

A fresh perspective on distribution modelling: a bayesian framework to understand the distribution of plant species along an environmental gradient

Bernat Bramon Mora^{1,*} and Jake M. Alexander¹

¹Institute of Integrative Biology, ETH Zürich, Zürich, Switzerland; *bernat.bramon@gmail.com

1 **Methods**

2 **Empirical data**

3 We studied the distribution of alpine plant communities along an elevation gradient. To do
4 so, we combined two different datasets: i) one describing the co-occurrence of species across
5 multiple open grasslands in the Swiss Alps, and ii) an extensive floristic database containing
6 environmental and physiological traits for all vegetation across Switzerland (Landolt *et al.*,
7 2010).

8 *Distribution data*

9 We studied the distribution of 798 species across 912 sites covering most of the mountain
10 region of the Western Alps in the Canton de Vaud (Switzerland; Scherrer & Guisan 2019).
11 Each of these sites is a 8 × 8 m plot placed somewhere along an elevation range from 375 m to
12 3210 m. In all sites, presence/absence data as well as Braun-Blanquet abundance-dominance
13 classes were recorded for all species. Additionally, following 30 years (1961–1990) of meteo-
14 rological data from national weather stations, Scherrer & Guisan (2019) calculated multiple
15 climatic variables for each site at high spatial resolution (25 m). Here, we focussed on 9
16 climatic variables, including: daily minimum, maximum and average temperature; sum of
17 growing degree-days above 5°C; mean temperature of wettest quarter; annual precipitation,
18 precipitation seasonality, and precipitation of driest quarter.

19 *Floristic data*

20 To complement the aforementioned distribution data, we used a floristic database of most
21 vegetation across Switzerland. This database was build based on expert knowledge and
22 field experience of botanists and ecologists, and contains information regarding species’
23 environmental preferences and physiological traits. Species’ environmental preferences in
24 this database can be used to inform distribution models—e.g. as an informative prior in
25 a Bayesian framework. These are characterized following the ecological indicator values
26 developed by Landolt *et al.* (2010), providing both an estimate of the average conditions in
27 which a species can be found and a broad description of their range of variation. These val-
28 ues are provided for a range of 10 climatic variables, including temperature, continentality,
29 light conditions, as well as moisture, acidity and nutrient content of the soil (see a full list
30 and description of the ecological indicators in the Supplementary Methods; Landolt *et al.*
31 2010). On the other hand, the information regarding species’ physiological traits represent
32 general descriptions of species’ growth and life strategies—examples include their growth
33 forms, nature of the storage organs, dispersal ability and pollinator agents. In total, we
34 identify more than 120 binary traits that characterize the physiology of species (see a full
35 list and description of the ecological indicators in the Supplementary Methods; Landolt
36 *et al.* 2010).

37 *[Trait data]*

38 This could be Tom’s data if we end up using it.

39 **Distribution model**

40 There is a long list of model structures well suited to characterize species’ distributions (see
41 XX for a review); however, we were interested in a model that explicitly incorporates all in-
42 formation regarding plant’s environmental preferences found in the floristic database. More
43 specifically, we wanted to account for the climatic indicator values and range of variation
44 registered for all plants in our dataset. These two values provide basic information regard-
45 ing plant’s optimal environmental conditions and width of their distributions. Therefore,

we first formulated a baseline model that directly accounts for such prior information.

Baseline model

Given y_{ij} the presence/absence of any species i in any given site j , and a set of k environmental variables x_{jk} , we estimate species' distributions as:

$$\begin{aligned}
y_{ij} &\sim \text{Binomial}(1, p_{ij}) \\
\log(p_{ij}) &= -\alpha_i - \sum_k \lambda_{ik} (x_{jk} - \beta_{ik})^2 \\
\log(\alpha) &\sim \text{MVNormal}(\hat{\alpha}, \Sigma^\alpha) \\
\beta_{ik} &\sim \text{MVNormal}(\hat{\beta}_k, \Sigma^{\beta_k}) \\
\log(\lambda_{ik}) &\sim \text{MVNormal}(\hat{\lambda}_k, \Sigma^{\lambda_k}) \\
\hat{\alpha}, \hat{\lambda}^k, \hat{\beta}^k &\sim \text{Normal}(0, 1)
\end{aligned} \tag{1}$$

