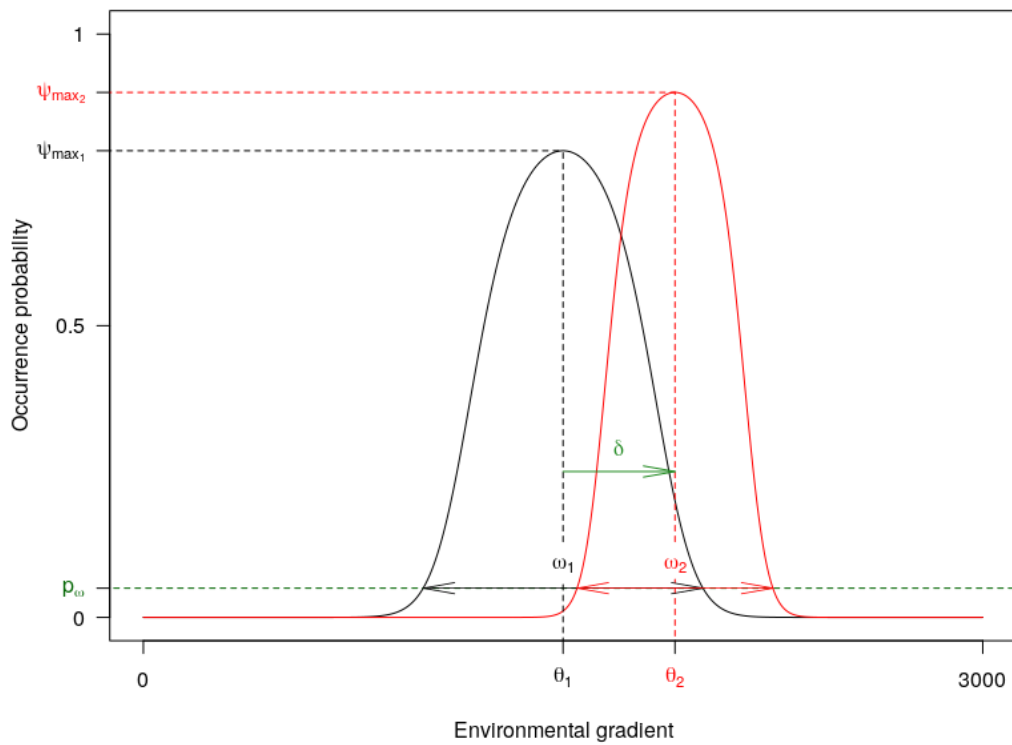


Figure S1.1: Schematic representation of ecological parameters that describes species-environment relationship.

References

Jamil T., Kruk C., & Braak C.J.F. ter (2014) A Unimodal Species Response Model Relating Traits to Environment with Application to Phytoplankton Communities. *PLoS ONE*, **9**, e97583.

Michaelis J. & Diekmann M.R. (2017) Biased niches – Species response curves and niche attributes from Huisman-Olff-Fresco models change with differing species prevalence and frequency. *PLOS ONE*, **12**, e0183152.

Appendix S2: Description of performance metrics used in the simulation study

The four performance metrics used (bias, RMSE, interval coverage and interval score) are computed for each model at different levels (species, scenario, overall). First, for each species i included in a particular scenario s , we computed average performance metrics over the thirty replications:

- Bias of species i in scenario s based on method m is defined as the average of differences between absolute¹ estimates based on method m and absolute true shift over the thirty replications:

$$Bias_{i,s,m} = \frac{1}{R} \sum_{r=1}^R |\hat{\delta}_{i,s,m,r}| - |\delta_{i,s}| ;$$

- The RMSE at species-specific level is defined by:

$$RMSE_{i,s,m} = \sqrt{\frac{1}{R} \sum_{r=1}^R (\hat{\delta}_{i,s,m} - \delta_{i,s})^2} ;$$

- The interval coverage for a species in a particular scenario is the ratio between the number of replications in which confidence interval estimate contain the true shift and the total number of replications;
- The interval score of species i in scenario s based on method m is defined as the average of interval scores obtained at each replication r . At the replication level, interval score is computed as follows:

$$IS_{i,s,m,r}(l_{i,s,m,r}, u_{i,s,m,r}; \delta_{i,s}) = \begin{cases} (u_{i,s,m,r} - l_{i,s,m,r}) + \frac{2}{\alpha}(l_{i,s,m,r} - \delta_{i,s}) & \text{if } \delta_{i,s} < l_{i,s,m,r} \\ (u_{i,s,m,r} - l_{i,s,m,r}) & \text{if } u_{i,s,m,r} > \delta_{i,s} > l_{i,s,m,r} \\ (u_{i,s,m,r} - l_{i,s,m,r}) + \frac{2}{\alpha}(\delta_{i,s} - u_{i,s,m,r}) & \text{if } \delta_{i,s} > u_{i,s,m,r} \end{cases}$$

where α is the confidence level of the interval, here set at 0.05, $l_{i,s,m,r}$ and $u_{i,s,m,r}$

¹Absolute values are took because true shifts could either be positive or negative, thus a negative bias always represents an underestimation of the true shif magnitude whatever its direction.

are the lower and upper limits of the interval estimated by model m for species i in scenario s and replicate r .

At the scenario level, metrics are averaged over the twenty species:

$$\overline{Bias}_{s,m} = \frac{1}{N} \sum_{i=1}^N Bias_{i,s,m}$$

$$\overline{RMSE}_{s,m} = \frac{1}{N} \sum_{i=1}^N RMSE_{i,s,m}$$

$$\overline{Coverage}_{s,m} = \frac{1}{N} \sum_{i=1}^N Coverage_{i,s,m}$$

$$\overline{IS}_{s,m} = \frac{1}{N} \frac{1}{R} \sum_{i=1}^N \sum_{r=1}^R IS(\hat{l}_{i,s,m,r}, \hat{u}_{i,s,m,r}; \delta_{i,s})$$

Appendix S3: Complementary results of the simulation study

Table S3. 1: Number of simulation replications per scenario that at least one model parameter had an R-hat higher than 1.1 after the maximum of iterations was reached

Scenario	EHMOS	cGLMM
Uniform sampling		
A1xB1xC1	0	26
A1xB1xC2	0	2
A1xB2xC1	0	30
A1xB2xC2	0	30
Unbalanced sampling		
A2xB1xC1	0	0
A2xB1xC2	0	0
A2xB2xC1	0	26
A2xB2xC2	0	5

Table S3. 2: Average computation times in hours.

Scenario	EHMOS	cGLMM
Uniform sampling		
A1xB1xC1	5.7	39.9
A1xB1xC2	5.0	17.8
A1xB2xC1	5.1	42.9
A1xB2xC2	4.4	34.4
Unbalanced sampling		
A2xB1xC1	5.6	8.5
A2xB1xC2	5.6	6.7
A2xB2xC1	5.7	35.8
A2xB2xC2	3.8	19.2

Table S3. 3: Average of species specific performance metrics for the eight simulated scenarios after 30 replications. Numbers in brackets represent standard deviations.

Scenario	EHMOS	cGLMM	t-test
Bias			
A1xB1xC1	1.26 (1.3)	-0.49 (1.61)	-0.64 (3.86)
A1xB1xC2	0.25 (2.48)	-0.29 (2.65)	2.1 (5.3)
A1xB2xC1	-3.52 (3.99)	42.59 (18.76)	-12.72 (7.61)

A1xB2xC2	-5.52 (6.58)	21.84 (11.52)	-13.05 (8.45)
A2xB1xC1	1.47 (0.98)	1.02 (1.08)	-9.26 (3.28)
A2xB1xC2	2.91 (2.29)	3.31 (2.5)	-29.87 (7.74)
A2xB2xC1	3.15 (1.85)	6.21 (4.41)	-10.41 (3.85)
A2xB2xC2	0.1 (5.11)	25.15 (11.47)	-21.78 (6.96)

Coverage

A1xB1xC1	0.94 (0.01)	0.94 (0.01)	0.96 (0.02)
A1xB1xC2	0.96 (0.01)	0.95 (0.01)	0.96 (0.03)
A1xB2xC1	0.96 (0.01)	0.95 (0.01)	0.76 (0.09)
A1xB2xC2	0.97 (0.01)	0.96 (0.01)	0.82 (0.07)
A2xB1xC1	0.95 (0.01)	0.96 (0.01)	0.9 (0.04)
A2xB1xC2	0.95 (0.01)	0.92 (0.01)	0.68 (0.08)
A2xB2xC1	0.96 (0.01)	0.96 (0.01)	0.9 (0.04)
A2xB2xC2	0.96 (0.01)	0.96 (0.01)	0.83 (0.06)

IS

A1xB1xC1	69.48 (3.5)	72.32 (3.58)	100.63 (4.01)
A1xB1xC2	113.39 (12.56)	116.17 (11.72)	148.46 (22.83)
A1xB2xC1	136.17 (24.66)	4558.69 (1929.05)	346 (120.27)
A1xB2xC2	192.37 (31.53)	5054.91 (3201.19)	281.66 (90.63)
A2xB1xC1	53.57 (2.61)	53.28 (2.62)	90.11 (10.58)
A2xB1xC2	117.07 (14.94)	125.35 (16.29)	359.86 (90.49)
A2xB2xC1	85.46 (14.41)	288 (179.88)	109.22 (13.4)
A2xB2xC2	150.14 (24.76)	945.45 (507.25)	216.9 (55.42)

