

Model linearity breeds contempt: using Bayesian non-linear models to uncover broad macroecological patterns

Bernat Bramon Mora^{1,*}, Antoine Guisan^{2,3} and Jake M. Alexander¹

¹Institute of Integrative Biology, ETH Zürich, Zürich, Switzerland; ²Department of Ecology and Evolution, University of Lausanne, Lausanne, Switzerland; ²Department of Ecology and Evolution, University of Lausanne, Lausanne, Switzerland; ³Institute of Earth Surface Dynamics, University of Lausanne, Lausanne, Switzerland; *bernat.bramon@gmail.com

1 Abstract

2 Species' realized niches are classically pictured as bell-shaped probability distributions. These
3 distributions, however, can actually take many different forms. For example, fat-tailed or
4 skewed responses are very common across fields, as these can naturally emerge as a result of
5 several ecological processes. While one does not need to know the shape of species' distribu-
6 tions to effectively model them, studying their basic form can teach us a lot about the ways
7 climatic processes and historical contingencies have shaped ecological communities. Unfor-
8 tunately, we still lack a general understanding of the basic properties describing the shape of
9 species' distributions, and much less is known about how these compare to each other across
10 gradients. Here, we use a set of Bayesian non-linear models to uncover such properties.
11 These models account for all prior knowledge we have regarding species' realized niches, in-
12 cluding expert knowledge of their environmental preferences and ecological strategies. With
13 this approach, we are able to distil the shape of empirical plant distributions, which helps
14 us tackle long-standing hypotheses regarding the way ecological communities are assembled
15 across space. In particular, we study the relationship between several properties of distri-
16 butions, such as the link between species' range size and elevation, revealing the existence
17 of broad macroecological patterns along environmental gradients. Moreover, we are able to
18 shed light on the extent to which some aspects of the shape of observed realized niches—
19 such as kurtosis and skewness of the distributions—could be intrinsic properties of species'

historical contexts. Overall, our approach offers a useful statistical framework to understand the shape of species' distributions, and our results provide an unprecedented perspective of the way systems of many species are distributed along environmental gradients.

Introduction

One of the central goals of ecology is to understand the ways species are distributed across space and time. While many ecological textbooks assume the shape of species' realized niches to be unimodal and symmetric along environmental gradients (Krebs, 1972), some have warned that empirical distributions can take many different forms (Austin, 1987, 2002; Sagarin *et al.*, 2006). In practice, there is a strong argument to be made in favour of assuming these to be bell shaped. Namely, if all that we are willing to assume about species' distributions is that these occupy finite geographic ranges, the most conservative statistical approach is to model their distribution as Gaussian (i.e. the corresponding maximum entropy distribution; Frank 2009). That said, there is currently no general agreement on the basic shape of species' realized niches. Indeed, many factors can play a role in defining their shape, and several natural processes can lead to non-normal distributions.

Fat-tailed and skewed distributions are very common across fields. The former naturally emerges as a result of processes involving seasonality (e.g. in communications patterns; Malmgren *et al.* 2008) or some stochastic events (e.g. in the spread of infectious diseases; Wong & Collins 2020). Indeed, species' dispersal patterns have been shown to have fat tails due to the natural variability among individuals (Petrovskii *et al.*, 2009). This is important because one might expect environmental and individual variation to also be crucial factors determining the presence and absence of species along gradients, and fat tails are therefore a plausible property of species' realized niches. Similarly, several processes can lead to skewed distributions. For example, species might present asymmetric environmental tolerances along altitudinal gradients, allowing them to withstand different temperature extremes (Sunday *et al.*, 2011). Species might also experience abiotic and biotic pressures that increase or decrease along a temperature gradient, which could result in species' distributions presenting steeper declines towards warmer or colder environments (Normand *et al.*, 2009). Overall, many different properties could characterize species' realized niches, and every new shape

entails different underlying hypotheses regarding the way communities are assembled over time (D'Amen *et al.*, 2017).

Comparing these properties across species allows us to study broad macroecological patterns that could be critical from a conservation and management perspective (Stevens, 1992; Channell & Lomolino, 2000a). For example, the Rapoport's rule, a classic biogeographical hypothesis, predicts species' ranges to increase with latitude or elevation (Stevens, 1992), hinting at the existence of general biogeographical constraints that shape species' distributions along gradients. This sort of macroecological patterns are interesting because they provide insights into the way different species assemble and establish in different environments (Linder *et al.*, 2000). That is, the differences in species' responses to the environment can shed light on how climatic processes and historical contingencies have differently shaped their distributions (Rohde, 1992; Helmuth *et al.*, 2004; Siefert *et al.*, 2015). Uncovering the shape of species' realized niches and the extent to which these vary across species is nevertheless a challenging statistical problem to solve. Indeed, to this date, we do not have an effective way to parsimoniously compare the shape of the realized niches of many species along environmental gradients.

Over the last two decades, ecologists have developed a plethora of distribution models to try to untangle the factors that play a role in defining species' realized niches (Guisan & Zimmermann, 2000). These models are fundamental to the scientific community for predicting changes in species' geographic distributions and the effects of environmental disturbances. Such frameworks, however, commonly assume an underlying linear relationship between covariates (but see 'semiparametric models'; Norberg *et al.* 2019). This is useful because it simplifies the optimization process, but it might not be ideal when studying and comparing the shape of species' distributions along environmental gradients. First and foremost, a linear relationship between covariates often comes with a set of implicit mathematical constraints that might not be biologically justified. While this might not hinder the predictive performance of the models (Norberg *et al.*, 2019), a direct biological interpretation of parameter estimates in linear models becomes increasingly difficult as one moves from unimodal and symmetric distributions (ter Braak & Looman, 1986; Jamil & ter Braak, 2013) to fat-tailed or skewed responses (Huisman *et al.*, 1993). Second, the aforementioned structural constraints also limit our ability to include any prior information to our parameter esti-

80 mates. Observations on species' geographic variation and optimal climatic conditions have
81 long been documented, with extensive databases compiled by botanists and field ecologists
82 documenting basic knowledge of species' realized niches (e.g. Landolt *et al.* 2010). That
83 said, this information is rarely accounted for in most modelling approaches, likely because
84 there is not a straightforward way to feed this information into the parameters of a linear
85 model (Scherrer & Guisan 2019; but see ter Braak & Looman 1986; Ovaskainen *et al.* 2017).
86 Finally, some have proposed several non-linear structures to characterize several features of
87 individual species' response curves (Huisman *et al.*, 1993). Setting aside the fact that the
88 interpretation and comparison of parameter estimates becomes challenging following most
89 of these model structures, these are generally not designed to jointly study different species,
90 taking full advantage of modern statistical approaches (e.g. sharing information among
91 species or accounting for parameter uncertainty; Evans *et al.* 2016).