Notice that this model structure assumes all plants to have a uni-modal distributions along each environmental axis (see the model's behaviour in Supplementary Figure XX), where parameters α_i , β_i^k , and λ_i^k describe amplitude of the probability p_{ij} , species' average climatic suitability and range of variation along the different environmental gradients, respectively. While potentially sacrificing predictive accuracy, this model structure allows us to explicitly incorporate all prior knowledge that we have regarding species' distributions via Σ^α , Σ^{β_k} and Σ^{λ_k} . More specifically, we express β_i^k and $\log(\lambda_i^k)$ as multivariate normal distributions—i.e. Gaussian processes—such that Σ^{β_k} and Σ^{λ_k} are variance-covariance matrices describing species' similarity in terms of their average climatic suitability and range of variation along the different environmental gradients, respectively. Likewise, $\log(\alpha)$ is characterized as a Gaussian Process, where the corresponding variance-covariance matrix Σ^α is designed to also incorporate some of the prior information that we have with regards to species' physiological traits.

In all cases, all variance-covariance matrices are defined as follows:

$$\Sigma_{ij}^\chi = \eta_\chi \exp\left(-\rho_\chi D_{ij}^{\chi^2}\right) + \delta_{ij} \sigma_\chi, \tag{2}$$

64 where Σ_{ij}^χ describes the covariance between any pair of species i and j for any given
 65 parameter α_i , β_i^k , and λ_i^k . Following this expression, such covariance declines exponentially
 66 with the square of the different D_{ij}^χ , which are distance measures computed using the
 67 prior information that we have regarding species' distributions. Specifically, given α_i , β_i^k ,
 68 and λ_i^k , the distance measures are calculated using plants' physiological traits, ecological
 69 indicator values and range of variation, respectively (see below for further details). For each
 70 covariance matrix, the hyperparameter ρ_χ determines the rate of decline of the covariance
 71 between any two species, and η_χ defines its maximum value. The hyperparameter σ_χ
 72 describes the additional covariance between the different observations for any given species.
 73 For any given hyperparameter, we choose adaptive priors across covariance structures.
 74 That is, and taking ρ_χ as an example, we choose a prior $\log(\rho_\chi) \sim \text{Normal}(\hat{\rho}, \sigma_\rho)$ such that
 75 $\hat{\rho} \sim \text{Normal}(0, 1)$ and $\sigma_\rho \sim \text{Exponential}(1)$. Similar priors were chosen for both η_χ and
 76 σ_χ . We generated the posterior samples for the Bayesian models with the help of the R
 77 package 'rstan' to (Team *et al.*, 2019).

78 *Distance matrices*

79 The missing component in the description of model (1) is the distance matrices D^χ used
 80 to define the covariance matrices Σ^α , Σ^{β^k} and Σ^{λ^k} . In this model, such distance matrices
 81 characterize differences between plant species. In the floristic data, however, the prior infor-
 82 mation that we have for these differences is represented by a set of ordinal and categorical
 83 traits. More specifically, both the ecological indicator values and range of variation—which
 84 define the prior information that we have for β_i^k , and λ_i^k , respectively—are ordinal traits
 85 specified for all species. In contrast, the plants' physiological data—shaping the prior for
 86 the parameters α_i —are characterized by categorical data containing multiple missing en-
 87 tries. Therefore, we need to carefully compile this data into distance matrices in order to
 88 be able to feed this prior information into the model.

89 More generally, we want to understand the way N species are characterized by M categor-
 90 ical traits. One way to frame this problem is by using a network representation. Following
 91 the ideas presented by Godoy-Lorite *et al.* (2016), we assume that species can be connected
 92 to each of these traits by an interaction (i, j) that can be of any type $r \in R$. Notice that this

provides as with multiple ways to account for the information—and lack thereof—contained in the different categorical and ordinal traits M . That is, the R types of interactions can represent the lack of information for a particular link (i, j) , the absence or presence of such interaction, and any type of association between i and j .

Given a set of interactions R^* between N and M , we use a Mixed Membership Stochastic Block Model (MMSBM) to characterize these. In particular, we consider that plants and traits can be classified into K and L groups, respectively. For every species i , we assume that there is a probability $\theta_{i\alpha}$ for it to belong to any of the K species groups. Likewise, we also assume that any trait j has a probability $\phi_{j\beta}$ of belonging to any of the L trait groups. Finally, we define $p_{\alpha\beta}(r)$ as the probability of a species from group α interacting with a trait from group β by an association type r . Putting these together, the probability of an interaction (i, j) of type r can be calculated as:

$$Pr[r_{ij} = r] = \sum_{\alpha\beta} \theta_{i\alpha} \phi_{j\beta} p_{\alpha\beta}(r) \quad (3)$$

Following this definition, we want to find the group memberships that maximize the likelihood $P(R^*|\theta, \phi, p)$. Doing so is difficult optimization problem; however, it has been shown that one can estimate the different $\theta_{i\alpha}$, $\phi_{j\beta}$, and $p_{\alpha\beta}(r)$ parameters by maximizing the likelihood using an expectation-maximization algorithm (Godoy-Lorite *et al.*, 2016; Tarrés-Deulofeu *et al.*, 2019). In simple terms, one can iteratively find multiple local minima for the likelihood, and average over the estimated the parameter values (Godoy-Lorite *et al.*, 2016)[†].