RMSE

A1xB1xC1	14.84 (0.49)	15.08 (0.53)	20.58 (1.44)
A1xB1xC2	23.97 (2.43)	24.68 (2.44)	28.58 (3.23)
A1xB2xC1	26.23 (4.23)	78.62 (25.04)	30.42 (4.99)
A1xB2xC2	35.22 (5.41)	59.91 (14.12)	36.45 (5.38)
A2xB1xC1	11.38 (0.43)	11.34 (0.43)	17.69 (1.8)
A2xB1xC2	24.66 (3.25)	25.88 (3.32)	39.65 (5.91)
A2xB2xC1	16.89 (2.46)	25.45 (7.24)	20.3 (2.55)

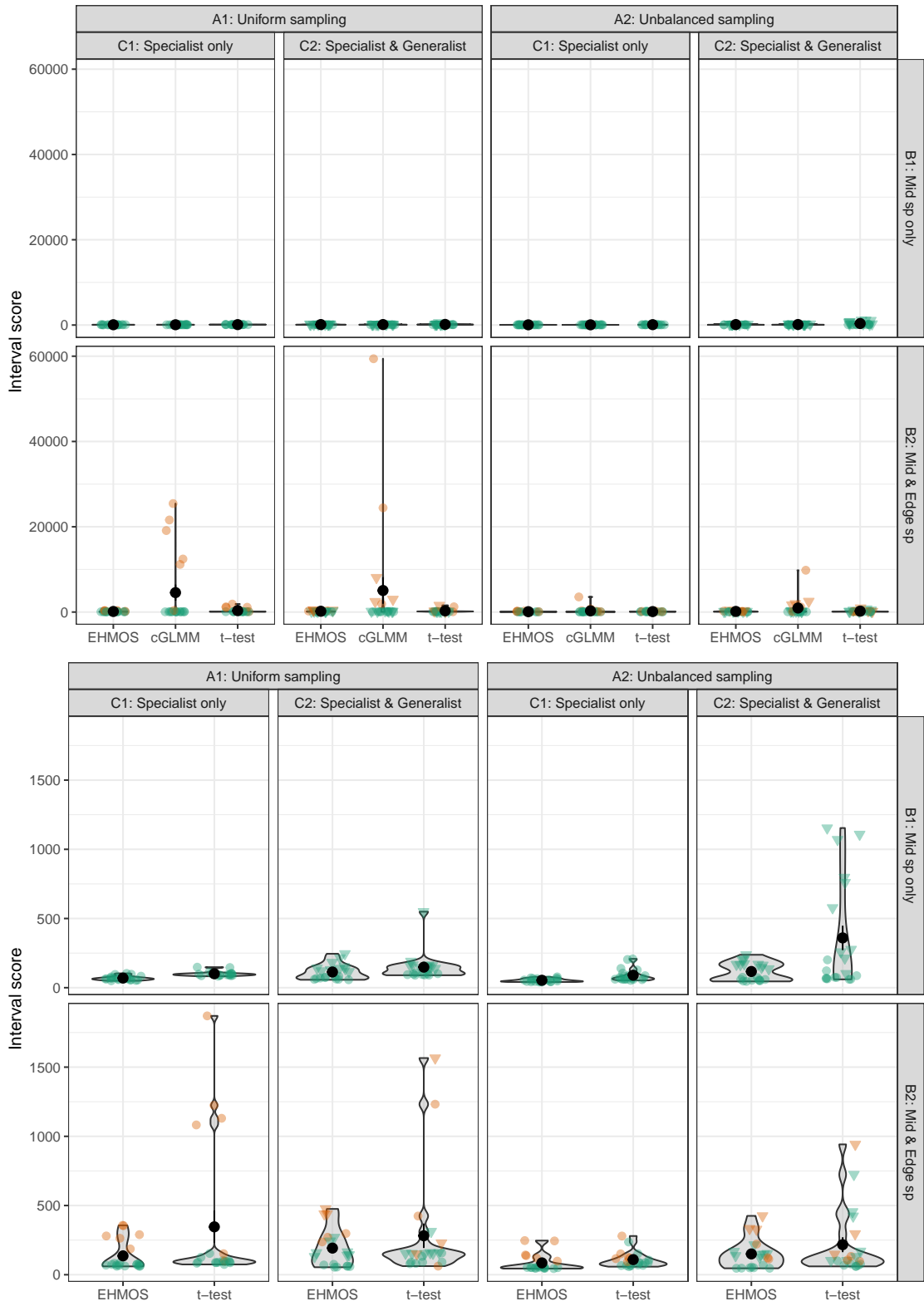


Figure S3.1: Distribution of species specific interval scores (IS) for the eight simulated scenarios after 30 replications, and for the three methods (top graph) or

without cGLMM (bottom graph). Colored points represent average IS in species group depending on their optimum position (placed in the edge or in the middle of sampling range) along the simulated gradient, and their ecological specialization type relative to the width of their ecological niche (generalist species with broad ecological niche or specialist species with narrower ecological niche).

Appendix S4: Departure from symmetric assumption

Models used in the manuscript make the assumption that species have symmetric unimodal relationships with environmental gradient (i.e. occupancy probabilities decrease at the same rate on both side of optimum, see Fig. S4.1). Species-environment relationships could however have other forms (Oksanen & Minchin, 2002; Jansen & Oksanen, 2013). For instance, species response curve (SRC) could be unimodal but skewed to one part of the gradient (i.e. asymmetric). Here we propose to investigate the effect of departure from models' assumption of symmetric unimodal species-gradient relationships on models accuracy by conducting a simulation study.

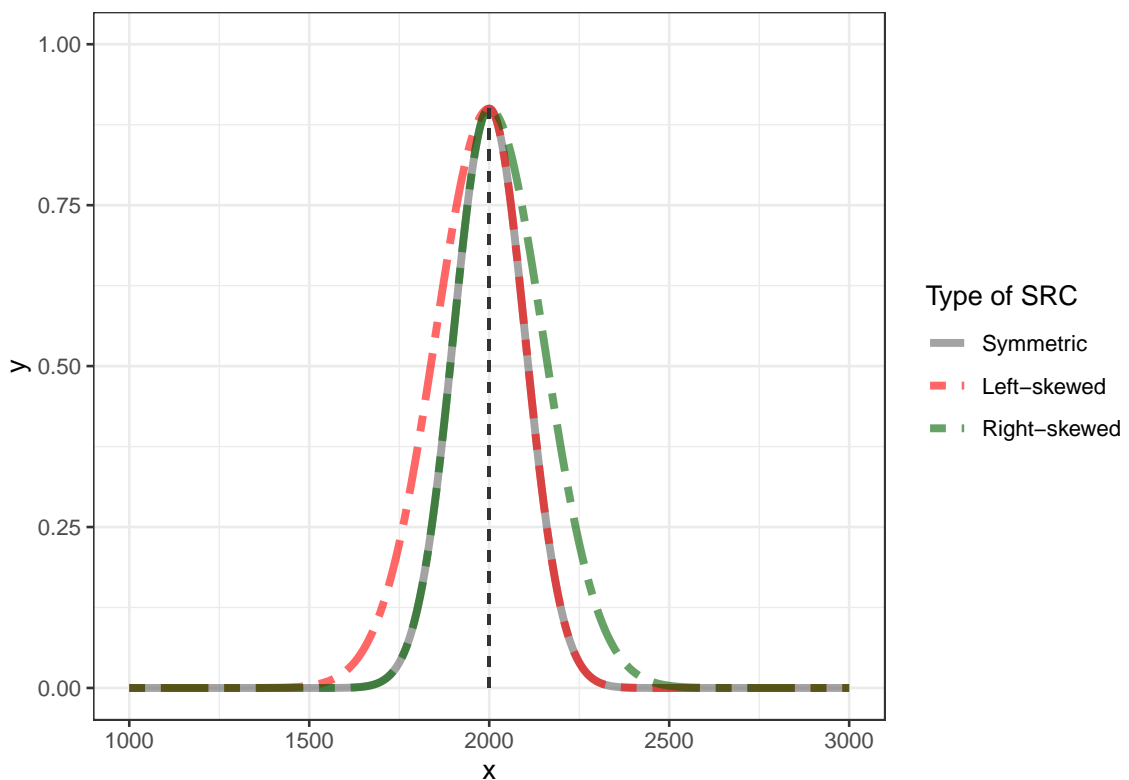


Fig S4.1: Three types of species response curves sharing a same optimum (dashed vertical segment).

Simulated datasets

We defined three potential shapes of species response curves: symmetric, asymmetric with a longer left tail and asymmetric with a longer right tail (Fig. S4.1). We simulated 9 species to describe all possible combinations of SRC that could occur

when two sampling occasions are considered (Fig. S4.2). To keep it simple we chose to only simulate specialist species, shifting upward and having their optima of both sampling occasions in the middle of the sampling range. We simulated curve widths ranged from 400 m to 600 m. Maximum occupancy probabilities were obtained by random draws from a uniform distribution bounded between 0.85 and 0.99. Species optima were randomly positioned between 1490 m and 2460 m, in the middle range of the unbalanced sampling design (i.e. sampling scenario A2 of the manuscript). We assigned only upward shifts following random draws from a uniform distribution bounded between 80 and 120 m (Fig. S4.2).

