

1 sampbias, a method for quantifying geographic sampling  
2 biases in species distribution data

3 Alexander Zizka<sup>1,2</sup>, Alexandre Antonelli<sup>3,4,5</sup>, Daniele Silvestro<sup>3,4,6</sup>

- 4 1. German Centre for Integrative Biodiversity Research Halle-Jena-Leipzig (iDiv), Univer-  
5 sity of Leipzig, Deutscher Platz 5e, 04103 Leipzig, Germany
- 6 2. Naturalis Biodiversity Center, Leiden University, Leiden, Darwinweg 2, 2333 CR Leiden  
7 The Netherlands
- 8 3. Gothenburg Global Biodiversity Centre, University of Gothenburg, Box 461, 405 30  
9 Gothenburg, Sweden
- 10 4. Department for Biological and Environmental Sciences, University of Gothenburg, Box  
11 461, 405 30 Gothenburg, Sweden
- 12 5. Royal Botanic Gardens Kew, TW9 3AE, Richmond, Surrey, United Kingdom
- 13 6. Department of Biology, University of Fribourg, Ch. du Musée 10, 1700 Fribourg,  
14 Switzerland

## 15 Abstract

16 Geo-referenced species occurrences from public databases have become essential to biodiversity  
17 research and conservation. However, geographical biases are widely recognized as a factor  
18 limiting the usefulness of such data for understanding species diversity and distribution. In  
19 particular, differences in sampling intensity across a landscape due to differences in human  
20 accessibility are ubiquitous but may differ in strength among taxonomic groups and datasets.  
21 Although several factors have been described to influence human access (such as presence of  
22 roads, rivers, airports and cities), quantifying their specific and combined effects on recorded  
23 occurrence data remains challenging. Here we present *sampbias*, an algorithm and software  
24 for quantifying the effect of accessibility biases in species occurrence datasets. *Sampbias* uses  
25 a Bayesian approach to estimate how sampling rates vary as a function of proximity to one  
26 or multiple bias factors. The results are comparable among bias factors and datasets. We  
27 demonstrate the use of *sampbias* on a dataset of mammal occurrences from the island of  
28 Borneo, showing a high biasing effect of cities and a moderate effect of roads and airports.  
29 *Sampbias* is implemented as a well-documented, open-access and user-friendly R package  
30 that we hope will become a standard tool for anyone working with species occurrences in  
31 ecology, evolution, conservation and related fields.

## 32 Keywords

33 Collection effort, Global biodiversity Information Facility (GBIF), Presence only data, Road-  
34 side bias, Sampling intensity

## 35 Background

36 Publicly available datasets of geo-referenced species occurrences, such as provided by the  
37 Global Biodiversity Information Facility ([www.gbif.org](http://www.gbif.org)) have become a fundamental resource  
38 in biological sciences, especially in biogeography, conservation, and macroecology. However,  
39 these datasets are typically not collected systematically and rarely include information on  
40 collection effort. Instead, they are often compiled from a variety of sources (e.g. scientific  
41 expeditions, census counts, genetic barcoding studies, and citizen-science observations).  
42 Species occurrences are therefore often subject to multiple sampling biases (Meyer et al.  
43 2016).

44 Sampling biases that may affect the recording of species occurrences (presence, absence and  
45 abundance, Isaac and Pockock 2015, Boakes et al. 2010) include the under-sampling of specific  
46 taxa (“taxonomic bias”, e.g., birds *vs.* nematodes), specific geographic regions (“geographic  
47 bias”, i.e. easily accessible *vs.* remote areas), and specific temporal periods (“temporal bias”,  
48 i.e. wet season *vs.* dry season). In particular geographic sampling bias—the fact that sampling  
49 effort is spatially biased, rather than equally distributed over the study area—is likely to be  
50 widespread in all non-systematically collected datasets of species distributions.