92 In this work, we rethink traditional modelling approaches and develop a conceptually
93 simple—and yet statistical and computationally complex—statistical framework to revisit
94 some classic hypothesis in ecology and biogeography. In particular, we develop a Bayesian
95 hierarchical model that accounts for all prior information that we have regarding the dis-
96 tribution of plant species along an elevation gradient in the Swiss Alps, including expert
97 knowledge of species environmental indicator values, range sizes, and plant ecological strate-
98 gies. We start by considering species' response curves as Gaussian distributed, and then we
99 adapt our model to allow non-linear responses characterizing skewed and long-tailed distri-
100 butions. Using this statistical framework, we are able to compare the basic properties of
101 the realized niche of multiple species, testing for the existence of broad macroecological pat-
102 terns. Comparing the posterior distribution of those parameters that control for the shape
103 of distributions, we are also able to showcase variation in the way different types of species,
104 such as native or neophytes, might respond to the environment. More generally, we are able
105 to uncover the approximate shape of empirical plant distributions and answer fundamental
106 questions regarding the way systems of many species are distributed along environmental
107 gradients.

Methods

Empirical data

We studied the distribution of plant communities along an elevation gradient. To do so, we combined two different datasets: i) one describing the co-occurrence of species across multiple open grasslands in the Swiss Alps (Randin *et al.*, 2009), and ii) an extensive floristic database containing environmental and physiological traits for all vegetation across Switzerland (Landolt *et al.*, 2010).

Distribution data

We used data describing the distribution of 798 species across 912 sites covering most of the mountain region of the Western Alps in the Canton de Vaud (Switzerland; Scherrer & Guisan 2019). Each of these sites is a 8×8 m plot placed somewhere along an elevation range from 375 m to 3210 m. In all sites, presence/absence data as well as Braun-Blanquet abundance-dominance classes were recorded for all species. Additionally, we used meteorological data provided by Scherrer & Guisan (2019), containing multiple variables characterizing the climate in each site at high spatial resolution (25 m). This dataset was compiled based on 30 years (1961–1990) of records from national weather stations. Since most of the data is highly correlated, we calculated the main axes of variation of the following scaled variables: daily minimum, maximum and average temperature; sum of growing degree-days above 5°C ; mean temperature of wettest quarter; annual precipitation, precipitation seasonality, and precipitation of driest quarter (Supplementary Fig. 1).

Floristic data

To complement the aforementioned distribution data, we used a floristic database of around 5500 vascular plants across Switzerland. Some of the information in this database has been previously shown to account for unexplained variation when used as explanatory variables in species' distribution models (Scherrer & Guisan, 2019). It was built based on expert knowledge and phytosociological field experience of botanists and ecologists, and contains

information regarding plants’ environmental preferences and ecological strategies.

Species’ environmental preferences in this database can be used to inform distribution models—e.g. as an informative prior in a Bayesian framework. These are characterized following the ecological indicator values developed by Landolt *et al.* (2010), providing both an estimate of the average conditions in which a species can be found as well as a broad description of their range of variation. These values are provided for a range of 8 environmental variables, including temperature, continentality, light conditions, as well as moisture, acidity and nutrient content of the soil (see a full list and description of the ecological indicators in the Supplementary Table 1; Landolt *et al.* 2010). In addition to species’ environmental preferences, the floristic data also contains information on species introduction status (e.g. identifying those species that are recent and historical range expanders) and change tendency (e.g. indicating species that have shown decline or increase in their populations over the recent decades). We describe this information in more detail in the Supplementary Table 1.

Baseline model

There is a long list of model structures well suited to characterizing species’ distributions (see Norberg *et al.* 2019). As a baseline model, however, we were interested in a hierarchical model that does not make any assumptions regarding the shape of the distributions, and yet explicitly incorporates all information that we have regarding plant’s environmental preferences. More specifically, we wanted to account for the climatic indicator values and range of variation registered in the floristic database for all plants in our dataset. These two values provide basic information regarding plant’s optimal environmental conditions and width of their distributions.

Response curve

To choose an appropriate response curve, we first need to agree on what we truly know about the system. Given the prior information that we have about the system, we know that species occupy specific geographic ranges; therefore, we know that their distributions have finite variance. While we could also assume that many other factors might influence

species' presence in a given site—e.g. the biotic interactions among species in the site—we do not necessarily have an *a priori* expectation of how exactly these factors will influence the shape of species' distributions. Therefore, for this baseline model, if all that we are willing to assume about species' realized niches is that these have finite variance, the most conservative assumption and the safest bet—i.e. the one with the largest entropy—is that they follow a Gaussian distribution (Fig 1a). That is, given the presence/absence or abundance y_{ij} of any species i in any given site j , and an environmental variable x_j , we can define species' responses to the environment as

$$y_{ij} \sim F(p_{ij})$$

$$\log(p_{ij}) = -\alpha_i - \gamma_i(x_j - \beta_i)^2, \quad (1)$$

where F is the likelihood function, and α_i , β_i^k , and γ_i describe amplitude of the probability p_{ij} , species' average climatic suitability and range of variation along the environmental gradient, respectively. Notice that F characterizes a Binomial distribution when considering binary data, and it characterizes an ordered categorical likelihood function when we consider Braun-Blanquet abundance-dominance classes as response variables (see the full description of both models in the Supplementary Methods). For the sake of simplicity, we use only one environmental variable to characterize the species' probability distribution. That said, this model can easily be generalized to account for multiple predictors (see Supplementary Methods).

Model priors

The model structure described above allows us to explicitly incorporate all prior knowledge that we have regarding species' distributions contained in the floristic database. To do so,

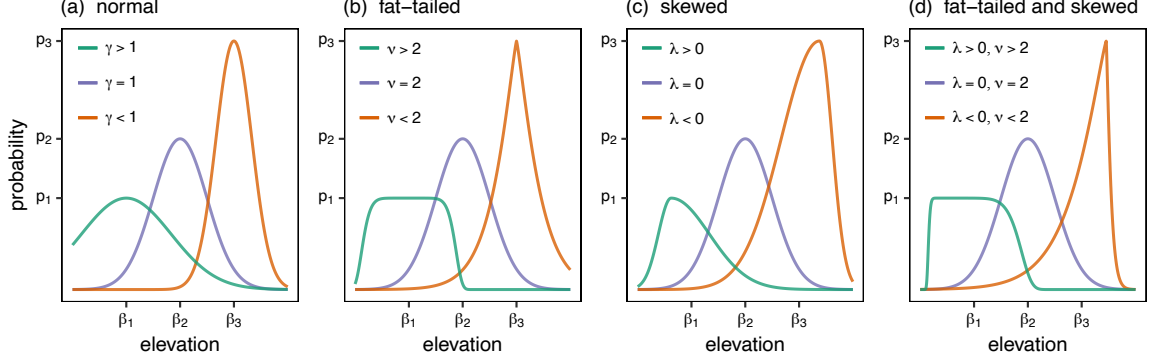


Figure 1: Different response curves. Panel (a) shows the distribution shapes characterized by Eq. (1) for different values of β , α and γ . Panel (b) shows the distribution shapes characterized by Eq. (4) for different values of β , α and ν , when $\gamma = 1$. Panel (c) shows the distribution shapes characterized by Eq. (5) for different values of β , α and λ , when $\gamma = 1$. Panel (d) shows the distribution shapes characterized by Eq. (6) for different values of β , α , λ and ν , when $\gamma = 1$. Notice that for all panels, we chose α values such that $p_i = \exp(-\alpha_i)$.