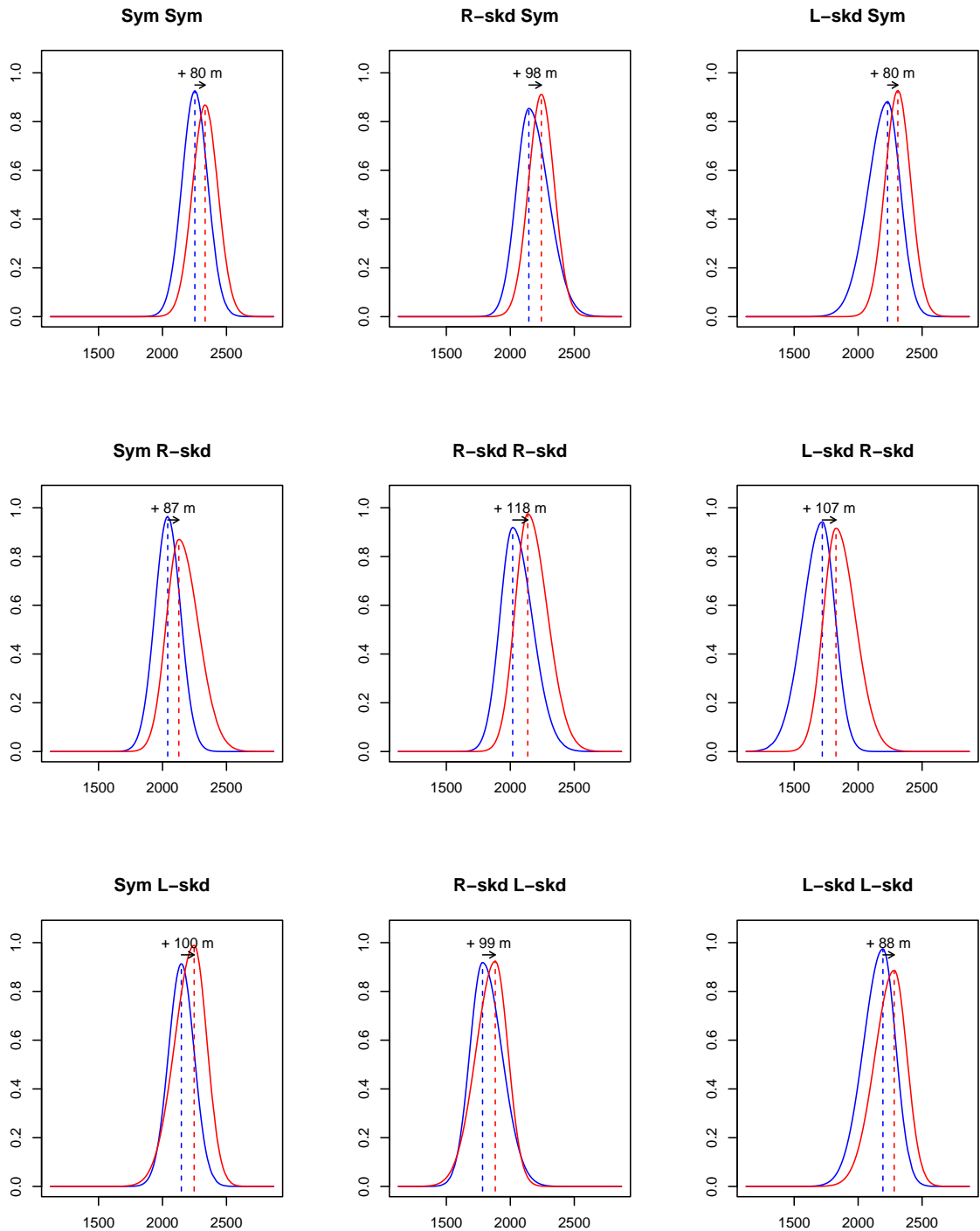


Fig S4.2: Species response curves of the nine simulated species. Blue curves represent the SRCs of the first sampling occasion. Panel titles indicates the shapes of the two SRCs, with sym, r_skd and l_skd standing for symmetric, right-skewed and left-skewed.

Species response curves were simulated using a mixture of two normal distributions with a common mean and two variance parameters (see Eq. (1) below). The common mean represents the species optimum. The two variances are related to

the width of the SRC at each side of the species optimum. The first distribution (resp. second) is used to simulate the left (resp. right) side of the optimum of the SRC. Hence, if the first normal distribution has a greater standard deviation than the second normal distribution, the curve will be left-skewed. If both distributions have the same standard deviation, the SRC is symmetric. Note that we used the density function of the normal distribution multiplied by a scaling parameter to reach a specified probability. The scaling parameter is defined as the ratio of the maximum occupancy probability intended and the maximum density of the normal distribution. We simulated the site-specific probability of occupancy for each species as follows:

$$\psi_{i,j,k}^s = \begin{cases} a_i \frac{\exp\left(-\frac{(x_j^s - (\theta_i + (S_k \times \delta_i)))^2}{2\sigma_1^2}\right)}{\sigma_1 \sqrt{2\pi}} & \text{for } x_j^s \leq \theta_i \\ b_i \frac{\exp\left(-\frac{(x_j^s - \theta_i + (S_k \times \delta_i))^2}{2\sigma_2^2}\right)}{\sigma_2 \sqrt{2\pi}} & \text{for } x_j^s > \theta_i \end{cases} \quad (6)$$

with a_i and b_i the scaling parameters related to the maximum probability, θ the species optimum, δ the species shift and σ the standard deviations related to the width.

Once SRCs were simulated for each species and sampling occasion, we generated presence/absence data, $Y_{i,j,k}^s$, from random draws of a Bernoulli distribution: $Y_{i,j,k}^s \sim \text{Bernoulli}(\psi_{i,j,k}^s)$. We replicated the process 30 times for the two sampling scenarios. Contrary to the manuscript simulation design, here the simulated parameters of the SRCs ($a_i, \theta_i, \delta_i, \sigma_i$) were identical for the two sampling scenarios.

Accuracy of optimum shift estimates under violation of the symmetry assumption

In our simulations, all three methods respond in the same way to departure from assumption of symmetric unimodal species-gradient relationships (Fig. S4.3, Table S4.1). More than the departure of symmetry assumption, the models are sensitive to changes in SRC shapes between the two sampling occasions. The more difference between SRC shapes, the larger the bias for all methods. For instance, the highest bias were obtained for species having a SRC skewed toward one side of the gradient

at the first sampling occasion and a SRC skewed toward the other side in the next sample (e.g. left-skewed to right-skewed). In contrast, methods produced accurate optimum shift estimates for species having same skewness during the two sampling occasions (e.g. species 5).

Table S4. 1: Average of bias in optimum shift estimates for the two simulated scenarios after 30 replications. Numbers in brackets represent standard deviations.

SRC 1	SRC 2	True shift	EHMOS	cGLMM	t-test
Sampling scenario A2					
Symmetric	Symmetric	80	-1 (18)	-2 (16)	-8 (16)
Right-skewed	Symmetric	98	-36 (20)	-31 (18)	-34 (19)
Left-skewed	Symmetric	80	36 (21)	39 (20)	40 (22)
Symmetric	Right-skewed	87	37 (21)	30 (20)	23 (18)
Right-skewed	Right-skewed	118	3 (18)	1 (17)	-8 (14)
Left-skewed	Right-skewed	107	69 (13)	62 (14)	62 (13)
Symmetric	Left-skewed	100	-29 (14)	-36 (15)	-51 (16)
Right-skewed	Left-skewed	99	-75 (17)	-71 (20)	-80 (16)
Left-skewed	Left-skewed	88	0 (17)	-3 (17)	-9 (18)
Sampling scenario A1					
Symmetric	Symmetric	80	2 (21)	-1 (17)	-13 (19)
Right-skewed	Symmetric	98	-40 (17)	-37 (14)	-37 (18)
Left-skewed	Symmetric	80	40 (21)	30 (17)	34 (20)
Symmetric	Right-skewed	87	38 (19)	27 (19)	56 (20)
Right-skewed	Right-skewed	118	-1 (15)	-7 (14)	19 (16)
Left-skewed	Right-skewed	107	70 (23)	52 (24)	65 (22)
Symmetric	Left-skewed	100	-35 (22)	-33 (19)	-33 (26)
Right-skewed	Left-skewed	99	-66 (22)	-54 (21)	-74 (18)
Left-skewed	Left-skewed	88	0 (17)	-1 (15)	4 (20)

