

¹ Spatially continuous identification of beta diversity
² hotspots using species distribution models

³ Advisory Committee Document

⁴ Gabriel Dansereau

⁵ November 25, 2019

⁶ **Advisory Committee Document**

⁷ **Abstract**

⁸ Beta diversity is an essential measure to describe the organization of biodiversity in space.
⁹ The calculation of local contributions to beta diversity (LCBD), specifically, allows for the
¹⁰ identification of sites with exceptional diversity within a region of interest, which is useful
¹¹ for both community ecology and conservation purposes. However, beta diversity implies a
¹² comparison among the sites of a given region, thus, its use is restricted to sites with known
¹³ species composition, and to discontinuous spatial scales. We therefore propose a method
¹⁴ to calculate LCBD indices on continuous scales for a whole region of interest, including
¹⁵ unsampled sites. First, species distributions can be predicted on continuous scales using
¹⁶ species distribution models (SDM). These models, such as the BIOCLIM method, use the
¹⁷ environmental conditions at sampled sites to predict the presence or absence of each species at
¹⁸ unsampled locations. Second, LCBD statistics can then be computed on the SDM predictions.
¹⁹ We therefore show that it is possible to identify beta diversity hotspots on spatially continuous
²⁰ and extended scales. Our results confirm that LCBD values are related to species richness,
²¹ and that species-poor sites contribute most to beta diversity.

²² **Introduction**

²³ Beta diversity, defined as the variation in species composition among sites in a geographic
²⁴ region of interest (Legendre, Borcard, and Peres-Neto 2005), is an essential measure to
²⁵ describe the organization of biodiversity in space. Total beta diversity within a community
²⁶ can be partitioned into local contributions to beta diversity (LCBD) (Legendre and De Cáceres
²⁷ 2013), which allows for the identification of sites with exceptional species composition,
²⁸ hence exceptional biodiversity. Such a method is useful for both community ecology and
²⁹ conservation biology, as it highlights sites that are most important for their research or
³⁰ conservation values. However, LCBD calculation methods require complete information
³¹ on community composition, such as a community composition matrix Y , thus they are
³² inappropriate for partially sampled or unsampled sites. To our knowledge, these methods
³³ have mostly been applied on community data from sampled sites, thus on discontinuous
³⁴ spatial scales, e.g. at intervals along a river stream (Legendre and De Cáceres 2013). This
³⁵ raises the following questions: 1) could LCBD indices be extended to continuous spatial
³⁶ scales, and 2) could this provide novel ecological insights in poorly sampled regions? We
³⁷ aim to answer these questions by combining the LCBD calculation methods with predictive
³⁸ biogeography approaches, and suggest that this would allow for the identification of sites
³⁹ with high conservation value in poorly sampled regions.

⁴⁰ Species distribution models (SDMs) already allow to make predictions on continuous spatial
⁴¹ scales which could be used to calculate LCBD indices. These methods, also known as
⁴² bioclimatic envelope models (Araújo and Peterson 2012), aim to predict species presence or
⁴³ absence based on observation of occurrences at known locations (Poisot et al. 2019). This way,
⁴⁴ they generate novel ecological insights, and represent an approach yet to be applied to LCBD.
⁴⁵ We believe that such an approach of generating novel ecological insights for unsampled or
⁴⁶ lesser-known locations could be an interesting new perspective in the study. Through them,
⁴⁷ we would be able to expand community information already available, and thus work on a
⁴⁸ much larger community matrix than in typical LCBD studies.

⁴⁹ Appropriate data to expand measures of exceptional biodiversity through space is increasingly
⁵⁰ available online. For instance, the Worldclim 2.0 database (Fick and Hijmans 2017) provides
⁵¹ interpolated climate data for global land areas at very high spatial resolution, and the eBird

52 platform (Sullivan et al. 2009) provides a growing citizen-contributed database of worldwide
53 bird observations. Both of these are commonly used in SDMs, and offer relevant information
54 on extended spatial scales. Hence, we believe that we could use them to predict community
55 composition and calculate LCBD indices on continuous spatial scales, and that the result
56 would be representative of the true community structure.

57 The predictive approach we suggest would be especially useful in poorly sampled regions,
58 or in regions with only sparse sampling. While it doesn't replace a full sampling within the
59 community, it does provide relevant ecological insights. For instance, the method could help
60 identify unsampled sites with potential conservation value which should be targeted as soon
61 as possible in future studies. We also believe that our method could also be combined with
62 IPCC climate change scenarios, which provide projections for climate variables, in a way that
63 would allow us to model beta diversity changes with climate change and to identify the sites
64 where the changes in the community will be most important. Again, this method would be
65 more relevant as an informative approach to suggest sites to prioritize for future conservation
66 and more structured research.

