

1 **Title**

2 A cautionary note on the use of hypervolume kernel density estimators in
3 ecological niche modeling

4

5 **Authors**

6 Huijie Qiao^{1*}, Luis E. Escobar^{2,3*}, Erin E. Saupe⁴, Liqiang Ji¹ and Jorge
7 Soberón⁵

8

9 **Affiliations**

10 ¹*Key Laboratory of Animal Ecology and Conservation Biology, Institute of*
11 *Zoology, Chinese Academy of Sciences, Beijing, 100101 China.*

12 ²*Department of Veterinary Population Medicine, University of Minnesota, St.*
13 *Paul, Minnesota, 55108 USA.*

14 ³*Minnesota Aquatic Invasive Species Research Center, University of*
15 *Minnesota, St. Paul, Minnesota, 55108 USA;*

16 ⁴*Department of Geology & Geophysics, Yale University, New Haven,*
17 *Connecticut, 06511 USA.*

18 ⁵*Biodiversity Institute and Department of Ecology and Evolutionary Biology,*
19 *University of Kansas, Lawrence, Kansas, 66045, USA.*

20

21 *Corresponding author: Luis E. Escobar, Phone: (+1) 315-313-8584. Email:

22 lescobar@umn.edu. Address: Veterinary Diagnostic Laboratory, University of

23 Minnesota, 1365 Gortner Avenue, St. Paul, MN 55108.

24

25 **Running title**

26 Kernel density methods in ecological niche modeling

27

28 **Abstract**

29 Blonder *et al.* (2014) introduced a new multivariate kernel density estimation
30 (KDE) method to infer Hutchinsonian hypervolumes to model ecological niches.

31 The authors argued their KDE method matches or outperforms several
32 methods for estimating hypervolume geometries and for conducting species
33 distribution modeling. Further clarification, however, is appropriate with respect
34 to the assumptions and limitations of KDE as a method for species distribution
35 modeling. Using virtual species and controlled environmental scenarios, we
36 show that KDE both under- and over-estimates niche volumes depending on
37 the dimensionality of the dataset and the number of occurrence records
38 considered. We suggest KDE may be a viable approach when dealing with
39 large sample sizes, limited sampling bias, and only a few environmental
40 dimensions.

41

42 **Key words**

43 Hutchinsonian hypervolumes; ecological space; multivariate kernel density
44 estimation; virtual species; minimum volume ellipsoid; niche

45

46 **Introduction**

47 In a recent contribution, Blonder *et al.* (2014) introduced a new hypervolume
48 multivariate kernel density estimation (KDE) method to delineate
49 Hutchinsonian hypervolumes (Hutchinson, 1957, 1978) in high-dimensional
50 ecological space. A hypervolume, in this formulation, is defined by a set of
51 points within an n -dimensional environmental or ecological space that reflects
52 suitable values of these n variables. According to the authors, KDE
53 outperforms several methods for estimating hypervolume geometries and for
54 conducting species distribution modeling (SDM).

55 Blonder *et al.* (2014) argued KDE is useful for fitting observed occurrences
56 to environmental values, and for recognizing clusters or holes in occurrence
57 datasets within environmental space. Here, we show that KDE recognizes
58 clusters or holes in occurrence datasets only when occurrence data are
59 numerous and when the dimensionality of the environmental space is not too
60 large. Indeed, the KDE method may have difficulty identifying holes, gaps,
61 and/or clusters in environmental space with limited occurrences, since in this
62 case the method tends to produce broad niche estimates that smooth out

63 these clusters and holes. Caution is warranted when applying KDE in
64 high-dimensional space because of the curse of dimensionality (Hastie *et al.*,
65 2009, sections 2.5 and 6.3) and the empty space phenomenon (Silverman,
66 1986, section 4.5). That is to say, as dimensionality increases, the number of
67 samples required to estimate a shape with accuracy will also increase
68 dramatically.