we define the prior distributions for the parameters in model (1) as:

$$\begin{aligned}
 \beta_i &\sim \text{MVNormal}(\hat{\beta}, \Sigma^\beta) \\
 \log(\gamma_i) &\sim \text{MVNormal}(\hat{\gamma}, \Sigma^\gamma) \\
 \log(\alpha_i) &\sim \text{Normal}(\hat{\alpha}, \sigma_\alpha) \\
 \hat{\beta}, \hat{\gamma}, \hat{\alpha} &\sim \text{Normal}(0, 1) \\
 \sigma_\alpha &\sim \text{Exponential}(1)
 \end{aligned} \tag{2}$$

where parameters γ_i and β_i are expressed as multivariate normal distributions—i.e. Gaussian processes—such that Σ^β and Σ^γ are variance-covariance matrices describing species' similarity in terms of their average climatic suitability and range of variation along the different environmental gradients, respectively. We define these variance-covariance matrices as follows:

$$\Sigma_{ij} = \eta \exp(-\rho D_{ij}^2) + \delta_{ij} \sigma, \tag{3}$$

where Σ_{ij} characterizes the covariance between any pair of species i and j , and δ_{ij} is the Kronecker delta. Such a covariance structure declines exponentially with the square of a distance matrix D_{ij} , which characterize differences between species computed using our prior information. In the floristic database, this information is represented by the set of ordinal traits specified for the different species. While there are many different ways to turn

ordinal data into distance matrices, we choose to use a mixed-membership stochastic block model because it allows us to deal with cases of missing data (see Supplementary Methods for extended details; Godoy-Lorite *et al.* 2016). In each covariance matrix, the hyperparameter ρ determines the rate of decline of the covariance between any two species, and η defines its maximum value. The hyperparameter σ describes the additional covariance between the different observations for any given species. For all these hyperparameters, we choose weakly informative priors such that $\sigma, \eta \sim \text{Exponential}(1)$ and $\rho \sim \text{Exponential}(0.5)$. Notice that other structures can be used to define the covariance matrices of the different Gaussian processes (McElreath, 2020), including structures that account for multiple distance matrix D_{ij} for any given parameter.

Sampling the posterior

We generated the posterior samples for the Bayesian models with the Hamiltonian Monte Carlo algorithm implementation provided by the R packages ‘rstan’ and ‘cmdstanr’ (Stan Development Team, 2021). Sampling models like the ones described above can be computationally very intensive. This is especially true when using ordered categorical likelihood functions (see Stan Development Team 2021). Therefore, we focus on those species for which we have at least 20 occurrences when modelling both binary data and ordinal data.

To test the performance of the model as well as our choice of prior distributions, we modelled simulated data and compared the sampled posterior distributions to the data-generating parameters (e.g. Supplementary Fig. 2; see Code Availability section). Notice that using the link function in Eq. (1) could cause problems when sampling the model, and some adjustments need to be made when specifying the model (see Code Availability section). To perform the data analysis and generate the figures, we used some of the functions available with the R package ‘rethinking’ (McElreath, 2020).

Modifying the baseline model

We proposed a baseline model that is naive regarding how the data is distributed, and yet accounts for all prior information that we have about the system. Now, we want to modify this model to test the extent to which empirical species’ distributions showcase different

shapes. We focused on two properties: fat-tailed and skewed responses. While there are several model structures that could account for these properties, we propose new species' response curves following three criteria. First, the probability distribution of a species along an environmental gradient must have a defined mean and variance. This is important because we know that species naturally have different environmental preferences as well as finite geographic ranges. Second, the Gaussian shape must be a special case of the probability distribution, allowing species to showcase variation regarding the presence (or lack thereof) of any given pattern. Finally, there must be a re-parametrization of the model that allows us to keep the same prior information and interpretable parameters.

Fat-tailed response curve

Fat-tailed distributions represent distributions with relatively high representation of extreme events. While many different distributions exhibit this property, we decided to accommodate this feature into our baseline model by considering a response curve that follows a generalized error distribution. Such a distribution is useful because the Gaussian shape is a special case of it, and it contains a parameter that regulates the level of kurtosis—ranging from longer to shorter tails than the Gaussian case (Fig 1b). In particular, we can adapt Eq. (1) to present this non-linear form as follows:

$$\log(p_{ij}) = -\alpha_i - \gamma'_i |x_j - \beta_i|^{\nu_i}, \quad (4)$$

where $\gamma'_i = g(\gamma_i, \nu_i)$, and ν_i is a parameter that describes the kurtosis of the distribution, which we define as $\nu_i \in (1, \infty)$. Following this, we choose an adaptive prior for this set of new parameter such that $\log(\nu_i - 1) \sim \text{Normal}(\hat{\nu}, \sigma_\nu)$, where $\hat{\nu} \sim \text{Normal}(0, 1)$ and $\sigma_\nu \sim \text{Exponential}(2)$. Given the relationship between γ'_i and γ_i , we can re-parametrize the model and follow Eq. (2) to define the prior distributions (see Supplementary Table 2; Nadarajah 2005). Notice that the Gaussian distribution will naturally emerge when $\nu_i = 2$.

Alternatively, we could have used other distributions that present fat tails and fulfil the selection criteria described above. For example, the non-standardized Student's t-distributions is an interesting distribution because, as opposed to the generalized error distribution, it allows for fat tails without generating a cusp at the center (see Fig 1b). However, we avoided

using the non-standardized Student’s t-distributions because it does not allow for tails that are lighter than normal (e.g. $\nu_i > 2$ in Eq. 4; Fig 1b), and the sampling of the model can be somewhat more challenging (ref).

Skewed response curve

Skewed responses present steeper declines towards either side of the distribution. One way to accommodate this feature in our models is by considering a skewed normal distribution. As for the case described above, the Gaussian is a special case of this distribution, and it contains a parameter that controls for the level and direction of ‘skewness’ (Fig 1c). Importantly, this distribution presents normal-like tails; therefore, the added skewness does not make additional assumptions regarding how species are distributed along the gradient. To test for the existence of this feature, we modified Eq. (1) as

$$\log(p_{ij}) = -\alpha_i - \gamma'_i \left(\frac{x_j - \beta'_i}{1 + \lambda_i \operatorname{sgn}(x_j - \beta'_i)} \right)^2, \quad (5)$$

where $\gamma'_i = q_1(\gamma_i, \lambda_i)$, $\beta'_i = q_2(\gamma_i, \beta_i, \lambda_i)$, and λ_i is a parameter that describes the skewness of the distribution such that $\lambda_i \in (-1, 1)$. The function $\operatorname{sgn}(x)$ characterizes the sign function. We chose λ_i to have an adaptive prior such that $\operatorname{logit}\left(\frac{\lambda_i + 1}{2}\right) \sim \operatorname{Normal}(\hat{\lambda}, \sigma_\lambda)$, where $\hat{\lambda} \sim \operatorname{Normal}(0, 1)$ and $\sigma_\lambda \sim \operatorname{Exponential}(1)$. Notice that this model can be re-parametrized following q_1 and q_2 , allowing us to set the rest of the prior distributions as described for the baseline model (see Supplementary Table 2; Code Availability section). In this case, the Gaussian distribution is a special case of Eq. (5) when $\lambda_i = 0$ (Ashour & Abdel-hameed, 2010).