51 Many aspects can lead to sampling biases, including socio-economic factors (i.e. national  
52 research spending, history of scientific research; [www.bio-dem.surge.sh](http://www.bio-dem.surge.sh), Meyer et al. 2015,  
53 Daru et al. 2018), political factors (armed conflict, democratic rights; Rydén et al. 2019),  
54 and physical accessibility (i.e. distance to a road or river, terrain conditions, slope; Yang  
55 et al. 2014, Botts et al. 2011). Especially physical accessibility by people is omnipresent

56 as a bias factor (e.g. Lin et al. 2015, Kadmon et al. 2004, Engemann et al. 2015), across  
57 spatial scales, as the commonly used term “roadside bias” testifies. In practice, this means  
58 that most species observations are made in or near cities, along roads, paths, and rivers, and  
59 near human settlements. Relatively fewer observations are expected to be available from  
60 inaccessible areas in e.g. a tropical rainforest or a mountain top. Since the recording of  
61 different taxonomic groups poses different challenges, geographic sampling bias and the effect  
62 of accessibility may differ among taxonomic groups (Vale and Jenkins 2012).

63 The implications of not considering geographic sampling biases in biodiversity research are  
64 likely to be substantial (Rocchini et al. 2011, Barbosa et al. 2013, Yang et al. 2013,  
65 Kramer-Schadt et al. 2013, Shimadzu and Darnell 2015, Meyer et al. 2016). The presence of  
66 geographic sampling biases is broadly recognized (e.g. Kadmon et al. 2004), and approaches  
67 exist to account for it in some analyses—such as for species-richness estimates (Engemann  
68 et al. 2015) species distribution models (Beck et al. 2014, Varela et al. 2014, Warren et al.  
69 2014, Boria et al. 2014, Fourcade et al. 2014, Fithian et al. 2015, Stolar and Nielsen 2015,  
70 Monsarrat et al. 2019), occupancy models (Kery and Royle 2016), and abundance estimates  
71 (Shimadzu and Darnell 2015). In contrast, few attempts have been made to explicitly quantify  
72 the overall bias (Hijmans et al. 2000, Kadmon et al. 2004) or to discern and quantify different  
73 sources of bias (Fithian et al. 2015, Fernández and Nakamura 2015, Ruete 2015). To our  
74 knowledge, no tools exist for comparing the strength of bias factors or datasets. We define as  
75 *bias factors* any anthropogenic or natural features that facilitate human access and sampling,  
76 such as roads, rivers, airports, and cities.

77 It is unrealistic to expect that accessibility bias in biodiversity data will ever disappear even

78 after more automated observation technologies are developed. It is therefore crucial that  
79 researchers realise the intrinsic biases associated with the data they deal with. This is the  
80 first step towards estimating to which extent these biases may affect their analyses, results,  
81 and conclusions. Any study dealing with species occurrence data should arguably assess the  
82 strength of accessibility biases in the underlying data. Such a quantification can also help  
83 researchers to target further sampling efforts.

84 Here, we present *sampbias*, a probabilistic method to quantify accessibility bias in datasets  
85 of species occurrences. *Sampbias* is implemented as a user-friendly R-package and uses a  
86 Bayesian approach to address three questions:

- 87 1) How strong is the accessibility bias in a given dataset?
- 88 2) How strong is the effect of different bias factors in causing the overall accessibility bias?
- 89 3) How is accessibility bias distributed in space, i.e. which areas are a priority for targeted  
90 sampling?

91 *Sampbias* is implemented in R (R Core Team 2019), based on commonly used packages for  
92 data handling (*ggplot*, Wickham 2009, *forcats*, 2019, *tidyr*, Wickham and Henry 2019,  
93 *dplyr*, Wickham et al. 2019, *magrittr*, Bache and Wickham 2014, *viridis*, Garnier 2018),  
94 handling geographic information and geo-computation (*raster*, Hijmans 2019, *sp*, Pebesma  
95 and Bivand 2005, Bivand et al. 2013) and statistical modelling (*stats*, R Core Team 2019).  
96 *Sampbias* offers an easy and largely automated means for biodiversity scientists and non-  
97 specialists alike to explore bias in species occurrence data, in a way that is comparable across

98 datasets. The results may be used to identify priorities for further collection or digitalization  
99 efforts, improve species distribution models (by providing bias surfaces in the analyses), or  
100 assess the reliability of scientific results based on publicly available species distribution data.