69 In situations where KDE is able to recognize correctly clusters or holes in
70 occurrence datasets, we argue that doing so is useful only to the extent that
71 the realized niche (RN) is sought and not the fundamental niche (FN). Blonder
72 and colleagues indicated their KDE method estimates ‘holey’ Hutchinsonian
73 hypervolumes without *a priori* reason to assume that a hypervolume (or niche)
74 should be normally or uniformly distributed in multiple dimensions (Fig. 1e in
75 Blonder *et al.* 2014). We argue, however, that traditional Hutchinsonian
76 hypervolumes would not fit tightly to available occurrence data, especially if
77 one seeks the FN (a point also noted by Blonder and colleagues). Empirical
78 and theoretical arguments suggest the FN has a convex shape (Birch, 1953;
79 Maguire, 1973; Austin *et al.*, 1984; Colwell & Rangel, 2009; Araújo & Peterson,
80 2012; Drake, 2015), and, consequently, convex-hulls or ellipsoids (multivariate
81 Gaussian shapes) may often be the simplest proxy (Peterson *et al.*, 2011).

82 Our argument is theoretical and emphasizes choosing the appropriate
83 method for a particular application: if the RN is desired, the KDE method of
84 Blonder *et al.* (2014) may be a good candidate, assuming low occurrence

85 density and high dimensionality does not prevent its practical application
86 (Franklin, 2005; Hastie *et al.*, 2009, sections 2.5 and 6.3). However, if the FN is
87 to be estimated, the KDE method may not be ideal.

88 If the KDE method functions as Blonder and colleagues propose,
89 producing strict estimates of the environmental space occupied by a species,
90 model transferability to different regions or time periods – a common goal in
91 SDM – will be limited. For example, say available occurrences for a species
92 are distributed in temperatures of 15°, 16°, 17°, 19°, and 20°C. In this scenario,
93 ignoring potential suitability for the species at 18°C, the 'hole' in the series,
94 may be biologically unrealistic. As is likely the case in this simplistic example,
95 many environmental holes in occurrence data may be due to biases in
96 sampling, the availability of existing environmental conditions, and/or biological
97 constraints, and do not reflect real suitability requirements.

98 Based on the considerations above, we re-evaluated the experiments of
99 Blonder *et al.* (2014) using diverse FN shapes, including range-boxes (RB;
100 Birch, 1953), convex-hulls (CH; Godsoe, 2010; Qiao *et al.*, 2015), and
101 minimum-volume ellipsoids (MVE; Maguire, 1973; Qiao *et al.*, 2015), which
102 have been previously invoked and employed in ecological studies. This
103 reassessment identifies those tools that best fit with a particular and diverse
104 set of research questions, and provides users with a rich source of information
105 for selecting model approaches.

106 **Methods**

107 **1. KDE performance**

108 We illustrate the functionality of the KDE method using a virtual environmental
109 space, E , composed of 10,000 unique random observations in two dimensions.
110 Different configurations and densities of occurrences were sampled from a
111 virtual FN within this environmental space, defined as a range box (red
112 rectangle in Figs. 1, S1, S2). Note that the FN is easily estimated with virtual
113 species based on controlled occurrence data, but observed occurrences from
114 real species will most likely capture the RN and not the FN, which is
115 constrained by biotic interactions, accessibility, and the available environment
116 (Peterson *et al.* 2011). Within this virtual FN, we collected independent
117 occurrence datasets of three different sample sizes m ($m = 10, 100, \text{ and } 1000$)
118 and four different sampling configurations: (i) evenly distributed or unbiased
119 (Figs. 1.a, S1.a, S2.a), clustered or biased (Figs. 1.b, S1.b, S2.b), (ii) absent
120 from the center of the FN or “holey” (Figs. 1.c, S1.c, S2.c), and (iv) distributed
121 in two distinct environmental clusters (Figs. 1.d, S1.d, S2.d). We repeated the
122 sampling process 10 times to capture variation. Using these 120 sampling
123 datasets (i.e., three sample sizes m x four sample configurations x 10
124 replicates), we estimated the virtual FN using the KDE approach and assessed
125 the quality of these estimates based on the resulting type I (i.e., false presence,
126 or incorrect rejection of a true null hypothesis) and type II error (i.e., false
127 absence, or the failure to reject a false null hypothesis). Error was quantified as

128 the number of observations in the virtual space that were incorrectly predicted.