67 In this document, we cover in more details the methods that we suggest for this research
68 project. The preparation part of the project, including data collection and manipulation, has
69 already been done, and a workflow for the analyses, including code implementation, has been
70 defined as well. We also detail preliminary analyses and results intended as proof-of-concept
71 for the approach, which of course needs to be refined. Finally, we discuss methods that we
72 intend to use in future analyses, and whose feasibility is not as clearly stated.

73 **Methods**

74 **1. Data Collection**

75 We decided to focus our analyses on bird species and collected the data available on eBird
76 for the Warblers family. The complete database contains nearly 600 million observations,
77 and presents two main advantages over other large scale datasets (Johnston et al. 2019): 1)
78 data is structured as checklist and users can explicitly specify their observations as "complete
79 checklists" when all detected species were reported, which allows to infer information on
80 species absences, 2) the dataset is semi-structured and checklists are associated with metadata

describing sampling effort, such as duration of search, distance travelled and number of observers, which can be used as controls in the analyses. We chose to focus specifically on the Warblers family, as it is a diverse group, popular among birders, with over 30 million observations.

We decided to restrict our analyses to North America and collected climate data available in the WorldClim 2 database (Fick and Hijmans 2017). We believe North America represents a suitable scale, large enough to cover a lot of variation in environmental variables and community structure, as well as phenomena such as species migration. We also expect such extent of the spatial scale to cover for imprecision in estimated species ranges. The WorldClim data consists of spatially interpolated monthly climate data for global areas, available for resolutions from 10 arc-minutes to 30 arc-seconds. The variables used are provided in Table 1, and consists of different measures of temperature and precipitation. We chose to use the coarser 10 arc-minutes resolution in our analyses, again to cover for imprecision, and because we believe it is sufficient for proof of concept.

Table 1: WorldClim 2 climate variables used in the analyses

Variable	Description
1	Annual Mean Temperature
2	Mean Diurnal Range (Mean of monthly (max temp - min temp))
3	Isothermality (BIO2/BIO7) (* 100)
4	Temperature Seasonality (standard deviation *100)
5	Max Temperature of Warmest Month
6	Min Temperature of Coldest Month
7	Temperature Annual Range (BIO5-BIO6)
8	Mean Temperature of Wettest Quarter
9	Mean Temperature of Driest Quarter
10	Mean Temperature of Warmest Quarter
11	Mean Temperature of Coldest Quarter
12	Annual Precipitation
13	Precipitation of Wettest Month

Variable	Description
14	Precipitation of Driest Month
15	Precipitation Seasonality (Coefficient of Variation)
16	Precipitation of Wettest Quarter
17	Precipitation of Driest Quarter
18	Precipitation of Warmest Quarter
19	Precipitation of Coldest Quarter

95 **2. Data Manipulation**

96 WorldClim variables and eBird occurrence data are provided in different formats, so they
 97 require some manipulation to be combined together. WorldClim variables are provided in
 98 a 2-dimensional grid format, useful for large scale analyses and visualization, where each
 99 cell or pixel corresponds to the resolution of 10 arc-minutes. Each of the 19 variables forms a
 100 different grid. On the other hand, eBird records are occurrence-based, so each entry in the
 101 dataset corresponds to an observation of a single species at a given location. These entries can
 102 easily be matched to the 2-D grid format of the WorldClim variables through their spatial
 103 coordinates, which we found more useful for large scale analyses and visualization. Hence,
 104 for each species, we matched all occurrences in eBird to the grid format of the WorldClim
 105 variables, and later created a presence-absence community matrix Y , with the sites being the
 106 grid cells. We also applied the Hellinger transformation on the raw presence-absence data,
 107 although the most appropriate method remains to be determined, especially since the data
 108 has to be compared with the SDM predictions. All data manipulations and further analyses
 109 were realized in *Julia v1.2.0* (Bezanson et al. 2017) with the basic structure built around the
 110 soon-to-be-released `SimpleSDMLayers.jl` package.

111 **3. SDM – The BIOCLIM method**

112 We used the BIOCLIM method to predict species distributions. BIOCLIM, first introduced by
 113 (Nix 1986), is considered as the classic “climate-envelope-model”, and is now available to
 114 users through the `dismo` package in R (Hijmans et al. 2017). It has long been outperformed
 115 by other methods (Elith et al. 2006), but it is still commonly used for its simplistic approach

116 and ease of understanding, as well as its simple relation to niche theory (Booth et al. 2014;
117 Hijmans et al. 2017). It is also a method designed for presence-only data, which does not
118 require information on absences, nor take them into account if provided (as in our case).
119 Despite that, we chose this method for our preliminary analyses as it was easier to implement
120 and because we believe it to be sufficient for proof-of-concept. We discuss possible alternatives
121 in the “Alternative methods” section below.