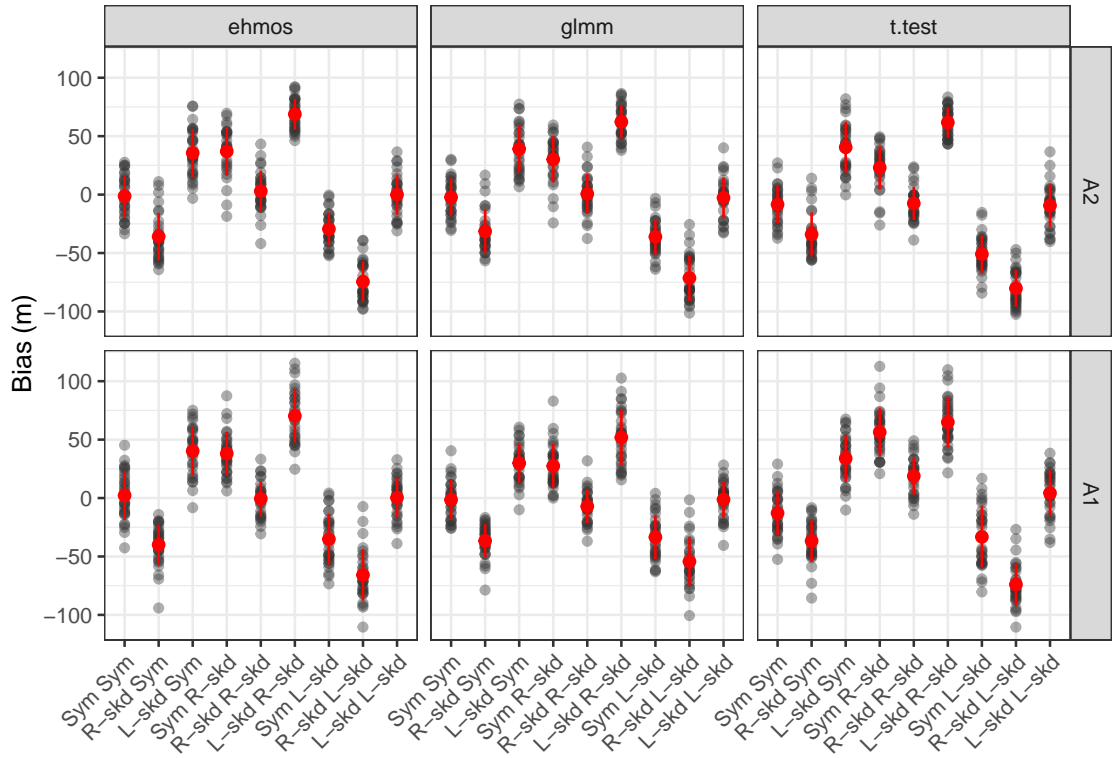


Figure S4.3: Bias in optimum shift estimates with simulated species having different kinds of unimodal species response curves (symmetric, right-skewed or left-skewed). Grey points correspond to results from one replication. Mean bias are represented with their standard deviations by red points and red segments.

Detection of departure from model assumption

We tried two approaches to test the departure from the symmetry assumption made by the three presented models: 1) a Bayesian P-value approach that is often used by ecologists to detect a lack-of-fit in Bayesian models (Conn et al., 2018), and 2) a HOF approach previously used to test for the symmetry assumption in range shift studies (Wilson et al., 2005).

Goodness-of-fit We assess quality of model fit using a Bayesian P-value approach (Gelman et al., 2013) for each species. We used the discrepancy measure defined by Zipkin et al. (2009): $D(y_i) = \sum_k (y_{ik} - \hat{\psi}_{ik})^2$ for the presence-absence observations y_{ik} of species i and their expected values $\hat{\psi}_{ik}$ under the model at MCMC replication k . This discrepancy statistic is computed independently for each species i . A reference distribution is computed by simulating data sets from the posterior distribution, $y_{ik}^{sim} \sim p(\theta_{ik}|y)$, and computing the discrepancy measure, $D(y_{ik}^{sim})$,

for the simulated data sets. The Bayesian p-value for species i is defined as the probability: $p_{Bi} = \Pr(D(y_i) > D(y_i^{sim}))$. Extreme values (e.g. less than 0.05 or greater than 0.95) are interpreted as a lack of fit.

Table S4. 2: Number of times where lack-of-fit were detected after 30 replications for the 9 species simulated with different shapes of response.

SRC.1	SRC.2	cGLMM	EHMOS
Symmetric	Symmetric	1	0
Right-skewed	Symmetric	0	0
Left-skewed	Symmetric	0	0
Symmetric	Right-skewed	0	0
Right-skewed	Right-skewed	0	0
Left-skewed	Right-skewed	0	0
Symmetric	Left-skewed	0	0
Right-skewed	Left-skewed	0	0
Left-skewed	Left-skewed	0	0

We detected no lack-of-fit according to the Bayesian p-value approach, except for one species in one replication of the simulation (Table S4.2).

Bayesian p-value approach seems conservative in regards of skewed species response. Thus, it should not be used to investigate departure from model assumption on shape of species response.

HOF approach HOF approach consists in applying seven models and select the model that best fit the data (Huisman et al., 1993; Jansen & Oksanen, 2013). The five models describe five shapes of response curves assumed to be ecologically meaningful (Oksanen & Minchin, 2002; Jansen & Oksanen, 2013): I) flat, II) monotone, III) plateau, IV) unimodal symmetric, V) unimodal skewed, VI) bimodal symmetric and VII) bimodal skewed. This approach has been used to determine the shape of species response along environmental gradients in many studies (e.g. Oksanen & Minchin, 2002; Jansen & Oksanen, 2013; Michaelis & Diekmann, 2017), and especially in range shift studies Wilson et al. (2005). For instance, Wilson et

al. (2005) used the HOF approach to verify the assumption of symmetrical species response curves made by their two-GLMs approach. We tested this method on our simulated data sets. We applied the HOF approach on each replicated data set and each species for sampling occasion one using the `HOF` function of the R-package `eHOF` (Jansen & Oksanen, 2013).

Table S4. 3: Model choices of the HOF approach for the 9 simulated species after 30 replications.

species	true shape	HOF models		
		IV	V	VII
1	Symmetric (IV)	23	7	0
2	Right-skewed (V)	20	10	0
3	Left-skewed (V)	22	8	0
4	Symmetric (IV)	25	5	0
5	Right-skewed (V)	19	11	0
6	Left-skewed (V)	20	8	2
7	Symmetric (IV)	25	5	0
8	Right-skewed (V)	16	14	0
9	Left-skewed (V)	18	12	0

The HOF approach failed to retrieve the correct shape for numerous species and iterations (Table S4. 3). It estimated symmetrical response in majority of replications even for species having skewed response. Species with symmetrical responses had been misclassified in at least 15% of the replications.

Discussion

We found that departure from assumption of symmetry in species response could induce bias in optimum shift estimates of the three methods. The estimates were biased when degree of skewness in species response changed between the two sampling occasions. Thus, readers should be cautious regarding the shape of species response curves, and more importantly in the potential changes in shapes of SRC between the samples studied, before applying any of the three methods presented

here.

We tried two approaches in order to detect departure from model assumption of symmetry in SRC: 1) Bayesian p-value approach as a measure of goodness-of-fit (Gelman et al., 2013) and 2) HOF approach that were used by Wilson et al. (2005) in this context. None was found to be a reliable evaluation of the symmetry assumption based on our simulations. This could be due to the relatively low degree of skewness of the simulated SRCs, but previous works have also shown that HOF approach could misclassify SRC shapes even with high degree of skewness (Jansen & Oksanen, 2013; Michaelis & Diekmann, 2017). Jansen & Oksanen (2013) found for instance that HOF approach could be sensitive to optimum position along the environmental gradient. Their results indicate numerous misclassifications of HOF approach (see Table 4 and Fig. A3 in Jansen & Oksanen (2013)), especially (but not only) for species having their optimum close to sampling limits (i.e. edge species). These results also stand for GAMs (Jansen & Oksanen, 2013; Citores et al., 2020). Besides, Oksanen & Minchin (2002) pointed out that defining symmetry or asymmetry with GAMs could be subjective and thus advice to use of HOF approach. For these reasons we do not used GAMs for testing the symmetry assumption.

Conclusion

It appears important to thoroughly inspect SRC shape before analysis due to potential bias when SRC shapes change over sampling occasions. More researches are needed to be able to test for asymmetrical SRC, especially simulation studies. Ecologists could also rely on more precise data, e.g. abundance data, when they are available to test the assumption of symmetrical species response curves.

References

- Citores L., Ibaibarriaga L., Lee D.-J., Brewer M.J., Santos M., & Chust G. (2020) Modelling species presence–absence in the ecological niche theory framework using shape-constrained generalized additive models. *Ecological Modelling*, **418**, 108926.
- Conn P.B., Johnson D.S., Williams P.J., Melin S.R., & Hooten M.B. (2018) A guide

to bayesian model checking for ecologists. *Ecological Monographs*, **88**, 526–542.

Gelman A., Carlin J.B., Stern H.S., Dunson D.B., Vehtari A., & Rubin D.B. (2013) *Bayesian Data Analysis*. Chapman; Hall, New York.

Huisman J., Olff H., & Fresco L. (1993) A hierarchical set of models for species response analysis. *Journal of Vegetation Science*, **4**, 37–46.

Jansen F. & Oksanen J. (2013) How to model species responses along ecological gradients—huisman—olff—fresco models revisited. *Journal of Vegetation Science*, **24**, 1108–1117.