Fat-tailed and skewed response curve

Finally, one could consider a response curve with both kurtosis and skewness. A convenient way to achieve this is by using a response curve that follows a skewed generalized error distribution. This is a combination of the two distributions described above, containing two parameters that control for both the level and direction of kurtosis and skewness (Fig 1d). The skewed generalized error distribution can be considered by modifying the species’ re-

273 sponse curve in Eq. (1) as

$$\log(p_{ij}) = -\alpha_i - \left(\frac{\gamma'_i |x_j - \beta'_i|}{1 + \lambda_i \operatorname{sgn}(x_j - \beta'_i)} \right)^{\nu_i}, \quad (6)$$

274 where $\gamma'_i = f_1(\gamma_i, \nu_i, \lambda_i)$, $\beta'_i = f_2(\gamma_i, \beta_i, \nu_i, \lambda_i)$, and ν_i and λ_i are parameters that control the
 275 kurtosis and skewness of the distribution, respectively. We define ν_i , λ_i and their prior dis-
 276 tributions as in Eq. 4 and 5, respectively. Again, we can re-parametrize the model following
 277 f_1 and f_2 , and set the rest of the prior distributions as in the baseline model (see Supple-
 278 mentary Table 2; Code Availability section). Notice that the generalized error distribution
 279 (Eq. 4) and the skew normal distribution (Eq. 5) are special cases of Eq. (6) when $\lambda_i = 0$
 280 and $\nu_i = 2$, respectively.

281 Evaluating the log-likelihood

282 One way to understand the shape of distributions is by evaluating the computed log-likelihood
 283 values across the environmental gradient. Comparing these values, one can understand what
 284 aspects of the shape of distributions are missing. To do so, for every sample of the model,
 285 we computed the log-likelihood values and the normalized probability distribution. This
 286 normalized probability is defined such that its maximum is set to 1 for all species in our
 287 dataset. In particular, for a heavy-tailed and skewed response, the normalized probability
 288 distribution was calculated for every sample of the Bayesian model using Eq. (6), where
 289 α_i was set to 0 for any value of x_j . Notice that the normalized probability distribution is
 290 interesting when comparing the log-likelihood values across species because it can be used
 291 to understand whether the model errors are at the tails of the distributions or their center.

292 Results

293 We studied the distribution data to characterize species' realized niches along the main
 294 axis of variation of all environmental variables. Using the presence and absence of species
 295 across sites as the response variable, we sampled the posterior distributions of the baseline
 296 model, accounting for the information in the floristic database regarding species' indicator
 297 values and range of variation. This allowed us to map the center and variance of species'

distributions along the environmental gradient (Fig. 2). Studying the relationship between these properties, we found these to be negatively correlated (i.e. β_i and γ_i in the baseline model were positively correlated; Fig. 2). This means that species found at lower elevations have generally wider distributions than those at higher elevations. The same relationship was found when using instead elevation or mean temperature as explanatory variables (Supplementary Fig. 3) as well as when using ordinal data (Supplementary Fig. 4); however, the pattern was not present along the second main axis of variation of our environmental variables (Supplementary Fig. 5). The comparison between the other parameter estimates revealed additional, somewhat more expected, relationships. In particular, we found the amplitude of distributions to be positively and negatively correlated with their mean and variance, respectively (i.e. α_i is positively correlated with β_i and γ_i ; Supplementary Fig. 6). This implies that, at higher elevations, species' distributions generally have lower amplitudes.

Maintaining the symmetry of species' distributions, we then allowed the kurtosis—or shape of the tails—of these to vary in different ways. To do so, we changed the response curve of our Bayesian model to follow a generalized error distribution (Eq. 4). A comparison of the WAIC values showed this non-linear regression to outperform the baseline model (Supplementary Fig. 7). Studying the resulting posterior distributions, we found the average kurtosis of the distributions to be slightly greater than zero, which corresponds to distributions with longer tails than the Gaussian case (Fig. 3). However, the parameter controlling for the kurtosis ν_i displayed a lot of variation across species (Supplementary Fig. 8), which might indicate that the shape of the tails is species-specific.

Using Eq. (5), we next studied the skewness of species' distributions. Based on the estimates for the WAIC values, this model outperformed the first two (Supplementary Fig. 7), which sheds light on the naturally skewed nature of species' distributions. Perhaps most importantly, studying the mean value of the skewness across species, we found this to be consistently below zero (Fig. 3). This indicates that species' distributions generally present steeper declines towards higher elevations (i.e. $\hat{\lambda} < 0$; Fig. 1). The same was true when using a model that allowed for both fat-tailed and skewed response curves (Eq. 6). This model outperformed the rest, presenting Akaike weights close to 1 (Supplementary Fig. 7), suggesting that both the kurtosis and skewness are useful properties to describe empirical distributions (Fig. 3). A study of how the prior knowledge we had regarding species' en-