122 Briefly, the BIOCLIM method defines species potential range as a multidimensional envi-
123 ronmental hypervolume bounded by the minimum and maximum values of all presences
124 (Franklin 2010). For each species, the algorithm establishes the percentile distribution of
125 the values of each environmental variables at the known locations of occurrences (Hijmans
126 et al. 2017). The environmental variables of all sites are then compared to those percentile
127 distributions and given scores between 0 (1st percentile) and 1 (100th percentile). The median
128 or 50th percentile is considered as the most suitable location and both tails (e.g. 10th and
129 90th percentile) are not distinguished, the values larger than 0.5 being subtracted from 1.
130 The minimum percentile score across all environmental variables is selected as the prediction
131 value for each site and multiplied by 2 so values are between 0 and 1 (Hijmans et al. 2017).
132 It should be noted that the limiting variable is thus not necessarily the same for all sites.
133 Values of 1 are rare, as it would mean a perfectly median site on all variables, and values
134 of 0 are frequent, since they are assigned whenever an environmental value is outside the
135 range of the observed values (Hijmans et al. 2017). Finally, before calculating richness or beta
136 diversity metrics, we transformed the predictions back to a presence-absence format, where
137 all predictions greater than one are considered as presence. This might tend to overestimate
138 species ranges and create some sort of border effect, but we believe the effects will be mitigated
139 given the spatial extent and coarse scale of our study.

140 4. LCBD calculation

141 We calculated the LCBD statistics through the total variance of the matrix Y for both the raw
142 data and SDM predictions. Legendre and De Cáceres (2013) showed that LCBD coefficients
143 can be calculated directly through the total variance of matrix Y , or through a matrix of
144 dissimilarities among sampling units. We chose the first approach, as it also allows to
145 compute species contributions to beta diversity (SCBD), although we did not investigate it for

now. First, the presence-absence matrix Y had to be transformed in an appropriate way, as mentioned earlier. We chose to apply the Hellinger transformation to the raw data and no transformation on the SDM predictions for now, as the most appropriate one still needs to be determined. We then computed a matrix S of squared deviations from column means and summed all the values of S to obtain the total sum of squares (SS) of the species composition data (Legendre and De Cáceres 2013). LCBD are then computed as $LCBD_i = SS_i/SS_{Total}$, where SS_i is the sum of squares of a sampling unit i . Finally, since our matrix Y is very large, the LCBD coefficients are very small, so we scaled them to the maximum value.

5. Prediction validity

The exact way of testing the validity of the predictions remains to be determined, and will also depend on the exact methods used to make the SDM predictions. A key element to note is that both SDM predictions and LCBD values will have to be validated, so will likely require different methods. Many metrics are well documented in the literature to test SDM predictions, such as the Kappa index (Franklin 2010), and could be used for the BIOCLIM predictions. Another possible way would be to separate the data into a training and testing dataset, with 70% and 30% of the data for instance, which is a common approach in machine learning techniques. However, this approach reduces the amount of data that can be used in the model, and raises the issue of making sure that the datasets are both random and representative of the data, as well as the community dynamics. Also, in this framework, the testing data cannot be considered as independent, which prevents using it in certain tests of significance. One interesting approach, suggested by (Elith et al. 2006) for SDMs, would be to find independent, well-structured presence-absence datasets for validation, on which beta diversity metrics has or could be calculated. This validation might not cover the entire extent of the predictions, but it might bring interesting perspectives if combined with other validation methods, mostly because it would bring a closer comparison to the way LCBD metrics are used at the moment.