The average estimates for the group memberships provide us with a different scale to classify species based on the traits these have. In short, for any species i , we can estimate a K -dimensional vector $\vec{\theta}_i$ that describes the extend to which i belong to each group

[†]While this averaging is trivial for the estimated probabilities $Pr[r_{ij} = r]$, it is non-trivial if one wants to find averages for the group memberships. The reason for this is related to the stochastic nature of the expectation-maximization algorithm. This algorithm initially assigns random group memberships to both species and traits. While this random labelling is irrelevant when studying the probabilities $Pr[r_{ij} = r]$, it is instead crucial for averaging $\theta_{i\alpha}$, $\phi_{j\beta}$, and $p_{\alpha\beta}(r)$. Therefore, before averaging the group membership estimates, one needs to find the bijective relationship for the labellings of different iterations of the optimization algorithm. In a nutshell, for every iteration, I do this by using a simulated annealing algorithm on the estimated $p_{\alpha\beta}(r)$, matching the corresponding labelling to a reference iteration. A full section for this is required in the Supplementary Information.

membership—i.e. the extend to which a species is of one type or another. This classification is useful because it can be used to compare species, defining a way to measure the distance between species based on an arbitrary—and potentially incomplete—set of categorical or ordinal traits M . The simplest case is to define the distance as $D_{ij} = |\vec{\theta}_i - \vec{\theta}_j|$. Alternatively, one could also define K distance matrices based on the different group memberships $D_{ij}^\alpha = |\theta_{i\alpha} - \theta_{j\alpha}|$.

121 *Modifying the variance-covariance structures*

The model structure defined in Eq. (1) allows us to test the effect of adding new information. Specifically, we can do this by modifying Eq. (2). For example, imagine that we have multiple matrices D^k characterizing species' differences along different axis of variation—i.e. two matrices characterizing ecological and environmental traits, or multiple matrices resulting from the different group memberships estimated using the MMSBM. One could modify Eq. (2) for a particular parameter—e.g. parameter α_i —such that

$$\Sigma_{ij}^\alpha = \eta_\alpha \exp \left(- \sum_k \rho_{\alpha k} D_{ij}^{k,2} \right) + \delta_{ij} \sigma_\alpha, \quad (4)$$

where now $\rho_{\alpha k}$ are separate relevance hyperparameters for each distance matrix in the total variance of α_i . Notice that the same is true for the covariance of parameters β_i^k and λ_i^k . Finally, for all hyperparameters and as described for the baseline model, we use adaptive priors across covariance structures.

132 **Results**

133 **Discussion**

134 **References**

135 Godoy-Lorite, A., Guimerà, R., Moore, C. & Sales-Pardo, M. (2016). Accurate and scalable
136 social recommendation using mixed-membership stochastic block models. *PNAS*, 113,
137 14207–14212.

- 138 Landolt, E., Bäumler, B., Ehrhardt, A., Hegg, O., Klötzli, F., Lämmli, W., Nobis,
139 M., Rudmann-Maurer, K., Schweingruber, F. H., Theurillat, J.-P., Urmi, E., Vust, M.
140 & Wohlgemuth, T. (2010). *Flora indicativa: Ökologische Zeigerwerte und biologische*
141 *Kennzeichen zur Flora der Schweiz und der Alpen*. Haupt, Bern. ISBN 978-3-258-07461-
142 0.
- 143 Scherrer, D. & Guisan, A. (2019). Ecological indicator values reveal missing predictors of
144 species distributions. *Scientific Reports*, 9, 1–8.
- 145 Tarrés-Deulofeu, M., Godoy-Lorite, A., Guimerà, R. & Sales-Pardo, M. (2019). Tensorial
146 and bipartite block models for link prediction in layered networks and temporal networks.
147 *Phys. Rev. E*, 99, 032307.
- 148 Team, S. D. *et al.* (2019). RStan: the R interface to Stan. R package version 2.19.1.