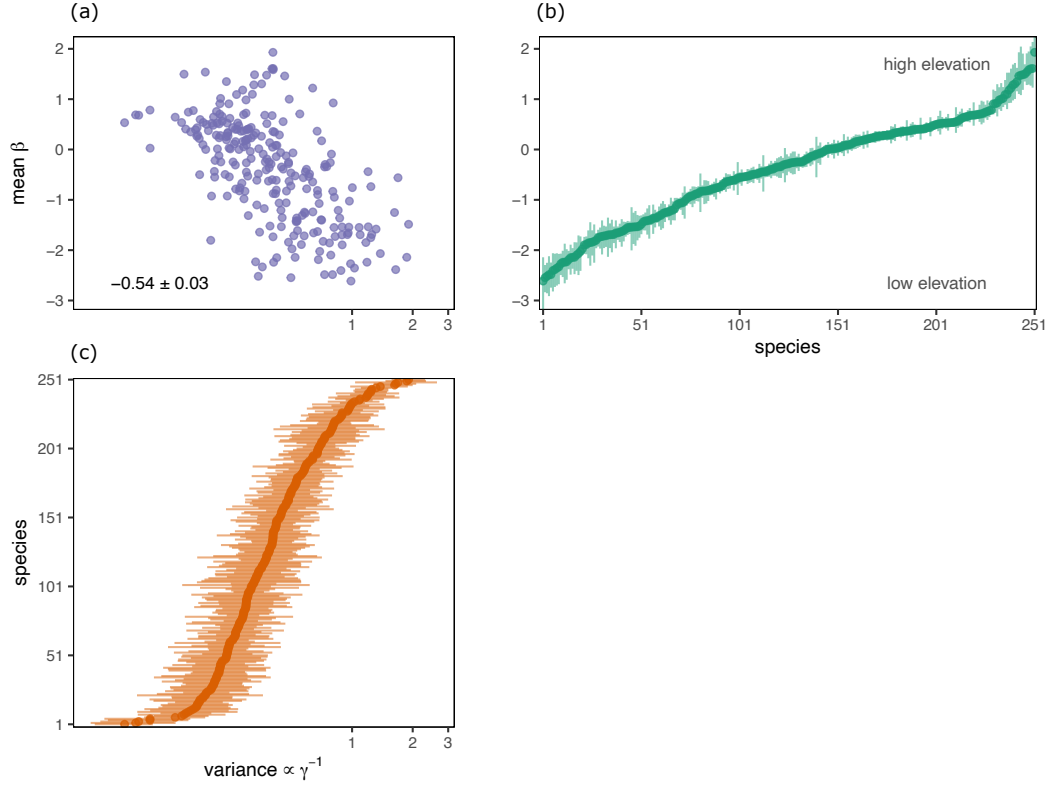


Figure 2: Relationship between the posterior distributions for parameters β_i and γ_i from Eq. (1) across species. Panel (a) describes the relationship between the mean (β_i) and variance ($\propto \gamma_i^{-1}$) of distributions. Every dot characterizes the average value of the corresponding posterior distributions for any given species. The value in the bottom-left corner of the plot displays the Pearson's correlation coefficient between the parameters calculated across all samples. Panel (b) displays the estimates for the center of species' distributions along the environmental gradient. Panel (c) displays the estimates for the variance of species' distributions along the environmental gradient. In (b) and (c), the points represent the mean of the posterior distributions, and the corresponding lines characterize the 89% confidence intervals.

environmental preferences, range of variation and ecological strategies informed the different parameters of the model is presented in the Supplementary Note 1 and Supplementary Fig. 9.

The model characterizing fat-tailed and skewed distributions allowed us to study the posterior distributions for the parameters describing the mean, variance, amplitude, kurtosis and skewness of species realized niches altogether. We observed that different types of species seem to present characteristically different distributions (Fig. 4). Focussing on the negative correlation between the mean and variance of species' distributions, we found some species to escape some of the aforementioned macroecological constraints (Supplementary Fig. XX). Moreover, recent and historical range expanders are often found at lower altitudes, presenting higher amplitudes and distributions that appear to showcase steeper declines towards

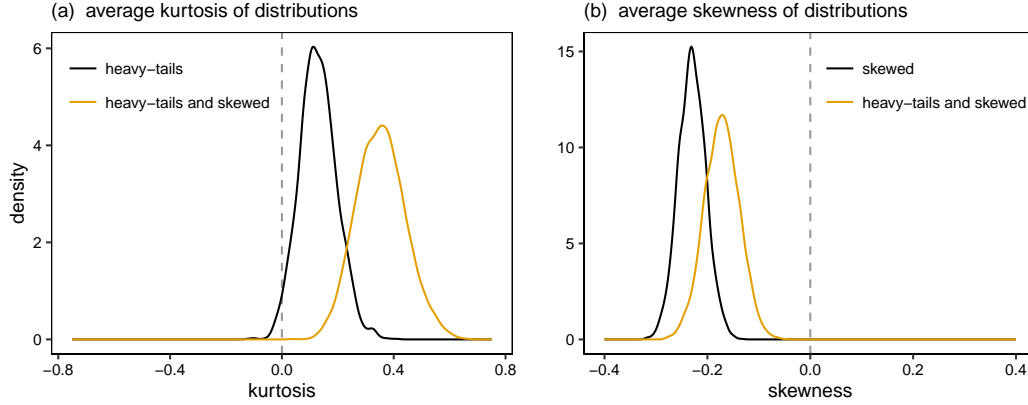


Figure 3: Average kurtosis and skewness of species' distributions. Calculated using the posterior distributions of parameters $\hat{\nu}$ and $\hat{\lambda}$ from the models (see Supplementary Table 2 and ?), the two panels describe the average (a) kurtosis and (b) skewness of distributions. Panel (a) displays the results obtained using response curves that follow a generalized error distribution (black line) and a skewed generalized error distribution (yellow line). Panel (b) displays the results obtained using response curves that follow a skewed normal distribution (black line) and a skewed generalized error distribution (yellow line). In both cases, the gray dotted line indicates the conditions by which species are normally distributed along the environmental axis.

lower elevations. Notice the nature of these results does not depend on the presence or absence of a species at the edge of the sampling area, as the same model produced comparable results when using simulated and bootstrapped data (Supplementary Note XX and Supplementary Fig. XX). Moreover, these results did not substantially change when using ordinal data (Supplementary Fig. XX).

Finally, we wanted to identify what aspects of the shape of distributions we were still missing. We used the computed log-likelihood values and normalized probabilities to understand where our best performing model fails to capture the variation in empirical plant distributions. We found most data points to be located at the tails of distributions (normalized probability ≈ 0) and to present high log-likelihood values (Supplementary Fig. XX). This is not surprising as the study area spans an extensive altitude gradient, and species' distributions are generally narrow relative to it; therefore, the model accurately predicts that species are usually absent in those sampling sites that fall relatively far from the center of their distributions. However, studying instead only those points for which the model did not perform well (with a likelihood ≤ 0.5), we found these to generally be associated with high normalized probabilities (Fig. 5). This indicates that the unexplained variation is often found at the center of species' distributions. Similar results were found when using ordinal data (Supplementary