6. Alternative methods

Other methods could possibly outperform BIOCLIM for the predictions, as have already proven by Elith et al. (2006). Better predictions will come by two different means: 1)

175 approaches that are better than BIOCLIM to model the relationship between species presence-
176 absence (or even abundance) and environmental variables, and 2) approaches that account
177 for other drivers of species distributions, such as ecological interactions for instance. The
178 most obvious alternative to BIOCLIM is MAXENT (Phillips, Anderson, and Schapire 2006),
179 another presence-only method that has come to be one of the most widely used methods.
180 Machine learning methods would be also be interesting alternatives that have been proven to
181 outperform BIOCLIM (Franklin 2010). Random Forests, especially are simple methods to put
182 in place, allow for quantification of the variables importance in explaining variation, and offer
183 intrinsic testing metrics. Neural networks could also be an interesting alternative. However,
184 while those methods might return more accurate predictions, they do not implicitly model
185 other drivers of species distribution, among which species interactions and functional niche.
186 Integrating those factors might prove more difficult given our dataset and our focus Warblers
187 species, as no appropriate information on their interaction is available to our knowledge.
188 Joint species distribution models (JSDMs) might be an interesting way to encompass those,
189 as they attempt to model species cooccurrence, rather than the distribution of single species
190 distributions (Pollock et al. 2014). A different taxonomic group and data datasets could also
191 be used with more details on interactions could also be used, though having a method that
192 can be applied to any taxonomic group would be more interesting. Yet, such an approach
193 might prove to be beyond the scope of the present research.

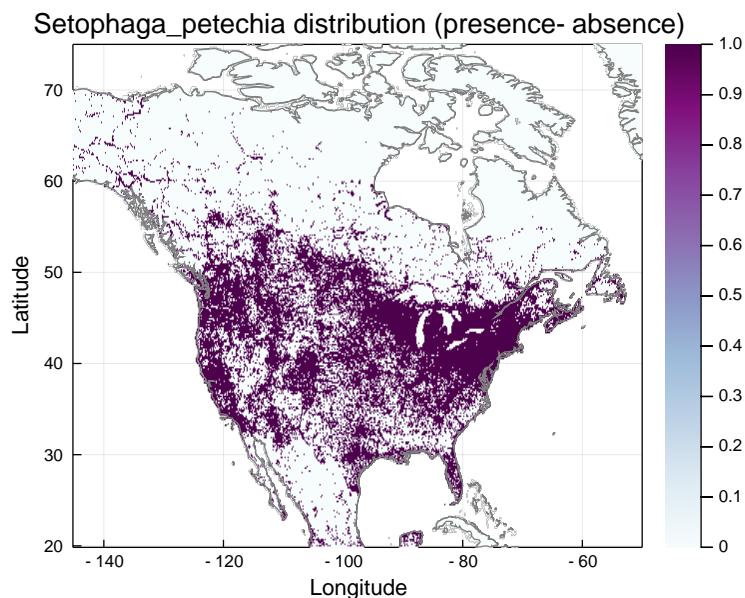
194 7. Climate change scenarios & temporal beta diversity

195 We aim to apply our method to environmental conditions from climate change scenarios, first
196 to model community compositions after climate change on continuous scales through SDMs,
197 and then to identify the sites where the community has changed in the most exceptional ways.
198 This can be done through LCBD values, but also through temporal beta diversity indices
199 (TBI) (Legendre 2019), which allow to study changes in community composition through
200 time from repeated surveys at given sites. Whereas LCBD values essentially measure the
201 contribution to beta diversity of each site compared to all other ones, TBI measure changes in
202 community composition for a single site between two surveys, and can also be decomposed
203 into species losses and gains. Moreover, TBI can be tested for significance using a permutation
204 test. An approach similar to that of (Legendre and Condit 2019) would be most interesting to

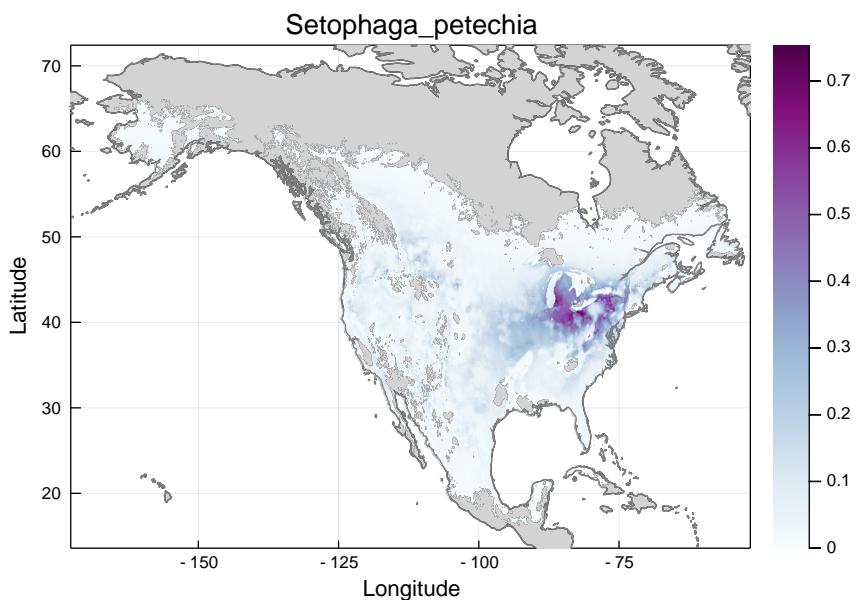
follow: they first computed LCBD indices and compared the sites that were significant for two surveys 30 years apart, highlighting a swamp region where important changes seemed to have occurred, and then used TBI indices to confirm the sites with significant changes, decompose those into losses and gains and identify the species that had changed the most. Such an approach could be highly informative with our data, although the permutation tests and corrections to apply might cause problems given the number of sites that would be implied in our study. The possibility of using climate change scenarios in the SDMs also needs to be investigated in more details. We did not try to download nor find the appropriate data for now, but we found that the interpolated variables are sometimes different than those used in Worldclim 2.0. The SDM models and predictions might therefore be slightly different than those used for the LCBD calculations, and potentially less reliable. Nonetheless, we believe it will be possible to do some kind of time analysis linking beta diversity, climate change, and species distribution modelling, which could return highly informative results for conservation purposes.