## 101 **Methods and Features**

### 102 **General concept**

103 Under the assumption that organisms exist across the entire area of interest, we can expect the  
104 number of sampled occurrences in a restricted area, such as a single biome, to be distributed  
105 uniformly in space (even though, of course, the density of individuals and the species diversity  
106 may be heterogeneous). With *sampbias* we assess to which extent variation in sampling rates  
107 can be explained by distance from bias factors.

108 *Sampbias* works at a user-defined spatial scale, and any dataset of multi-species occurrence  
109 records can be tested against any geographic gazetteer. Reliability increases with increasing  
110 dataset size. Default global gazetteers for airports, cities, rivers and roads are provided  
111 with *sampbias*, and user-defined gazetteers can be added easily. Species occurrence data as  
112 downloaded from the data portal of GBIF can be directly used as input data for *sampbias*.  
113 The output of the package includes measures of the sampling rates across space, which are  
114 comparable between different gazetteers (e.g. comparing the biasing effect of roads and rivers),  
115 different taxa (e.g. birds *vs.* flowering plants) and different data sets (e.g. specimens *vs.*  
116 human observations).

## 117 **Distance calculation**

118 *Sampbias* uses gazetteers of the geographic location of bias factors (hereafter indicated with  
119 B) to generate a regular grid across the study area (the geographic extent of the dataset).  
120 For each grid cell  $i$ , we then compute a vector  $X_i(j)$  of minimum distances (straight aerial  
121 distance, “as the crow flies”) to each bias factor  $j \in B$ . The resolution of the grid defines the  
122 precision of the distance estimates, for instance a 1x1 degree raster will yield approximately  
123 a 110 km precision at the equator. Due to the assumption of homogeneous sampling and a  
124 computational trade-off between the resolution of the regular grid and the extent of the study  
125 area (for instance, a 1 second resolution for a global dataset would become computationally  
126 prohibitive in most practical cases), *sampbias* is best suited for local or regional datasets at  
127 high resolution (c. 100 – 10,000 m).

## 128 **Quantifying accessibility bias using a Bayesian framework**

129 We describe the observed number of sampled occurrences  $S_i$  within each cell  $i$  as the result of  
130 a Poisson sampling process with rate  $\lambda_i$ . We model the rate  $\lambda_i$  as a function of a parameter  
131  $q$ , which represents the expected number of occurrences per cell in the absence of biases,  
132 i.e. when  $\sum_{j=1}^B X_i(j) = 0$ . Additionally, we model  $\lambda_i$  to decrease exponentially as a function  
133 of distance from bias factors, such that increasing distances will result in a lower sampling  
134 rate. For a single bias factor the rates of cell  $i$  with distance  $X_i$  from a bias is:

$$\lambda_i = q \times \exp(-wX_i)$$

135 where  $w \in \mathbb{R}^+$  defines the steepness of the Poisson rate decline, such that  $w \approx 0$  results in a  
136 null model of uniform sampling rate  $q$  across cells. In the presence of multiple bias factors  
137 (e.g. roads and rivers), the sampling rate decrease is a function of the cumulative effects of  
138 each bias and its distance from the cell:

$$\lambda_i = q \times \exp \left( - \sum_{j=1}^B w_j X_i(j) \right) \quad (1)$$

139 where a vector  $\mathbf{w} = [w_1, \dots, w_B]$  describes the amount of bias attributed to each specific factor.

140 To quantify the amount of bias associated with each factor, we jointly estimate the parameters  
141  $q$  and  $\mathbf{w}$  in a Bayesian framework. We use Markov Chain Monte Carlo (MCMC) to sample  
142 these parameters from their posterior distribution:

$$P(q, \mathbf{w} | \mathbf{S}) \propto \prod_{i=1}^N Poi(S_i | \lambda_i) \times P(q) P(\mathbf{w}) \quad (2)$$

143 where the likelihood of sampled occurrences  $S_i$  within each cell  $Poi(S_i | \lambda_i)$  is the probability  
144 mass function of a Poisson distribution with rate per cell defined as in Eqn. (1). The  
145 likelihood is then multiplied across the  $N$  cells considered. We used exponential priors on  
146 the parameters  $q$  and  $\mathbf{w}$ ,  $P(q) \sim \Gamma(1, 0.01)$  and  $P(\mathbf{w}) \sim \Gamma(1, 1)$ , respectively.