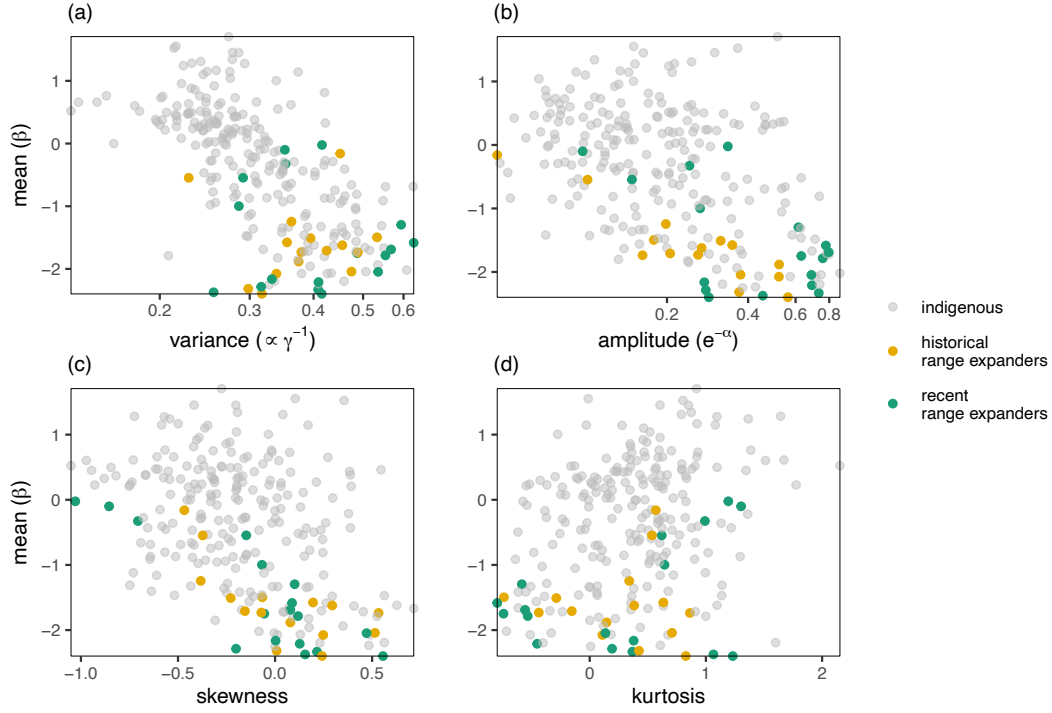


Figure 4: Comparing the distributions of different types of species. Focussing on the differences between indigenous, historical range expanding and recent range expanding species, the panels describe the relationship between the basic properties of their distributions. Panels (a-d) characterize the relationship between the mean, and the variance, amplitude, skewness and kurtosis of the species' distributions (Supplementary Table XX). The points in every panel are calculated as the average value across all samples of the model.

Fig. XX).

Discussion

In this work, we used non-linear response curves to model the distribution of species across an environmental gradient. First, we used a baseline model that considered these as bell-shaped, and we studied the relationship between the basic parameters characterizing them. We found both the amplitude and variance of distributions to be negatively correlated with elevation. Considering more complex response curves, we then found species' distributions to also present non-normal tails and skewed shapes. Specifically, we found species' distributions to generally be characterized by fat tails and steeper declines towards higher elevations. That said, the nature of these distributions was not homogeneous across species, as some of them presented singularly different properties. This is the case of rapid and historical range expanders, often found in warmer environments, with distributions presenting higher ampli-

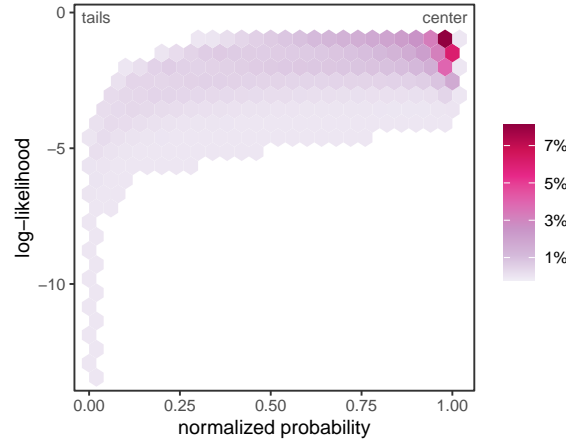


Figure 5: Studying the distribution of log-likelihood values. The graph maps the log-likelihood and normalized probability values for all species across all samples. Notice that in this figure there are only displayed those points that present a likelihood smaller than 0.5. The mapping of all log-likelihood values is presented in Supplementary Fig. XX.

tudes and skewed responses towards high altitudes. Finally, we studied the variation that remained unexplained by the best performing model. We found this unexplained variation to be generally located at the center of distributions, which identified potential general properties of empirical distributions that were missed by our model. Putting this all together, our results uncovered several aspects of the shape of empirical plant distributions and revealed crucial differences between the way species are assembled along environmental gradients.

Our approach allowed us to parsimoniously compare the shape of the species' realized niches along an altitude gradient, testing for the existence of several macroecological patterns. For example, the Rapoport's rule predicts wider ranges of species at higher latitudes and altitudes (Stevens, 1992); and therefore, one might expect a positive correlation between the mean and variance of species distributions. A common explanation for the Rapoport's rule is that climatic variability selects for species with greater climatic tolerances. But while this pattern has been largely studied for multiple systems and across gradients (McCain & Knight, 2013), contrasting evidence suggests that this rule is not pervasive across species (Ribas & Schoereder, 2006; Bhattarai & Vetaas, 2006; McCain & Knight, 2013). Our results seem to contradict the predictions of the Rapoport's rule, as we observed a negative correlation between species' range and elevation. Moreover, other properties of species' distributions—such as their amplitude—were also significantly correlated with species' ranges. This is interesting because it hints at the existence of some general macroecological constraints

that dictate the way different species assemble across environments. That said, our results also suggest that species such as neophytes and archeophytes might not obey this same constraints, as these were singularly positioned along the gradient (Supplementary Fig. XX).

The level of skewness of species' distribution as well as the variability in the shape of their tails diverged from traditionally assumed bell-shaped curves. This allowed us to focus on other interesting macroecological hypotheses. For instance, the so-called abiotic stress limitation hypothesis predicts species' distributions to present steeper declines towards stressful conditions (Austin, 1990). Normand *et al.* (2009) tested this for vegetation data using Huisman *et al.*'s statistical models for several independent species, finding no clear support for such a hypothesis (but see Ziffer-Berger *et al.* 2014). Our results, however, showcased species' distributions to generally present steeper declines towards higher elevations, providing clear evidence of this geographical pattern. Moreover, we were able to highlight the degree to which different species might present different levels of decline towards stressful conditions, as plants found at low elevations—such as recent and historical range expanding species—displayed contrasting levels of skewness. This is important because it could provide glimpses of the different stages of species' assembly processes, with range expanders' distributions trending towards higher elevations.

There are many other properties characterizing empirical distributions that might not have been captured by the different models. One possible way to untangle these properties is by studying the unexplained variation in the empirical data. We observed that this variation is often located at the center of distributions, which suggests that the aspects of their shape not picked up by the models involve those points at the peak of the distributions. This observation is directly linked to another macroecological pattern: the so called abundant-center hypothesis (Sagarin & Gaines, 2002). This hypothesis predicts species to be most abundant at the center of their distributions, and it is an implicit assumption at the core of most modelling approaches. Namely, if one is only willing to assume that species have finite geographic ranges, the abundant-center hypothesis is a consequence our state of ignorance (i.e. the maximum entropy distribution). That said, several studies have pointed out that the abundant-center hypothesis is not pervasive in empirical distributions (Wagner *et al.*, 2011; Pironon *et al.*, 2017; Dallas *et al.*, 2017), suggesting that population abundance could often be more strongly driven by interactions and community structure than the environment (Dallas

418 *et al.*, 2017). Our results, for both binary and ordinal data, support these observations,
 419 suggesting that the species' probability of appearance—as well as likelihood of presenting
 420 high abundance at a given site—might not ubiquitously be highest at the center of their
 421 distributions. Allowing species to showcase other distribution shapes, such as those including
 422 multimodal or plateau peaks, could potentially resolve some of the unexplained variation.
 423 Indeed, studying the tails of species' distributions, we observed several species presenting low
 424 kurtosis levels. While this implied that these distributions had shorter tails than normal, it
 425 also reflected plateau-shaped response curves (e.g. $\nu > 2$ in Fig. 1b).