219 Preliminary Results

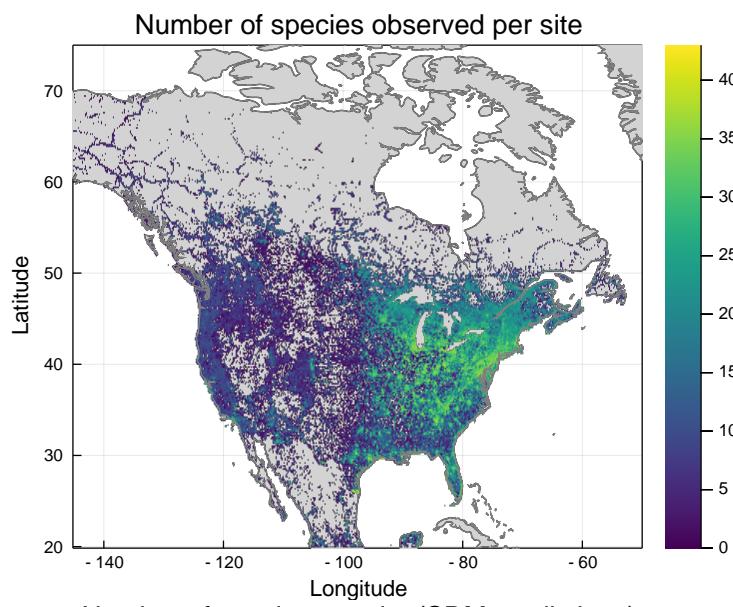
220 Our preliminary results mainly compare raw data statistics to prediction statistics. (Raw &
221 SDM figures will be presented side-by-side)