Michaelis J. & Diekmann M.R. (2017) Biased niches – Species response curves and niche attributes from Huisman-Olff-Fresco models change with differing species prevalence and frequency. *PLOS ONE*, **12**, e0183152.

Oksanen J. & Minchin P.R. (2002) Continuum theory revisited: What shape are species responses along ecological gradients? *Ecological Modelling*, **157**, 119–129.

Urli M., Delzon S., Eyermann A., Couallier V., García-Valdés R., Zavala M.A., & Porté A.J. (2014) Inferring shifts in tree species distribution using asymmetric distribution curves: A case study in the Iberian mountains. *Journal of Vegetation Science*, **25**, 147–159.

Wilson R.J., Gutiérrez D., Gutiérrez J., Martínez D., Agudo R., & Monserrat V.J. (2005) Changes to the elevational limits and extent of species ranges associated with climate change. *Ecology letters*, **8**, 1138–1146.

Zipkin E.F., DeWan A., & Royle J.A. (2009) Impacts of forest fragmentation on species richness: A hierarchical approach to community modelling. *Journal of Applied Ecology*, **46**, 815–822.

4.3 Conclusions

In this study, we developed a new Bayesian hierarchical model, the Explicit Hierarchical Model of Optimum Shifts (EHMOS). Using simulations, we showed that EHMOS was more accurate than two widely used methods to estimate optimum shifts (i.e., the mean comparison method and GLM-based approach). In addition, EHMOS only allowed accurate optimum shift estimates for species having a partially observed relationship with the gradient. EHMOS estimates of optimum shifts of 24 true Orthoptera species were consistent with what is expected under ongoing climate change, with mostly upward shifts, which further improved confidence in the accuracy of the EHMOS method. We finally discussed EHMOS potential extensions that can lead to new insights about ecological processes driving inter-specific variability in observed optimum shifts.

5 Discussion

This thesis aimed to investigate potential bias and related issues that can arise when using partially observed data to estimate species-environment relationships (SERs). We placed our research in the context of three particular studies with different types of partially observed data: 1) partially observed response data linked to imperfect detection issues, 2) partially observed covariate data related to a specific case of spatial misalignment, and 3) partially observed relationship in the particular context of optimum shift modelling.

In the first study, we fitted a multi-species occupancy model (MSOM) to Orthoptera species along an altitudinal gradient to estimate SER accounting for imperfect detection. We showed that despite a robust sampling design involving five spatial replicates, each surveyed with three sampling techniques, most species and most sampling sites had respectively detection probability and inventory completeness lower than one, indicating that the issue of imperfect detection remained even after an intensive sampling. In addition, we found that detection probability varied among species and with grass height. Both results suggest that imperfect detection might be unavoidable, even under well-designed sampling schemes (Guillera-Arroita *et al.*, 2014). It can be particularly true when studying species communities, as we can expect species to have different detectability patterns depending on their behaviour (Veech *et al.*, 2016). For instance, in our case study, some species were highly detectable with the sweep netting technique while most of them were undetectable with this method. We thus recommend designing monitoring schemes that allow the use of site occupancy models. To overcome reservations about those sampling designs due to the potential increase in sampling effort (Welsh *et al.*, 2013), we also proposed a method to investigate sampling efficiency and the corresponding effects of potential optimisation.

In the second study, we investigated the effect of area-to-point spatial misalignment (i.e., when environmental covariates are described at a coarser spatial resolution than the resolution at which they affect the species response) on fine-scale SER estimates of three models: a GLM, a spatial GLM and a Berkson measurement error model (BEM). We found that the BEM gave more accurate estimates of SERs relative to the GLM and the spatial GLM which produced flattened SER estimates under area-to-point spatial misalignment. The BEM thus appeared to be a potential solution to deal with area-to-point misalignment when estimating fine-scale SER. Moreover, GLM estimates should not be interpreted if the environmental data are only available at coarse spatial resolutions. Those insights might help to better guide management actions by encouraging and facilitating more accurate fine-scale SER estimates.

In the third study, we developed a new formulation of a Bayesian linear model that explicitly estimates optimum shifts along environmental gradients for multiple species. This model named

the Explicit Hierarchical Model of Optimum Shifts (EHMOS) proved to be more accurate in simulations than a GLMM and the mean comparison method, i.e. two methods currently in use. The EHMOS especially improved estimates of optimum shifts for edge species, i.e. species for which the environment that they can occupy had only been partially sampled. This advantage was accompanied by better accuracy under an unbalanced sampling design, relative to the mean comparison method. Finally, by explicitly estimating optimum shifts, the EHMOS can allow the investigation of ecological hypotheses (e.g., effects of ecological traits on the magnitude of shifts) and thus help to explain observed optimum shift variability among species (Lenoir & Svenning, 2015). These improvements in optimum shift modelling can lead to new insights into the effects of climate change that species are already experiencing.

5.1 Contribution

The first study concerned the imperfect detection of species presence. While bias can arise when detectability issues are not accounted for, models dealing with imperfect detection are seldom used, especially in entomological studies. Thus, we promoted the use of multi-species occupancy models (MSOMs) specifically for the study of insect communities. We provided an example code that could be easily adapted to specific entomological studies. In addition, we propose a methodology to investigate the effects of sampling optimisation scenarios on sampling efficiency. This method can offset the common belief that occupancy models require too much sampling effort and improve the dissemination of MSOMs.

The second study investigated the effects of mismatch between covariate and response scales on SER estimates. Such a spatial misalignment is known to bias SER estimates but, in practice, no method is currently in use to deal with it and spatial misalignment can be overlooked. We investigated the accuracy of two potential solutions, a spatial GLM and a Berkson error model, to estimate SERs under area-to-point misalignment. The Berkson error model gave promising results but needs further investigation for broad applications. However, we can already advise ecologists not to overlook spatial misalignment and not to interpret SER estimates produced by GLMs or spatial GLMs under spatial area-to-point misalignment.

The third study focused on modelling optimum shifts along environmental gradients, particularly for species having an observed gradient truncated (i.e., edge species). We formulated a new model, called the Explicit Hierarchical Model of Optimum Shifts (EHMOS), that improved the accuracy of estimated optimum shifts relative to two common methods (a GLM-based and a mean comparison method), especially for edge species. With this study, we improved the estimation of optimum shifts in SER that are especially studied in the very popular domain of climate change

effects on species.

Finally, all the models built during this work will be made available on GitHub, along with tutorials to facilitate their dissemination.

5.2 Perspectives

5.2.1 Making more of the Bayesian framework

I stated earlier (see Introduction) that I preferred Bayesian models over frequentist models for conceptual reasons. I am attached by priors in the Bayesian formulation. It was an anathema to always consider that we have no knowledge about the system studied, while ecological systems have been studied for centuries. Hence, incorporating past knowledge in the model through prior distributions seemed interesting.

Using more informative priors seems appealing. In many cases, the assumption of no knowledge about parameter values is inconsistent with reality. For instance, in chapter 2, we followed common practices in ecology (Banner *et al.*, 2020) and chose default, wide normal priors for GLM coefficients. However, we had some knowledge about values that the coefficients could or could not take. For instance, there was no ecological evidence supporting the possibility that a species can have a convex (i.e., U-shaped) response curve along an environmental gradient. Thus, even if we knew that the quadratic coefficient (i.e., β_2) could not be negative, our choice of priors said otherwise. Given evidence showing an increase in estimate precision when choosing more informative priors (Morris *et al.*, 2015), further investigations on how the choice of priors can improve (or not) our results present an interesting avenue for future research.

Investigations about the choice of priors should be made in any Bayesian analysis (i.e., prior sensitivity analysis; Banner *et al.* (2020)). More consideration about the choice of priors is required in general (Depaoli *et al.*, 2020), and in ecological modelling in particular (Lemoine, 2019). A common practice in ecology is to use wide normal distributions as default priors for location parameters (Banner *et al.*, 2020). This choice is often motivated by a search for relative objectivity, but it can induce the opposite (Lemoine, 2019; Wesner & Pomeranz, 2021). For example, in logistic regression (often used to model presence/absence data), putting a normal prior with a very large variance on the intercept in the logit scale, which seems a non-informative prior, translates to a very informative prior in the probability scale (see fig. 1 in Northrup & Gerber, 2018). That is why Lemoine (2019) advocated for the use of *weakly-informative priors* informed by ecological knowledge in parameter potential values. However, in some cases, it can be difficult to translate ecological knowledge into parameter potential values. When studying species-environment relationships, knowledge is often about optimum or ecological tolerance values.

Translating this kind of knowledge to model parameters of GLMs may not be straightforward as an optimum, for instance, depends on two GLM parameters (linear and quadratic coefficients). Ecologically formulated models, such as the one we developed in the third study, should make the use of more informative priors easier for the study of species-environment relationships. Further investigations of the benefits of using such models, with appropriate priors, on the precision of SER estimates and accuracy of prediction might be interesting.