147 We summarize the parameters by computing the mean of the posterior samples and their  
148 standard deviation. We interpret the magnitude of the elements in  $\mathbf{w}$  as a function of the  
149 importance of the individual biases. We note, however, that this test is not explicitly intended  
150 to assess the significance of each bias factor (for which a Bayesian variable selection method



151 could be used), particularly since several bias factors might be correlated (e.g. cities, and  
152 airports). Instead, these analyses can be used to quantify the expected amount of bias in the  
153 data that can be predicted by single or multiple predictors in order to identify under-sampled  
154 and unexplored areas.

155 We summarize the results by mapping the estimated sampling rates ( $\lambda_i$ ) across space. These  
156 rates represent the expected number of sampled occurrences for each grid cell and provide a  
157 graphical representation of the spatial variation of sampling rates. Provided that the cells are  
158 of equal size, the estimated rates will be comparable across data sets, regions, and taxonomic  
159 groups. Analysing different regions, biomes, or taxa in separate analyses allows to account  
160 for differences in over sampling rates, which are not linked with bias factors. For instance,  
161 the unbiased sampling rate  $q$  is expected to differ between a highly sampled clade like birds  
162 and under-sampled groups of invertebrates, but their sampling biases ( $\mathbf{w}$ ) might be similar  
163 across the two groups.

## 164 **Example and Empirical validation**

165 A default *sampbias* analysis can be run with few lines of code in R. The main function  
166 `calculate_bias` creates an object of the class "`sampbias`", for which the package provides  
167 a plotting and summary method. Based on a `data.frame` including species identity and  
168 geographic coordinates. Additional options exist to provide custom gazetteers, a custom grain  
169 size of the analysis, as well as some operators for the calculation of the bias distances. A  
170 tutorial on how to use *sampbias* is available with the package and in the electronic supplement  
171 of this publication (Appendix S1).

172 To exemplify the use and output of *sampbias*, we downloaded the occurrence records of all  
173 mammals available from the island of Borneo (n = 6,262, GBIF.org 2016), and ran *sampbias*  
174 using the default gazetteers as shown in the example code below, to test the biasing effect  
175 of the main airports, cities and roads in the dataset. The example dataset is provided with  
176 *sampbias*.

177 We found a strong effect of cities on sampling intensity, a moderate effect of roads and airports  
178 and negligible effect of rivers (Fig. 1). All models predict a low number of collection records  
179 in the centre of Borneo (Fig. 2), which reflects the original data, and where accessibility  
180 means are low (Figure S1 in Appendix S2). The empirical example illustrates the use of  
181 *sampbias*, for detailed analyses or a smaller geographic scale, higher resolution gazetteers,  
182 including smaller roads and rivers and a higher spatial resolution would be desirable. Results  
183 might change with increasing resolution, since roads and rivers might have a stronger effect  
184 on higher resolutions (facilitating most the access to their immediate vicinity), whereas cities  
185 and airports might have a stronger effect on the larger scale (facilitating access to a larger  
186 area).

```
library(sampbias)

#a data table with species identify, longitude, and latitude

example.in <- read.csv(system.file("extdata",
                                  "mammals_borneo.csv",
                                  package="sampbias"),
                      sep = "\t")
```

```
#running sampbias

example.out <- calculate_bias(x = example.in,
                             res = 0.05,
                             buffer = 0.5)

#summary

summary(example.out)

plot(example.out)

#projecting the bias effect in space

proj <- project_bias(example.out)

map_bias(proj)
```

## 187 Data accessibility

188 *Sampbias* is available under a GNU General Public license v3 from <https://github.com/azi>  
189 [zka/sampbias](https://github.com/azi/zka/sampbias), and includes the example dataset as well as a tutorial (Appendix S1) and a  
190 summary of possible warnings produced by the package (Appendix S3).