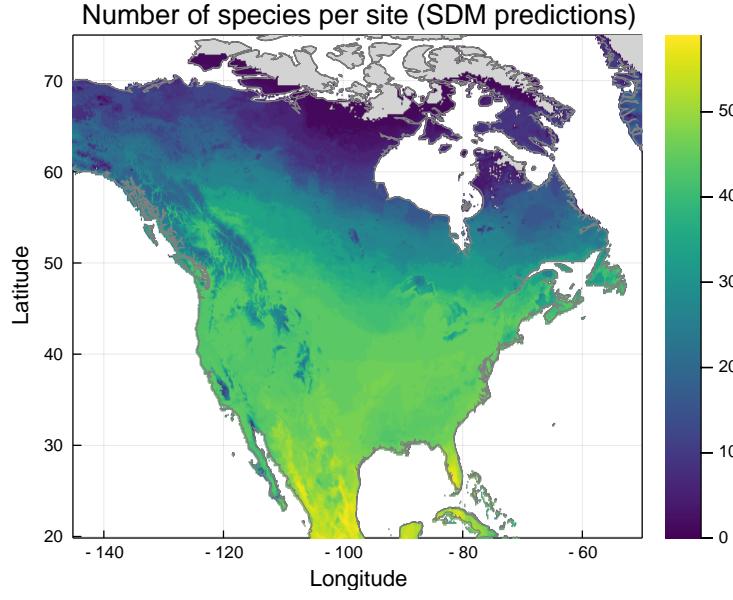
222



223

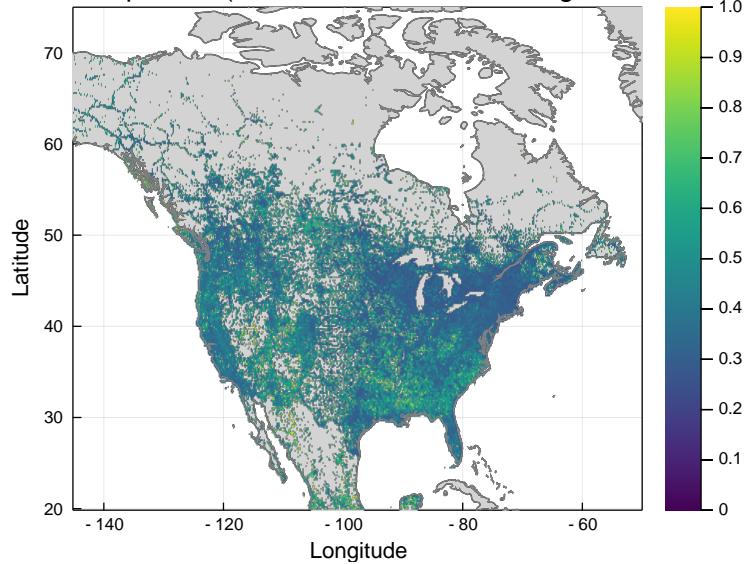


224



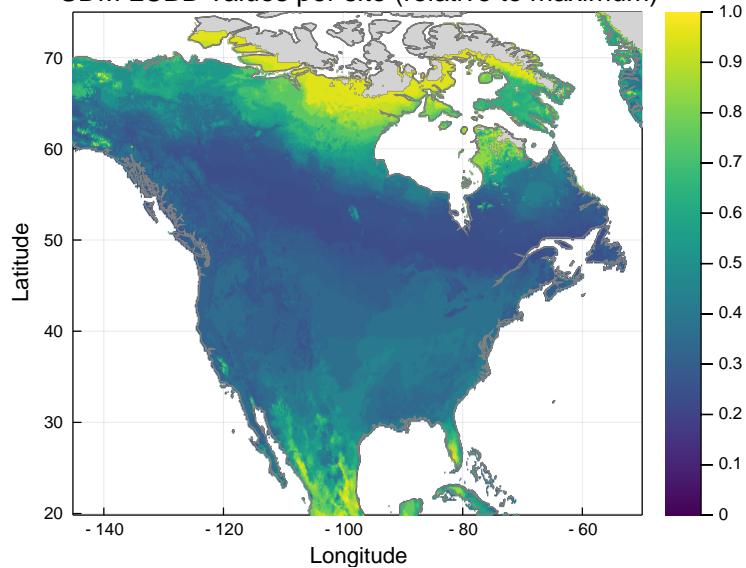
225

LCBD values per site (relative to maximum, hellinger transformed)