426 The different hypotheses regarding the shape of species' distributions address central top-
 427 ics in ecology and evolution (Sagarin & Gaines, 2002). These distributions are the result of
 428 environmental variability (Helmuth *et al.*, 2002; Butterfield, 2015), biotic interactions (Hast-
 429 ings *et al.*, 1997) and historical contingencies (Frick *et al.*, 2010), and their shape determines
 430 gene flow (Haldane & Ford, 1956; Lesica & Allendorf, 1995; Pironon *et al.*, 2017) and energy
 431 balances along gradients (Hall *et al.*, 1992). Perhaps most importantly, the shape of species'
 432 distributions will influence their responses to environmental changes (Channell & Lomolino,
 433 2000a), and it could therefore be used as an ecological compass to inform conservation and
 434 management decisions (Channell & Lomolino, 2000b). In this context, we identify two areas
 435 we feel represent key steps from which to move forward. First, trait data could crucially
 436 inform the different parameters controlling the shape of distributions. For example, if the
 437 skewness of species' distributions is the result of uneven environmental tolerances along the
 438 gradient (Sunday *et al.*, 2011), this information should be accounted for analogously to the
 439 way we used the expert knowledge on plants' environmental preferences. The same is true
 440 for species' ecological strategies, with aspects regarding their competitive ability potentially
 441 informing the shape of distributions. Second, from a performance standpoint, the models
 442 presented here will likely do a worse job at predicting species' occurrences than some of the
 443 distribution models developed over the recent years (Norberg *et al.*, 2019), including those
 444 accounting for spatial autocorrelation (Ovaskainen *et al.*, 2016), associations between species
 445 (Tikhonov *et al.*, 2020), and some non-parametric approximations (Harris, 2015). However,
 446 our models have clear interpretable parameters, and can be used to directly compare the
 447 shape of species' realized niches. These comparisons could be used to generate hypotheses
 448 regarding where and when different species might strongly interact with one another along

449 an environmental gradient (Louthan *et al.*, 2015), making ecologically-informed predictions
450 regarding the presence and absence of these relationships (?).

References

- Ashour, S. K. & Abdel-hameed, M. A. (2010). Approximate skew normal distribution. *Journal of Advanced Research*, 1, 341–350.
- Austin, M. P. (1987). Models for the analysis of species’ response to environmental gradients. *Vegetatio*, 69, 35–45.
- Austin, M. P. (1990). Community theory and competition in vegetation. *Community theory and competition in vegetation.*, 215–238.
- Austin, M. P. (2002). Spatial prediction of species distribution: An interface between ecological theory and statistical modelling. *Ecological Modelling*, 157, 101–118.
- Bhattarai, K. R. & Vetaas, O. R. (2006). Can Rapoport’s rule explain tree species richness along the Himalayan elevation gradient, Nepal? *Diversity and Distributions*, 12, 373–378.
- Butterfield, B. J. (2015). Environmental filtering increases in intensity at both ends of climatic gradients, though driven by different factors, across woody vegetation types of the southwest USA. *Oikos*, 124, 1374–1382.
- Channell, R. & Lomolino, M. V. (2000a). Dynamic biogeography and conservation of endangered species. *Nature*, 403, 84–86.
- Channell, R. & Lomolino, M. V. (2000b). Trajectories to extinction: Spatial dynamics of the contraction of geographical ranges. *Journal of Biogeography*, 27, 169–179.
- Dallas, T., Decker, R. R. & Hastings, A. (2017). Species are not most abundant in the centre of their geographic range or climatic niche. *Ecology Letters*, 20, 1526–1533.
- D’Amen, M., Rahbek, C., Zimmermann, N. E. & Guisan, A. (2017). Spatial predictions at the community level: From current approaches to future frameworks. *Biological Reviews*, 92, 169–187.
- Evans, M. E. K., Merow, C., Record, S., McMahon, S. M. & Enquist, B. J. (2016). Towards Process-based Range Modeling of Many Species. *Trends in Ecology & Evolution*, 31, 860–871.

477 Frank, S. A. (2009). The Common Patterns of Nature. *Journal of evolutionary biology*, 22,
478 1563–1585.

479 Frick, W. F., Pollock, J. F., Hicks, A. C., Langwig, K. E., Reynolds, D. S., Turner, G. G.,
480 Butchkoski, C. M. & Kunz, T. H. (2010). An Emerging Disease Causes Regional Popula-
481 tion Collapse of a Common North American Bat Species. *Science*, 329, 679–682.

482 Godoy-Lorite, A., Guimerà, R., Moore, C. & Sales-Pardo, M. (2016). Accurate and scalable
483 social recommendation using mixed-membership stochastic block models. *Proceedings of*
484 *the National Academy of Sciences*, 113, 14207–14212.

485 Guisan, A. & Zimmermann, N. E. (2000). Predictive habitat distribution models in ecology.
486 *Ecological Modelling*, 135, 147–186.

487 Haldane, J. B. S. & Ford, E. B. (1956). The relation between density regulation and natural
488 selection. *Proceedings of the Royal Society of London. Series B - Biological Sciences*, 145,
489 306–308.

490 Hall, C. A. S., Stanford, J. A. & Hauer, F. R. (1992). The Distribution and Abundance of
491 Organisms as a Consequence of Energy Balances along Multiple Environmental Gradients.
492 *Oikos*, 65, 377–390.

493 Harris, D. J. (2015). Generating realistic assemblages with a joint species distribution model.
494 *Methods in Ecology and Evolution*, 6, 465–473.

495 Hastings, A., Harrison, S. & McCann, K. (1997). Unexpected spatial patterns in an insect
496 outbreak match a predator diffusion model. *Proceedings of the Royal Society of London.*
497 *Series B: Biological Sciences*, 264, 1837–1840.

498 Helmuth, B., Harley, C. D. G., Halpin, P. M., O'Donnell, M., Hofmann, G. E. & Blanchette,
499 C. A. (2002). Climate Change and Latitudinal Patterns of Intertidal Thermal Stress.
500 *Science*, 298, 1015–1017.

501 Helmuth, B., Kingsolver, J. G. & Carrington, E. (2004). BIOPHYSICS, PHYSIOLOGICAL
502 ECOLOGY, AND CLIMATE CHANGE: Does Mechanism Matter? *Annual Review of*
503 *Physiology*, 67, 177–201.