## 191 **Acknowledgements**

192 We thank the organizers of the 2016 Ebben Nielsen challenge for inspiring and recognizing  
193 this research. We thank all data collectors and contributors to GBIF for their effort. AZ is  
194 thankful for funding by iDiv via the German Research Foundation (DFG FZT 118), specifically  
195 through sDiv, the Synthesis Centre of iDiv. AA is supported by grants from the Swedish  
196 Research Council, the Knut and Alice Wallenberg Foundation, the Swedish Foundation for  
197 Strategic Research and the Royal Botanic Gardens, Kew. DS received funding from the  
198 Swedish Research Council (2015-04748) and from the Swiss National Science Foundation  
199 (PCEFP3\_187012)

## 200 **Author contributions**

201 All authors conceived this study, AZ and DS developed the statistical algorithm and wrote  
202 the R-package, AZ and DS wrote the manuscript with contributions from AA.

203 **Figures**

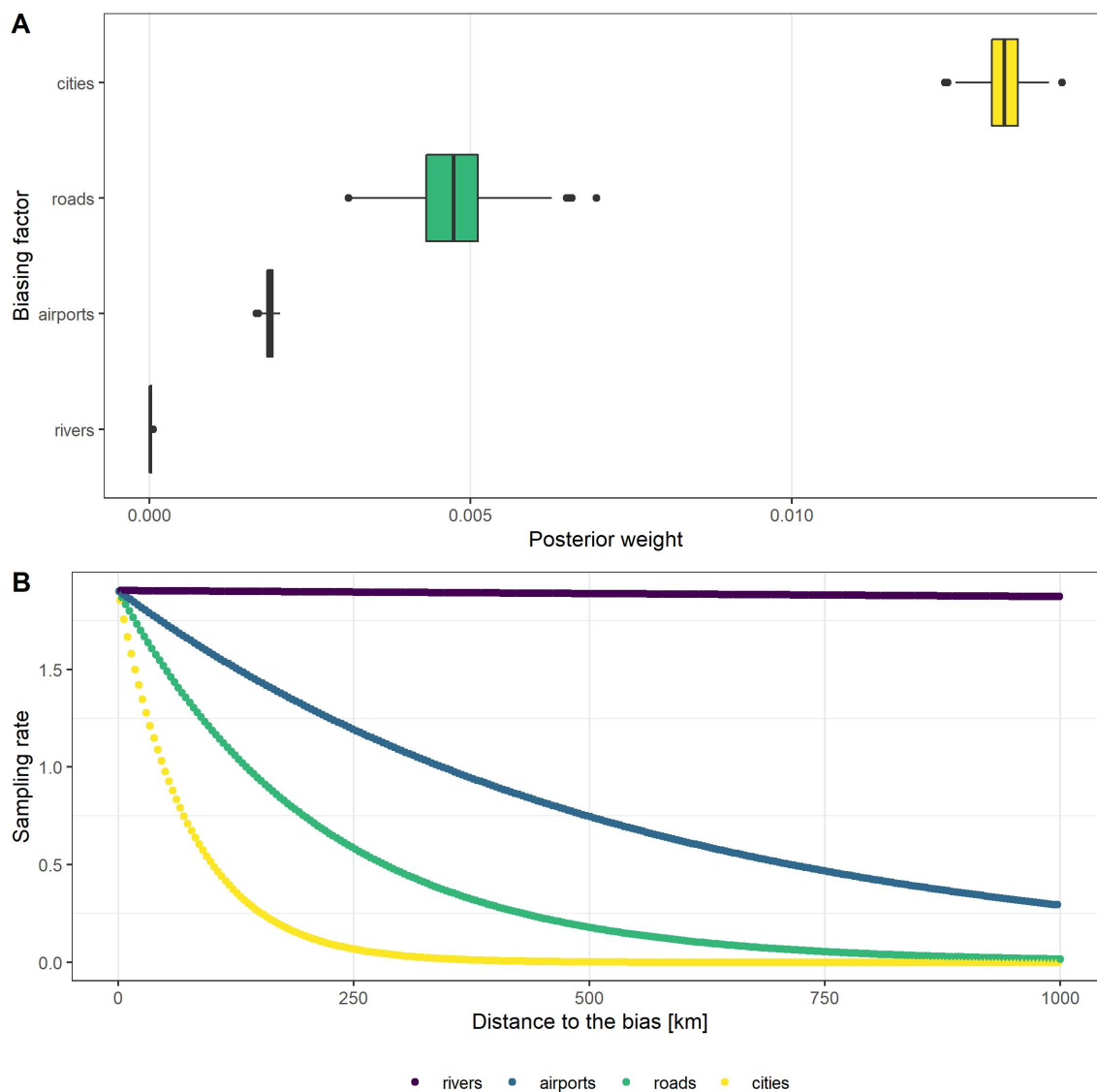


Figure 1: Results of the empirical validation analysis, estimating the accessibility bias in mammal occurrences from Borneo). A) bias weights ( $w$ ) defining the effects of each bias factor, B) sampling rate as function of distance to the closest instance of each bias factor (i.e. expected number of occurrences) given the inferred *sampbias* model. At the study scale of 0.05 degrees (c. 5km) *sampbias* finds the strongest biasing effect for the proximity of cities and roads.