5.2.2 Effects of species traits on SERs

In the first and third chapters, we modelled SERs for multiple species simultaneously with multi-species hierarchical models. In those models, we considered that all the observed species came from the same community and shared common attributes represented by the model hyperparameters. However, some situations can question this assumption. For instance, a community can be composed of multiple (functional) groups in which a species may have a response to the environment similar to species belonging to the same group but different from species in other groups (Pacifi *et al.*, 2014). Thus, considering a unique community level may not represent reality and lead to biased estimates (Poggiato *et al.*, 2021). In addition, species sharing ecological traits may have more similar responses to the environment than species with different characteristics (Pollock *et al.*, 2012; Valente & Betts, 2019). In such situations, we can expect better estimates if the multi-species model includes trait effects on species response (e.g., Pollock *et al.*, 2012) or a priori grouping species in functional guilds (e.g., Pacifi *et al.*, 2014). Indeed, closely related species should share more information with each other than with species not sharing the same trait or not belonging to the same guild (Pollock *et al.*, 2012; Valente & Betts, 2019). However, the advantages and limitations of such hierarchical models are not yet defined (Pacifi *et al.*, 2014; Poggiato *et al.*, 2021) and need further investigation.

Incorporating species traits or species groups in the hierarchy of the multi-species models we built also has the potential to bring new ecological knowledge (Estrada *et al.*, 2016; Yates *et al.*, 2018; Valente & Betts, 2019). For instance, including trait effects on detection in the MSOM can shed new insights on species detection patterns. To my knowledge, how species traits affect species detection has received little attention. However, we can expect that detection patterns vary with the ecological characteristics of species. For instance, in our study, we can assume that two Orthoptera species with high flying capacity might be less detectable by sweep netting than species with no wings and thus no escape. Furthermore, investigations on the effect of traits on species detection can improve functional diversity assessment (Jarzyna & Jetz, 2016; Si *et al.*, 2018; Palacio *et al.*, 2020), a popular field of research (Sutherland *et al.*, 2009).

In chapter 3, allowing the shift parameter of EHMOS to vary in function of species traits

can improve our understanding of inter-specific variability of observed optimum shifts. Strong ecological assumptions propose species traits as candidates for explaining differences in species range shifts (Kearney & Porter, 2009). However, a recent review highlighted a lack of consistency in empirical evidence about trait effects on range shifts (Beissinger & Riddell, 2021). Beissinger & Riddell (2021) proposed a series of potential explanations, but they neglected to mention statistical uncertainty in shift estimates that could blur the signal when testing for trait effects in *post-hoc* analysis (e.g., Felde *et al.*, 2012). The EHMOS can overcome this issue by allowing direct tests for trait effects on optimum shifts and thus account for the uncertainty of the estimates.

5.2.3 Room for more complexity

The hierarchical Bayesian models developed in this work assumed symmetrical unimodal species responses to environmental gradients. While this facilitates ecological interpretation of model estimates and can improve estimates in case of a partially observed gradient (see Chapter 3), observed SER can be asymmetric (i.e., skewed). For instance, a competitor can exclude the species from one part of its suitability range, leading to skewed observed SER (e.g., see Figure 1 in Poggiato *et al.*, 2021). One solution is to use more flexible models, such as GAM (Pedersen *et al.*, 2019) or HOF models (Oksanen & Minchin, 2002). However, more flexible models can better describe the observed SER but not explain the underlying ecological processes. For instance, such models can not distinguish a true skewed SER, due to particular physiological constraints (e.g., plants' response to temperature can be asymmetric Austin, 2007), from an observed skewed SER that results from the effect of an external factor than the environmental covariate studied (e.g., the presence of competitors or interactions between environmental factors). Failing to understand the ecological processes that drive the observed response may lower the transferability (i.e., the ability to predict in new environments) of such models (Bell & Schlaepfer, 2016; Yates *et al.*, 2018). In addition, concerns have arisen about the sensitivity of such flexible models to data quality (Michaelis & Diekmann, 2017) and sampling bias (Merow *et al.*, 2014), especially in the presence of truncated gradients (Bell & Schlaepfer, 2016; Citores *et al.*, 2020). Thus, further investigation is needed to distinguish estimated skewed responses due to data specificity from actual skewed responses before interpreting SER estimated by flexible models. Another solution is to increase the complexity of the hierarchical structure to describe ecological processes leading to skewed responses (e.g., including species interactions; Poggiato *et al.* (2021)). HBM can be extended to allow for such external effects (e.g., joint SDM Warton *et al.*, 2015). These HBMs have received much attention regarding their predictive ability in SDM (Poggiato *et al.*, 2021). However, to my knowledge, little is known about their ability to estimate SER accurately. Further investigation into how external factors (e.g., biotic interactions) can shape observed SER is crucial

to explain SER better and improve SDM predictions.

Bibliography

- Austin, M.P. (2002) Spatial prediction of species distribution: An interface between ecological theory and statistical modelling. *Ecological Modelling*, **157**, 101–118.
- Austin, M.P. (2007) Species distribution models and ecological theory: A critical assessment and some possible new approaches. *Ecological Modelling*, **200**, 1–19.
- Bailey, L.L., Hines, J.E., Nichols, J.D. & MacKenzie, D.I. (2007) Sampling design trade-offs in occupancy studies with imperfect detection: Examples and software. *Ecological Applications*, **17**, 281–290.
- Banks-Leite, C., Pardini, R., Boscolo, D., Cassano, C.R., Püttker, T., Barros, C.S., *et al.* (2014) Assessing the utility of statistical adjustments for imperfect detection in tropical conservation science. *Journal of Applied Ecology*, **51**, 849–859.
- Banner, K.M., Irvine, K.M. & Rodhouse, T.J. (2020) The use of Bayesian priors in Ecology: The good, the bad and the not great. *Methods in Ecology and Evolution*, **11**, 882–889.
- Beissinger, S.R. & Riddell, E.A. (2021) Why Are Species' Traits Weak Predictors of Range Shifts? *Annual Review of Ecology, Evolution, and Systematics*, **52**, null.
- 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.
- Bellard, C., Bertelsmeier, C., Leadley, P., Thuiller, W. & Courchamp, F. (2012) Impacts of climate change on the future of biodiversity. *Ecology Letters*, **15**, 365–377.
- Bolker, B.M., Brooks, M.E., Clark, C.J., Geange, S.W., Poulsen, J.R., Stevens, M.H.H., *et al.* (2009) Generalized linear mixed models: A practical guide for ecology and evolution. *Trends in Ecology & Evolution*, **24**, 127–135.
- Chandler, R. & Hepinstall-Cymerman, J. (2016) Estimating the spatial scales of landscape effects on abundance. *Landscape Ecology*, **31**, 1383–1394.
- Chauvier, Y., Descombes, P., Guéguen, M., Boulangeat, L., Thuiller, W. & Zimmermann, N.E. (2022) Resolution in species distribution models shapes spatial patterns of plant multifaceted diversity. *Ecography*, **n/a**, e05973.
- Chen, I.-C., Shiu, H.-J., Benedick, S., Holloway, J.D., Chey, V.K., Barlow, H.S., *et al.* (2009) Elevation increases in moth assemblages over 42 years on a tropical mountain. *Proceedings of the National Academy of Sciences*, **106**, 1479–1483.
- Chevalier, M., Broennimann, O., Cornuault, J. & Guisan, A. (2021) Data integration methods to account for spatial niche truncation effects in regional projections of species distribution. *Ecological Applications*, **31**, e02427.
- Citores, L., Ibaibarriaga, L., Lee, D.-J., Brewer, M.J., Santos, M. & Chust, G. (2020) Modelling species presence–absence in the ecological niche theory framework using shape-constrained

- generalized additive models. *Ecological Modelling*, **418**, 108926.
- Clavel, J., Julliard, R. & Devictor, V. (2011) Worldwide decline of specialist species: Toward a global functional homogenization? *Frontiers in Ecology and the Environment*, **9**, 222–228.
- Connor, T., Hull, V., Vina, A., Shortridge, A., Tang, Y., Zhang, J., *et al.* (2018) Effects of grain size and niche breadth on species distribution modeling. *Ecography*, **41**, 1270–1282.
- Coudun, C. & Gégout, J.-C. (2005) Ecological behaviour of herbaceous forest species along a pH gradient: A comparison between oceanic and semicontinental regions in northern France. *Global Ecology and Biogeography*, **14**, 263–270.
- 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, Predicting Species Distributions*, **199**, 164–175.
- Cressie, N., Calder, C.A., Clark, J.S., Hoef, J.M.V. & Wikle, C.K. (2009) Accounting for uncertainty in ecological analysis: The strengths and limitations of hierarchical statistical modeling. *Ecological Applications*, **19**, 553–570.
- De Knecht, H., Langevelde, F. van van, Coughenour, M., Skidmore, A., De Boer, W., Heitkönig, I., *et al.* (2010) Spatial autocorrelation and the scaling of species–environment relationships. *Ecology*, **91**, 2455–2465.
- Depaoli, S., Winter, S.D. & Visser, M. (2020) The Importance of Prior Sensitivity Analysis in Bayesian Statistics: Demonstrations Using an Interactive Shiny App. *Frontiers in Psychology*, **11**, 3271.
- Devarajan, K., Morelli, T.L. & Tenan, S. (2020) Multi-species occupancy models: Review, roadmap, and recommendations. *Ecography*, **43**, 1612–1624.
- Devictor, V., Julliard, R. & Jiguet, F. (2008) Distribution of specialist and generalist species along spatial gradients of habitat disturbance and fragmentation. *Oikos*, **117**, 507–514.
- Dorazio, R.M. (2016) Bayesian data analysis in population ecology: Motivations, methods, and benefits. *Population ecology*, **58**, 31–44.
- Dormann, C.F. (2007) Promising the future? Global change projections of species distributions. *Basic and Applied Ecology*, **8**, 387–397.
- Elith, J., Kearney, M. & Phillips, S. (2010) The art of modelling range-shifting species. *Methods in Ecology and Evolution*, **1**, 330–342.
- Ellison, A.M. (2004) Bayesian inference in ecology. *Ecology Letters*, **7**, 509–520.
- Estrada, A., Morales-Castilla, I., Caplat, P. & Early, R. (2016) Usefulness of Species Traits in Predicting Range Shifts. *Trends in Ecology & Evolution*, **31**, 190–203.
- Faurby, S. & Araújo, M.B. (2018) Anthropogenic range contractions bias species climate change forecast. *Nature Climate Change*, **8**, 252–256.