- Huisman, J., Olff, H. & Fresco, L. F. M. (1993). A hierarchical set of models for species response analysis. *Journal of Vegetation Science*, 4, 37–46.
- Jamil, T. & ter Braak, C. J. F. (2013). Generalized linear mixed models can detect unimodal species-environment relationships. *PeerJ*, 1, e95.
- Krebs, C. J. (1972). *Ecology: The Experimental Analysis of Distribution and Abundance/by Charles J. Krebs*. 4th edn. Harper & Row, New York.
- Landolt, E., Bäumler, B., Ehrhardt, A., Hegg, O., Klötzli, F., Lämmli, W., Nobis, M., Rudmann-Maurer, K., Schweingruber, F. H., Theurillat, J.-P., Urmi, E., Vust, M. & Wohlgemuth, T. (2010). *Flora indicativa: Ökologische Zeigerwerte und biologische Kennzeichen zur Flora der Schweiz und der Alpen*. Haupt, Bern. ISBN 978-3-258-07461-0.
- Lesica, P. & Allendorf, F. W. (1995). When Are Peripheral Populations Valuable for Conservation? *Conservation Biology*, 9, 753–760.
- Linder, E. T., Villard, M.-A., Maurer, B. A. & Schmidt, E. V. (2000). Geographic range structure in North American landbirds: Variation with migratory strategy, trophic level, and breeding habitat. *Ecography*, 23, 678–686.
- Louthan, A. M., Doak, D. F. & Angert, A. L. (2015). Where and When do Species Interactions Set Range Limits? *Trends in Ecology & Evolution*, 30, 780–792.
- Malmgren, R. D., Stouffer, D. B., Motter, A. E. & Amaral, L. A. N. (2008). A Poissonian explanation for heavy tails in e-mail communication. *Proceedings of the National Academy of Sciences*, 105, 18153–18158.
- McCain, C. M. & Knight, K. B. (2013). Elevational Rapoport’s rule is not pervasive on mountains. *Global Ecology and Biogeography*, 22, 750–759.
- McElreath, R. (2020). *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. CRC press.
- Nadarajah, S. (2005). A generalized normal distribution. *Journal of Applied Statistics*, 32, 685–694.

530 Norberg, A., Abrego, N., Blanchet, F. G., Adler, F. R., Anderson, B. J., Anttila, J., Araújo,
 531 M. B., Dallas, T., Dunson, D., Elith, J., Foster, S. D., Fox, R., Franklin, J., Godsoe, W.,
 532 Guisan, A., O'Hara, B., Hill, N. A., Holt, R. D., Hui, F. K. C., Husby, M., Kålås, J. A.,
 533 Lehtikainen, A., Luoto, M., Mod, H. K., Newell, G., Renner, I., Roslin, T., Soininen, J.,
 534 Thuiller, W., Vanhatalo, J., Warton, D., White, M., Zimmermann, N. E., Gravel, D. &
 535 Ovaskainen, O. (2019). A comprehensive evaluation of predictive performance of 33 species
 536 distribution models at species and community levels. *Ecological Monographs*, 89, e01370.

537 Normand, S., Treier, U. A., Randin, C., Vittoz, P., Guisan, A. & Svenning, J.-C. (2009).
 538 Importance of abiotic stress as a range-limit determinant for European plants: Insights
 539 from species responses to climatic gradients. *Global Ecology and Biogeography*, 18, 437–449.

540 Ovaskainen, O., Roy, D. B., Fox, R. & Anderson, B. J. (2016). Uncovering hidden spatial
 541 structure in species communities with spatially explicit joint species distribution models.
 542 *Methods in Ecology and Evolution*, 7, 428–436.

543 Ovaskainen, O., Tikhonov, G., Norberg, A., Blanchet, F. G., Duan, L., Dunson, D., Roslin,
 544 T. & Abrego, N. (2017). How to make more out of community data? A conceptual
 545 framework and its implementation as models and software. *Ecology Letters*, 20, 561–576.

546 Petrovskii, S., Morozov, A., Taylor, A. E. P. D. & DeAngelis, E. D. L. (2009). Dispersal in
 547 a Statistically Structured Population: Fat Tails Revisited. *The American Naturalist*, 173,
 548 278–289.

549 Pironon, S., Papuga, G., Villellas, J., Angert, A. L., García, M. B. & Thompson, J. D.
 550 (2017). Geographic variation in genetic and demographic performance: New insights from
 551 an old biogeographical paradigm. *Biological Reviews*, 92, 1877–1909.

552 Randin, C. F., Engler, R., Normand, S., Zappa, M., Zimmermann, N. E., Pearman, P. B.,
 553 Vittoz, P., Thuiller, W. & Guisan, A. (2009). Climate change and plant distribution:
 554 Local models predict high-elevation persistence. *Global Change Biology*, 15, 1557–1569.

555 Ribas, C. R. & Schoereder, J. H. (2006). Is the Rapoport effect widespread? Null models
 556 revisited. *Global Ecology and Biogeography*, 15, 614–624.

557 Rohde, K. (1992). Latitudinal Gradients in Species Diversity: The Search for the Primary
 558 Cause. *Oikos*, 65, 514–527.

559 Sagarin, R. D. & Gaines, S. D. (2002). The ‘abundant centre’ distribution: To what extent
560 is it a biogeographical rule? *Ecology Letters*, 5, 137–147.

561 Sagarin, R. D., Gaines, S. D. & Gaylord, B. (2006). Moving beyond assumptions to under-
562 stand abundance distributions across the ranges of species. *Trends in Ecology & Evolution*,
563 21, 524–530.

564 Scherrer, D. & Guisan, A. (2019). Ecological indicator values reveal missing predictors of
565 species distributions. *Scientific Reports*, 9, 1–8.

566 Siefert, A., Lesser, M. R. & Fridley, J. D. (2015). How do climate and dispersal traits limit
567 ranges of tree species along latitudinal and elevational gradients? *Global Ecology and*
568 *Biogeography*, 24, 581–593.

569 Stan Development Team (2021). RStan: The R interface to Stan.

570 Stan Development Team (2021). Stan Modeling Language Users Guide and Reference Man-
571 ual.

572 Stevens, G. C. (1992). The Elevational Gradient in Altitudinal Range: An Extension of
573 Rapoport’s Latitudinal Rule to Altitude. *The American Naturalist*, 140, 893–911.

574 Sunday, J. M., Bates, A. E. & Dulvy, N. K. (2011). Global analysis of thermal tolerance
575 and latitude in ectotherms. *Proceedings of the Royal Society B: Biological Sciences*, 278,
576 1823–1830.

577 ter Braak, C. J. F. & Looman, C. W. N. (1986). Weighted averaging, logistic regression and
578 the Gaussian response model. *Vegetatio*, 65, 3–11.

579 Tikhonov, G., Opedal, Ø. H., Abrego, N., Lehikoinen, A., de Jonge, M. M. J., Oksanen, J.
580 & Ovaskainen, O. (2020). Joint species distribution modelling with the r-package Hmsc.
581 *Methods in Ecology and Evolution*, 11, 442–447.

582 Wagner, V., von Wehrden, H., Wesche, K., Fedulin, A., Sidorova, T. & Hensen, I. (2011).
583 Similar performance in central and range-edge populations of a Eurasian steppe grass
584 under different climate and soil pH regimes. *Ecography*, 34, 498–506.

585 Wong, F. & Collins, J. J. (2020). Evidence that coronavirus superspreading is fat-tailed.
586 *Proceedings of the National Academy of Sciences*, 117, 29416–29418.

587 Ziffer-Berger, J., Weisberg, P. J., Cablk, M. E. & Osem, Y. (2014). Spatial patterns provide
588 support for the stress-gradient hypothesis over a range-wide aridity gradient. *Journal of*
589 *Arid Environments*, 102, 27–33.