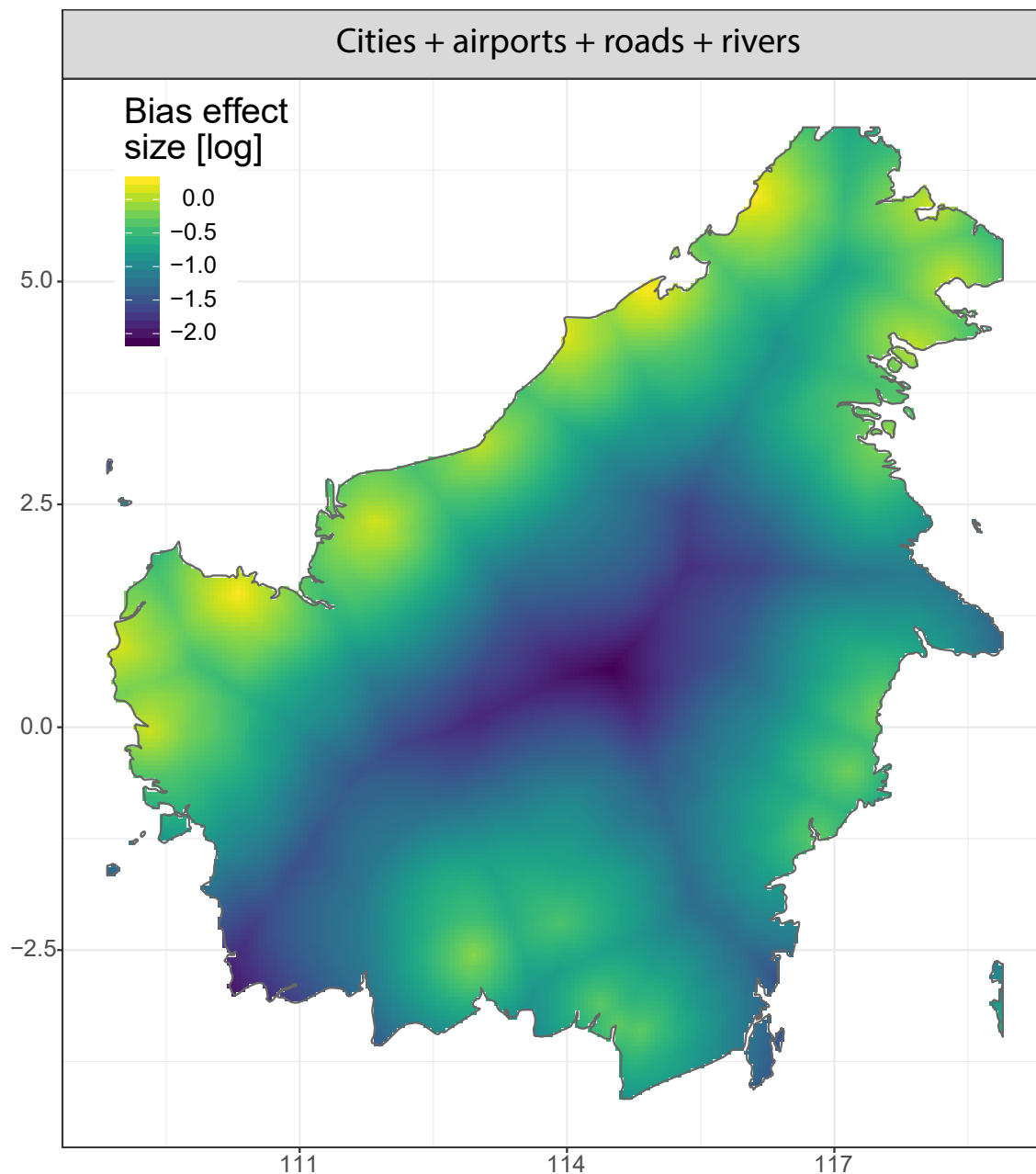


Figure 2: Spatial projection of the estimated sampling rates in an empirical example dataset of mammal occurrences on the Indonesian island of Borneo (downloaded from [www.gbif.org](http://www.gbif.org). GBIF.org, 2016). The colours show the projection of the sampling rates (i.e. expected number of occurrences per cell) given the inferred extitsampbias model. The highest undersampling is in the centre of the island.

204 **Supplementary material**

205 Appendix S1 - Tutorial running *sampbias* in R

206 Appendix S2 - Supplementary Figure S1

207 Appendix S3 - Possible warnings and their solutions

## 208 **References**

- 209 Bache, S. M. and Wickham, H. 2014. magrittr: A Forward-Pipe Operator for R.
- 210 Barbosa, A. M. et al. 2013. Species-people correlations and the need to account for survey  
211 effort in biodiversity analyses. - *Diversity and Distributions* 19: 1188–1197.
- 212 Beck, J. et al. 2014. Spatial bias in the GBIF database and its effect on modeling species'  
213 geographic distributions. - *Ecological Informatics* 19: 10–15.
- 214 Bivand, R. S. et al. 2013. *Applied spatial data analysis with R*, Second edition. - Springer.
- 215 Boakes, E. H. et al. 2010. Distorted views of biodiversity: Spatial and temporal bias in  
216 species occurrence data. - *PLoS Biology* 8: e1000385.
- 217 Boria, R. A. et al. 2014. Spatial filtering to reduce sampling bias can improve the performance  
218 of ecological niche models. - *Ecological Modelling* 275: 73–77.
- 219 Botts, E. A. et al. 2011. Geographic sampling bias in the South African Frog Atlas Project:  
220 Implications for conservation planning. - *Biodiversity and Conservation* 20: 119–139.
- 221 Daru, B. H. et al. 2018. Widespread sampling biases in herbaria revealed from large-scale  