- Felde, V.A., Kapfer, J. & Grytnes, J.-A. (2012) Upward shift in elevational plant species ranges in Sikkilsdalen, central Norway. *Ecography*, **35**, 922–932.
- Foster, S.D., Shimadzu, H. & Darnell, R. (2012) Uncertainty in spatially predicted covariates: Is it ignorable? *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, **61**, 637–652.
- Freeman, B.G., Scholer, M.N., Ruiz-Gutierrez, V. & Fitzpatrick, J.W. (2018) Climate change causes upslope shifts and mountaintop extirpations in a tropical bird community. *Proceedings of the National Academy of Sciences*, **115**, 11982–11987.
- Gotway, C.A. & Young, L.J. (2002) Combining Incompatible Spatial Data. *Journal of the American Statistical Association*, 632–648.
- Guillera-Arroita, G. (2017) Modelling of species distributions, range dynamics and communities under imperfect detection: Advances, challenges and opportunities. *Ecography*, **40**, 281–295.
- Guillera-Arroita, G., Lahoz-Monfort, J.J., Elith, J., Gordon, A., Kujala, H., Lentini, P.E., *et al.* (2015) Is my species distribution model fit for purpose? Matching data and models to applications. *Global Ecology and Biogeography*, **24**, 276–292.
- Guillera-Arroita, G., Lahoz-Monfort, J.J., MacKenzie, D.I., Wintle, B.A. & McCarthy, M.A. (2014) Ignoring Imperfect Detection in Biological Surveys Is Dangerous: A Response to ‘Fitting and Interpreting Occupancy Models’. *PLOS ONE*, **9**, e99571.
- Guillera-Arroita, G., Ridout, M.S. & Morgan, B.J.T. (2010) Design of occupancy studies with imperfect detection: Design of occupancy studies. *Methods in Ecology and Evolution*, **1**, 131–139.
- Guisan, A., Edwards, T.C. & Hastie, T. (2002) Generalized linear and generalized additive models in studies of species distributions: Setting the scene. *Ecological Modelling*, **157**, 89–100.
- Guisan, A., Graham, C.H., Elith, J., Huettmann, F. & Group, the N.S.D.M. (2007) Sensitivity of predictive species distribution models to change in grain size. *Diversity and Distributions*, **13**, 332–340.
- Guisan, A., Tingley, R., Baumgartner, J.B., Naujokaitis-Lewis, I., Sutcliffe, P.R., Tulloch, A.I.T., *et al.* (2013) Predicting species distributions for conservation decisions. *Ecology Letters*, **16**, 1424–1435.
- Heegaard, E. (2002) The outer border and central border for species–environmental relationships estimated by non-parametric generalised additive models. *Ecological Modelling*, **157**, 131–139.
- Hefley, T. & B. Hooten, M. (2016) Hierarchical Species Distribution Models. *Current Landscape Ecology Reports*.
- 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.
- Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G. & Jarvis, A. (2005) Very high resolution

- interpolated climate surfaces for global land areas. *International Journal of Climatology*, **25**, 1965–1978.
- Huisman, J., Olf, H. & Fresco, L.f.m. (1993) A hierarchical set of models for species response analysis. *Journal of Vegetation Science*, **4**, 37–46.
- Hutchinson, G.E. (1957) Concluding Remarks. *Cold Spring Harbor Symposia on Quantitative Biology*, **22**, 415–427.
- Jansen, F. & Oksanen, J. (2013) How to model species responses along ecological gradients – Huisman–Olf–Fresco models revisited. *Journal of Vegetation Science*, **24**, 1108–1117.
- Jarzyna, M.A. & Jetz, W. (2016) Detecting the Multiple Facets of Biodiversity. *Trends in Ecology & Evolution*, **31**, 527–538.
- Kearney, M. & Porter, W. (2009) Mechanistic niche modelling: Combining physiological and spatial data to predict species’ ranges. *Ecology Letters*, **12**, 334–350.
- Kellner, K.F. & Swihart, R.K. (2014) Accounting for Imperfect Detection in Ecology: A Quantitative Review. *PLOS ONE*, **9**, e111436.
- Kermorvant, C., D’Amico, F., Bru, N., Caill-Milly, N. & Robertson, B. (2019) Spatially balanced sampling designs for environmental surveys. *Environmental Monitoring and Assessment*, **191**, 524.
- Kéry, M. (2010) *Introduction to WinBUGS for ecologists: Bayesian approach to regression, ANOVA, mixed models and related analyses*. Academic Press.
- Kéry, M. & Royle, J.A. (2015) *Applied Hierarchical Modeling in Ecology: Analysis of distribution, abundance and species richness in R and BUGS*. AP.
- Kéry, M. & Royle, J.A. (2021) *Applied Hierarchical Modeling in Ecology: Analysis of Distribution, Abundance and Species Richness in R and BUGS: Volume 2: Dynamic and Advanced Models*. Academic Press.
- Lahoz-Monfort, J.J., Guillera-Aroita, G. & Wintle, B.A. (2014) Imperfect detection impacts the performance of species distribution models. *Global ecology and biogeography*, **23**, 504–515.
- Latimer, A.M., Wu, S., Gelfand, A.E. & Silander, J.A. (2006) Building Statistical Models To Analyze Species Distributions. *Ecological Applications*, **16**, 33–50.
- Lele, S.R. & Dennis, B. (2009) Bayesian methods for hierarchical models: Are ecologists making a Faustian bargain. *Ecological Applications*, **19**, 581–584.
- Lembrechts, J.J., Alexander, J.M., Cavieres, L.A., Haider, S., Lenoir, J., Kueffer, C., *et al.* (2017) Mountain roads shift native and non-native plant species’ ranges. *Ecography*, **40**, 353–364.
- Lembrechts, J.J., Nijs, I. & Lenoir, J. (2019) Incorporating microclimate into species distribution models. *Ecography*, **42**, 1267–1279.
- Lemoine, N.P. (2019) Moving beyond noninformative priors: Why and how to choose weakly

- informative priors in Bayesian analyses. *Oikos*, **128**, 912–928.
- Lenoir, J., Gégout, J.C., Marquet, P.A., Ruffray, P. de & Brisse, H. (2008) A Significant Upward Shift in Plant Species Optimum Elevation During the 20th Century. *Science*, **320**, 1768–1771.
- Lenoir, J. & Svenning, J.-C. (2015) Climate-related range shifts – a global multidimensional synthesis and new research directions. *Ecography*, **38**, 15–28.
- Lunn, D.J., Thomas, A., Best, N. & Spiegelhalter, D. (2000) WinBUGS-a bayesian modelling framework: Concepts, structure, and extensibility. *Statistics and Computing*, **10**, 325–337.
- MacKenzie, D.I., Nichols, J.D., Lachman, G.B., Droege, S., Royle, J.A. & Langtimm, C.A. (2002) Estimating Site Occupancy Rates When Detection Probabilities Are Less Than One. *Ecology*, **83**, 2248–2255.
- MacKenzie, D.I., Nichols, J.D., Royle, J.A., Pollock, K.H., Bailey, L.L. & Hines, J.E. (2006) *Occupancy Estimation and Modeling : Inferring Patterns and Dynamics of Species Occurrence*. Elsevier.
- Martínez-Minaya, J., Cameletti, M., Conesa, D. & Pennino, M.G. (2018) Species distribution modeling: A statistical review with focus in spatio-temporal issues. *Stochastic Environmental Research and Risk Assessment*, **32**, 3227–3244.
- Mata, L., Goula, M. & Hahs, A.K. (2014) Conserving insect assemblages in urban landscapes: Accounting for species-specific responses and imperfect detection. *Journal of insect conservation*, **18**, 885–894.
- McElreath, R. (2016) *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. Chapman; Hall/CRC, Boca Raton.
- McInerney, G.J. & Etienne, R.S. (2012) Ditch the niche – is the niche a useful concept in ecology or species distribution modelling? *Journal of Biogeography*, **39**, 2096–2102.
- McInerney, G.J. & Purves, D.W. (2011) Fine-scale environmental variation in species distribution modelling: Regression dilution, latent variables and neighbourly advice. *Methods in Ecology and Evolution*, **2**, 248–257.
- Meineri, E. & Hylander, K. (2017) Fine-grain, large-domain climate models based on climate station and comprehensive topographic information improve microrefugia detection. *Ecography*, **40**, 1003–1013.
- Menéndez, R., González-Megías, A., Jay-Robert, P. & Marquéz-Ferrando, R. (2014) Climate change and elevational range shifts: Evidence from dung beetles in two European mountain ranges. *Global Ecology and Biogeography*, **23**, 646–657.
- Merow, C., Smith, M.J., Edwards Jr, T.C., Guisan, A., McMahon, S.M., Normand, S., *et al.* (2014) What do we gain from simplicity versus complexity in species distribution models? *Ecography*, **37**, 1267–1281.