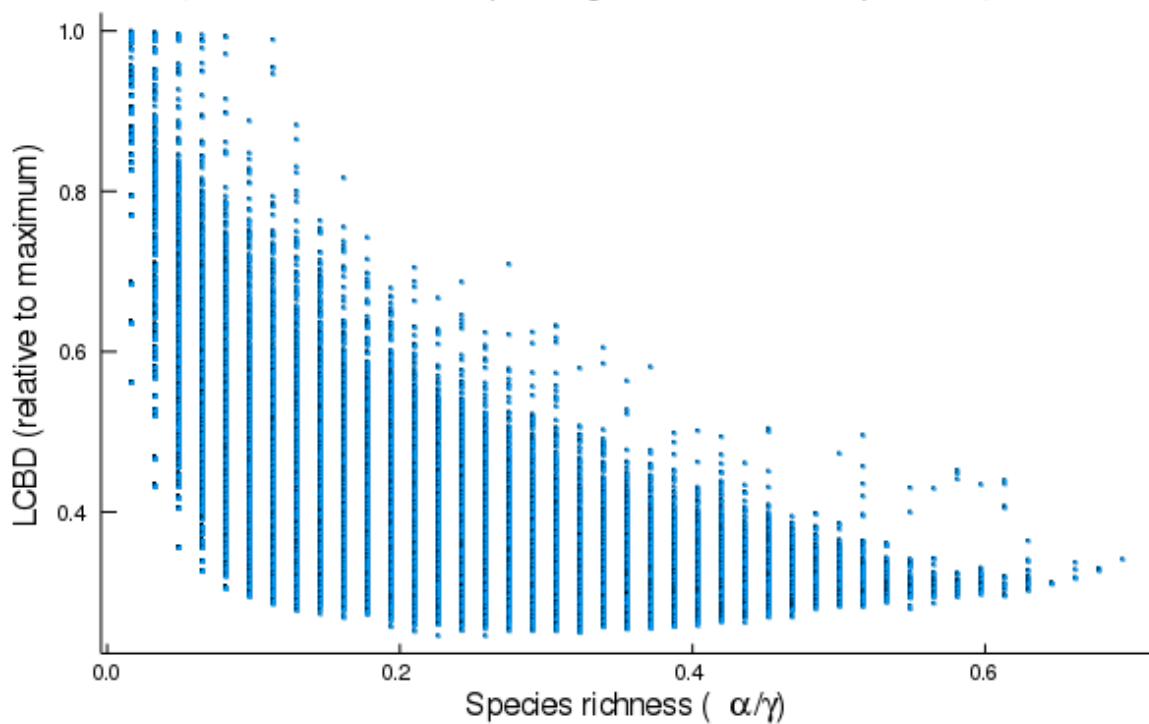
226

SDM LCBD values per site (relative to maximum)