222 digitization. - *New Phytologist* 217: 939–955.
- 223 Engemann, K. et al. 2015. Limited sampling hampers “big data” estimation of species  
224 richness in a tropical biodiversity hotspot. - *Ecology and Evolution* 5: 807–820.
- 225 Fernández, D. and Nakamura, M. 2015. Estimation of spatial sampling effort based on



- 226 presence-only data and accessibility. - *Ecological Modelling* 299: 147–155.
- 227 Fithian, W. et al. 2015. Bias correction in species distribution models: pooling survey and  
228 collection data for multiple species. - *Methods in Ecology and Evolution* 6: 424–438.
- 229 Fourcade, Y. et al. 2014. Mapping species distributions with MAXENT using a geographically  
230 biased sample of presence data: A performance assessment of methods for correcting sampling  
231 bias. - *PLoS ONE* 9: e97122.
- 232 Garnier, S. 2018. viridis: Default color maps from 'matplotlib'.
- 233 GBIF.org 2016. (08 September 2016) GBIF occurrence download, [doi.org/10.15468/dl.7fg4zx](https://doi.org/10.15468/dl.7fg4zx).
- 234 Hijmans, R. J. 2019. geosphere: Spherical Trigonometry.
- 235 Hijmans, R. et al. 2000. Assessing the geographic representativeness of Genbank collections:  
236 The case of Bolivian wild potatoes. - *Conservation Biology* 14: 1755–1765.
- 237 Isaac, N. J. B. and Pocock, M. J. O. 2015. Bias and information in biological records. -  
238 *Biological Journal of the Linnean Society* 115: 522–531.
- 239 Kadmon, R. et al. 2004. Effect of roadside bias on the accuracy of predictive maps produced  
240 by bioclimatic models. - *Ecological Applications* 14: 401–413.
- 241 Kery, M. and Royle, J. A. 2016. Applied hierarchical modeling in ecology - Analysis of  
242 distribution, abundance and species richness in R and BUGS: Volume 1: Prelude and Static  
243 Models. - Academic Press, Elsevier.