- Michaelis, J. & Diekmann, M.R. (2017) Biased niches – Species response curves and niche attributes from Huisman-Olff-Fresco models change with differing species prevalence and frequency. *PLOS ONE*, **12**, e0183152.
- Morris, W.K., Vesk, P.A., McCarthy, M.A., Bunyavejchewin, S. & Baker, P.J. (2015) The neglected tool in the bayesian ecologist's shed: A case study testing informative priors' effect on model accuracy. *Ecology and Evolution*, **5**, 102–108.
- Mouquet, N., Lagadeuc, Y., Devictor, V., Doyen, L., Duputié, A., Eveillard, D., *et al.* (2015) REVIEW: Predictive ecology in a changing world. *Journal of Applied Ecology*, **52**, 1293–1310.
- Northrup, J.M. & Gerber, B.D. (2018) A comment on priors for Bayesian occupancy models. *PLOS ONE*, **13**, e0192819.
- Oksanen, J. & Minchin, P.R. (2002) Continuum theory revisited: What shape are species responses along ecological gradients? *Ecological Modelling*, **157**, 119–129.
- Ovaskainen, O. & Soininen, J. (2011) Making more out of sparse data: Hierarchical modeling of species communities. *Ecology*, **92**, 289–295.
- Ovaskainen, O., Tikhonov, G., Norberg, A., Blanchet, F.G., Duan, L., Dunson, D., *et al.* (2017) How to make more out of community data? A conceptual framework and its implementation as models and software. *Ecology Letters*, **20**, 561–576.
- Pacifici, K., Zipkin, E.F., Collazo, J.A., Irizarry, J.I. & DeWan, A. (2014) Guidelines for a priori grouping of species in hierarchical community models. *Ecology and Evolution*, **4**, 877–888.
- Palacio, F.X., Maragliano, R.E. & Montalti, D. (2020) The costs of ignoring species detectability on functional diversity estimation. *The Auk*, **137**, ukaa057.
- Pearson, R.G. & Dawson, T.P. (2003) Predicting the impacts of climate change on the distribution of species: Are bioclimate envelope models useful? *Global Ecology and Biogeography*, **12**, 361–371.
- Pecl, G.T., Araújo, M.B., Bell, J.D., Blanchard, J., Bonebrake, T.C., Chen, I.-C., *et al.* (2017) Biodiversity redistribution under climate change: Impacts on ecosystems and human well-being. *Science*, **355**, eaai9214.
- Pedersen, E.J., Miller, D.L., Simpson, G.L. & Ross, N. (2019) Hierarchical generalized additive models in ecology: An introduction with mgcv. *PeerJ*, **7**, e6876.
- Plummer, M. (2003) JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. *Working Papers*, **8**.
- Poggiato, G., Münkemüller, T., Bystrova, D., Arbel, J., Clark, J.S. & Thuiller, W. (2021) On the Interpretations of Joint Modeling in Community Ecology. *Trends in Ecology & Evolution*, **36**, 391–401.
- Pollock, L.J., Morris, W.K. & Vesk, P.A. (2012) The role of functional traits in species distributions revealed through a hierarchical model. *Ecography*, **35**, 716–725.

- Potter, K.A., Woods, A.H. & Pincebourde, S. (2013) Microclimatic challenges in global change biology. *Global Change Biology*, **19**, 2932–2939.
- Royle, J.A. (2004) N-Mixture Models for Estimating Population Size from Spatially Replicated Counts. *Biometrics*, **60**, 108–115.
- Royle, J.A. & Dorazio, R.M. (2008) Hierarchical modeling and inference in ecology: The analysis of data from populations, metapopulations and communities.
- Rumpf, S.B., Hülber, K., Klöner, G., Moser, D., Schütz, M., Wessely, J., *et al.* (2018) Range dynamics of mountain plants decrease with elevation. *Proceedings of the National Academy of Sciences*, **115**, 1848–1853.
- Russell, R.E., Royle, J.A., Saab, V.A., Lehmkühl, J.F., Block, W.M. & Sauer, J.R. (2009) Modeling the effects of environmental disturbance on wildlife communities: Avian responses to prescribed fire. *Ecological Applications*, **19**, 1253–1263.
- Shoo, L.P., Williams, S.E. & Hero, J.-M. (2006) Detecting climate change induced range shifts: Where and how should we be looking? *Austral Ecology*, **31**, 22–29.
- Si, X., Cadotte, M.W., Zhao, Y., Zhou, H., Zeng, D., Li, J., *et al.* (2018) The importance of accounting for imperfect detection when estimating functional and phylogenetic community structure. *Ecology*, **99**, 2103–2112.
- Sutherland, W.J., Adams, W.M., Aronson, R.B., Aveling, R., Blackburn, T.M., Broad, S., *et al.* (2009) One Hundred Questions of Importance to the Conservation of Global Biological Diversity. *Conservation Biology*, **23**, 557–567.
- ter Braak, C.J.F. & Looman, C.W.N. (1986) Weighted averaging, logistic regression and the Gaussian response model. *Vegetatio*, **65**, 3–11.
- Thuiller, W., Brotons, L., Araújo, M.B. & Lavorel, S. (2004) Effects of restricting environmental range of data to project current and future species distributions. *Ecography*, **27**, 165–172.
- Tingley, M.W., Nadeau, C.P. & Sandor, M.E. (2020) Multi-species occupancy models as robust estimators of community richness. *Methods in Ecology and Evolution*, **11**, 633–642.
- Tyre, A.J., Tenhumberg, B., Field, S.A., Niejalke, D., Parris, K. & Possingham, H.P. (2003) Improving Precision and Reducing Bias in Biological Surveys: Estimating False-Negative Error Rates. *Ecological Applications*, **13**, 1790–1801.
- Urli, M., Delzon, S., Eyermann, A., Couallier, V., García-Valdés, R., Zavala, M.A., *et al.* (2014) Inferring shifts in tree species distribution using asymmetric distribution curves: A case study in the Iberian mountains. *Journal of Vegetation Science*, **25**, 147–159.
- Valente, J.J. & Betts, M.G. (2019) Response to fragmentation by avian communities is mediated by species traits. *Diversity and Distributions*, **25**, 48–60.
- Veech, J.A., Ott, J.R. & Troy, J.R. (2016) Intrinsic heterogeneity in detection probability and its

- effect on N-mixture models. *Methods in Ecology and Evolution*, **7**, 1019–1028.
- Veen, B. van der, Hui, F.K.C., Hovstad, K.A., Solbu, E.B. & O'Hara, R.B. (2021) Model-based ordination for species with unequal niche widths. *Methods in Ecology and Evolution*, **n/a**.
- Warton, D.I., Blanchet, F.G., O'Hara, R.B., Ovaskainen, O., Taskinen, S., Walker, S.C., *et al.* (2015) So Many Variables: Joint Modeling in Community Ecology. *Trends in Ecology & Evolution*, **30**, 766–779.
- Watson, J.E.M., Rao, M., Ai-Li, K. & Yan, X. (2012) Climate Change Adaptation Planning for Biodiversity Conservation: A Review. *Advances in Climate Change Research*, **3**, 1–11.
- Welsh, A.H., Lindenmayer, D.B. & Donnelly, C.F. (2013) Fitting and Interpreting Occupancy Models. *PLOS ONE*, **8**, e52015.
- Wesner, J.S. & Pomeranz, J.P.F. (2021) Choosing priors in Bayesian ecological models by simulating from the prior predictive distribution. *Ecosphere*, **12**, e03739.
- Wikle, C.K. (2003) Hierarchical Bayesian Models for Predicting the Spread of Ecological Processes. *Ecology*, **84**, 1382–1394.
- Yates, K.L., Bouchet, P.J., Caley, M.J., Mengersen, K., Randin, C.F., Parnell, S., *et al.* (2018) Outstanding Challenges in the Transferability of Ecological Models. *Trends in Ecology & Evolution*, **33**, 790–802.
- Yoccoz, N.G., Nichols, J.D. & Boulinier, T. (2001) Monitoring of biological diversity in space and time. *Trends in Ecology & Evolution*, **16**, 446–453.
- Zellweger, F., De Frenne, P., Lenoir, J., Vangansbeke, P., Verheyen, K., Bernhardt-Römermann, M., *et al.* (2020) Forest microclimate dynamics drive plant responses to warming. *Science*, **368**, 772–775.
- Zurell, D., König, C., Malchow, A.-K., Kapitza, S., Bocedi, G., Travis, J., *et al.* (2021) Spatially explicit models for decision-making in animal conservation and restoration. *Ecography*, **44**, 1–16.