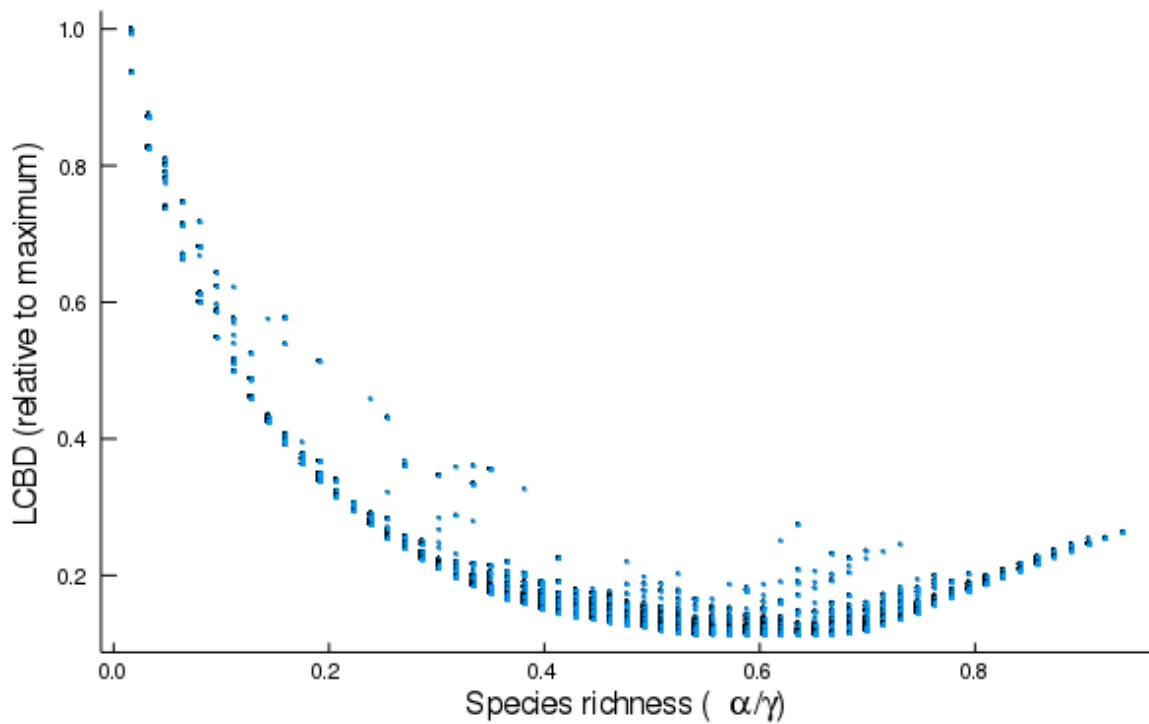
227

Relationship between LCBD (hellinger transformed) and species richness



228

Relationship between LCBD (hellinger transformed) and species richness



229

230 **References**

- 231 Araújo, Miguel B., and A. Townsend Peterson. 2012. "Uses and Misuses of Bioclimatic
232 Envelope Modeling." *Ecology* 93 (7): 1527–39. <https://doi.org/10.1890/11-1930.1>.
- 233 Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B. Shah. 2017. "Julia: A Fresh
234 Approach to Numerical Computing." *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- 236 Booth, Trevor H., Henry A. Nix, John R. Busby, and Michael F. Hutchinson. 2014. "BIOCLIM:
237 The First Species Distribution Modelling Package, Its Early Applications and Relevance to
238 Most Current MaxEnt Studies." *Diversity and Distributions* 20 (1): 1–9. <https://doi.org/10.1111/ddi.12144>.
- 240 Elith, Jane, Catherine H. Graham, Robert P. Anderson, Miroslav Dudík, Simon Ferrier, Antoine
241 Guisan, Robert J. Hijmans, et al. 2006. "Novel Methods Improve Prediction of Species'
242 Distributions from Occurrence Data." *Ecography* 29 (2): 129–51. <https://doi.org/10.1111/j.2006.0906-7590.04596.x>.
- 244 Fick, Stephen E., and Robert J. Hijmans. 2017. "WorldClim 2: New 1-Km Spatial Resolution
245 Climate Surfaces for Global Land Areas." *International Journal of Climatology* 37 (12):
246 4302–15. <https://doi.org/10.1002/joc.5086>.
- 247 Franklin, Janet. 2010. "Moving Beyond Static Species Distribution Models in Support of
248 Conservation Biogeography: Moving Beyond Static Species Distribution Models." *Diversity
249 and Distributions* 16 (3): 321–30. <https://doi.org/10.1111/j.1472-4642.2010.00641.x>.
- 250 Hijmans, Robert J., Steven Phillips, John Leathwick, and Jane Elith. 2017. *Dismo: Species
251 Distribution Modeling*. <https://CRAN.R-project.org/package=dismo>.
- 252 Johnston, A., W. M. Hochachka, M. E. Strimas-Mackey, V. Ruiz Gutierrez, O. J. Robinson,
253 E. T. Miller, T. Auer, S. T. Kelling, and D. Fink. 2019. "Best Practices for Making Reli-
254 able Inferences from Citizen Science Data: Case Study Using eBird to Estimate Species
255 Distributions." *bioRxiv*, March, 574392. <https://doi.org/10.1101/574392>.
- 256 Legendre, Pierre. 2019. "A Temporal Beta-Diversity Index to Identify Sites That Have
257 Changed in Exceptional Ways in Space-Time Surveys." *Ecology and Evolution* 9 (6): 3500–

- 258 3514. <https://doi.org/10.1002/ece3.4984>.
- 259 Legendre, Pierre, Daniel Borcard, and Pedro R. Peres-Neto. 2005. "Analyzing Beta Diver-
260 sity: Partitioning the Spatial Variation of Community Composition Data." *Ecological
261 Monographs* 75 (4): 435–50. <https://doi.org/10.1890/05-0549>.
- 262 Legendre, Pierre, and Richard Condit. 2019. "Spatial and Temporal Analysis of Beta Diversity
263 in the Barro Colorado Island Forest Dynamics Plot, Panama." *Forest Ecosystems* 6 (1): 7.
264 <https://doi.org/10.1186/s40663-019-0164-4>.
- 265 Legendre, Pierre, and Miquel De Cáceres. 2013. "Beta Diversity as the Variance of Community
266 Data: Dissimilarity Coefficients and Partitioning." *Ecology Letters* 16 (8): 951–63. <https://doi.org/10.1111/ele.12141>.
- 268 Nix, Henry A. 1986. "A Biogeographic Analysis of Australian Elapid Snakes." *Atlas of Elapid
269 Snakes of Australia* 7: 4–15.
- 270 Phillips, Steven J., Robert P. Anderson, and Robert E. Schapire. 2006. "Maximum Entropy
271 Modeling of Species Geographic Distributions." *Ecological Modelling* 190 (3): 231–59.
272 <https://doi.org/10.1016/j.ecolmodel.2005.03.026>.
- 273 Poisot, Timothée, Richard LaBrie, Erin Larson, Anastasia Rahlin, and Benno I. Simmons. 2019.
274 "Data-Based, Synthesis-Driven: Setting the Agenda for Computational Ecology." *Ideas in
275 Ecology and Evolution* 12 (July). <https://doi.org/10.24908/iee.2019.12.2.e>.
- 276 Pollock, Laura J., Reid Tingley, William K. Morris, Nick Golding, Robert B. O'Hara, Kirsten
277 M. Parris, Peter A. Vesk, and Michael A. McCarthy. 2014. "Understanding Co-Occurrence
278 by Modelling Species Simultaneously with a Joint Species Distribution Model (JSDM)." *Methods in Ecology and Evolution* 5 (5): 397–406. <https://doi.org/10.1111/2041-210X.12180>.
- 281 Sullivan, Brian L., Christopher L. Wood, Marshall J. Iliff, Rick E. Bonney, Daniel Fink, and
282 Steve Kelling. 2009. "eBird: A Citizen-Based Bird Observation Network in the Biological
283 Sciences." *Biological Conservation* 142 (10): 2282–92. <https://doi.org/10.1016/j.biocon.2009.05.006>.