- 244 Kramer-Schadt, S. et al. 2013. The importance of correcting for sampling bias in MaxEnt  
245 species distribution models. - *Diversity and Distributions* 19: 1366–1379.
- 246 Lin, Y.-p. et al. 2015. Uncertainty analysis of crowd-sourced and professionally collected  
247 field data used in species distribution models of Taiwanese moths. - *Biological Conservation*  
248 181: 102–110.
- 249 Meyer, C. et al. 2015. Global priorities for an effective information basis of biodiversity  
250 distributions. - *Nature Communications* 6: 8221.
- 251 Meyer, C. et al. 2016. Multidimensional biases, gaps and uncertainties in global plant  
252 occurrence information. - *Ecology Letters* 19: 992–1006.
- 253 Monsarrat, S. et al. 2019. Accessibility maps as a tool to predict sampling bias in historical  
254 biodiversity occurrence records. - *Ecography* 42: 125–136.
- 255 Pebesma, E. J. and Bivand, R. S. 2005. Classes and methods for spatial Data: the sp Package.  
256 - *R News* 5: 21–41.
- 257 R Core Team 2019. R: A language and environment for statistical computing.
- 258 Rocchini, D. et al. 2011. Accounting for uncertainty when mapping species distributions:  
259 The need for maps of ignorance. - *Progress in Physical Geography: Earth and Environment*  
260 35: 211–226.
- 261 Ruete, A. 2015. Displaying bias in sampling effort of data accessed from biodiversity databases  
262 using ignorance maps. - *Biodiversity Data Journal* 3: e5361.

- 263 Rydén, O. et al. 2019. Linking democracy and biodiversity conservation: Empirical evidence  
264 and research gaps. - *Ambio*: 15pp.
- 265 Shimadzu, H. and Darnell, R. 2015. Attenuation of species abundance distributions by  
266 sampling. - *Royal Society Open Science* 2: 140219.
- 267 Stolar, J. and Nielsen, S. E. 2015. Accounting for spatially biased sampling effort in presence-  
268 only species distribution modelling. - *Diversity and Distributions* 21: 595–608.
- 269 Vale, M. M. and Jenkins, C. N. 2012. Across-taxa incongruence in patterns of collecting bias.  
270 - *Journal of Biogeography* 39: 1744–1744.
- 271 Varela, S. et al. 2014. Environmental filters reduce the effects of sampling bias and improve  
272 predictions of ecological niche models. - *Ecography*: 1084–1091.
- 273 Warren, D. L. et al. 2014. Incorporating model complexity and spatial sampling bias into  
274 ecological niche models of climate change risks faced by 90 California vertebrate species of  
275 concern. - *Diversity and Distributions* 20: 334–343.
- 276 Wickham, H. 2009. *ggplot2 - Elegant graphics for data analysis*. - Springer.
- 277 Wickham, H. 2019. *forcats: Tools for working with categorical variables (Factors)*.
- 278 Wickham, H. and Henry, L. 2019. *tidyr: Tidy messy data*.
- 279 Wickham, H. et al. 2019. *dplyr: A grammar of data manipulation*.
- 280 Yang, W. et al. 2013. Geographical sampling bias in a large distributional database and its

281 effects on species richness-environment models. - *Journal of Biogeography* 40: 1415–1426.

282 Yang, W. et al. 2014. Environmental and socio-economic factors shaping the geography of

283 floristic collections in China. - *Global Ecology and Biogeography* 23: 1284–1292.