129

130 **2. Comparison of KDE to other algorithms**

131 We created three virtual FN configurations — RB, CH, and MVE — to explore

132 quantitatively the performance of different modeling algorithms in estimating

133 FNs. To create these virtual FNs, we first generated e uncorrelated virtual

134 environmental variables (with e taking one of four possible values, $e = 2, 4, 6,$

135 and 8), to create the environmental space, E , composed of 10,000 unique

136 random observations. Environmental values in E ranged between 0 and 1 in

137 each of the eight dimensions (Fig. S3). Next, we selected 10 random

138 observations (N) in E to define the vertices of the FNs under three shape

139 hypotheses, $N = RB, CH,$ and MVE. Environmental values used to define N

140 were constrained between 0.2 and 0.8 to avoid potential novel environmental

141 conditions (Fig. S3).

142 The environmental observations inside each of these virtual niches were

143 regarded as the species' presences. For each virtual FN (i.e., RB, CH and

144 MVE), we collected independent occurrence datasets of three sample sizes, m

145 ($m = 10, 100,$ and 1000) in e environmental dimensions ($e = 2, 4, 6,$ and 8). This

146 sampling process was repeated 10 times to generate random replicates of

147 species' occurrences, which resulted in 360 simulations from the combination

148 of three FN shape hypotheses x three sample sizes (m) x four environmental

149 dimensions (e) x 10 random replicates.

150 To model the virtual FNs, we used the methods proposed by Blonder and
151 colleagues (2014), including RB, CH, MVE, and KDE. As in Blonder *et al.*
152 (2014), KDEs were inferred using a Silverman bandwidth estimator (Silverman,
153 1986, section 4.5) and a quantile threshold of 0.5. Note that smaller
154 bandwidths (i.e., larger thresholds) will lead to smaller hypervolumes. As aptly
155 noted by Blonder and colleagues, analyses that have few observations ($m/e <$
156 10, as a rough guideline) will be sensitive to the choice of bandwidth.

157 Following Blonder *et al.* (2014), we used the volume of the niche to explore
158 the amount of E predicted by the models. We compared the volume of the
159 ‘estimated niche’ (n) with the known ‘true’ volume (N) of the virtual FN. Niche
160 size measured as volume, however, may be insensitive to type I error, such
161 that the ‘true niche’ and the ‘estimated niche’ may yield similar volumes but
162 have minimal or no environmental overlap. To avoid this problem, we
163 evaluated all models using sensitivity (Eq. S1) and specificity (Eq. S2) based
164 on omission error (Fielding & Bell, 1997), and the Jaccard index (Eq. S3)
165 based on comparisons between the known (N) and estimated (n) niche
166 volumes (Jaccard, 1912; Godsoe, 2014).

167

168 **Results**

169 **1. KDE performance**

170 The KDE method tended to overestimate the FN and extend beyond the

171 occurrence data (i.e., the RN) when using small sample sizes. The severity of
172 this overestimation, however, varied depending on the sample configuration
173 (Figs.S1, S2). KDE identified the 'hole' (black box; Figs. S1c, S2c) only under
174 the largest sample size (Fig. 1). The 'clusters' were identified with sample sizes
175 over 100 (Figs. 1, S2), but in these instances, KDE estimates extended
176 significantly beyond the FN and the RN. In general, type I error decreased and
177 type II error increased when more occurrences were used for model calibration
178 (Fig.2).

179

180 **2. Comparison of KDE to other algorithms**

181 In most cases, the KDE algorithm overestimated the volume of the true FN
182 when the shape of the niche was defined as RB (Fig. S4.a), a result congruent
183 with that of Blonder and colleagues (2014). The RB and CH algorithms
184 returned the most variable niche volume estimates, with consistent
185 underestimation of FN volumes. These two algorithms, however, obtained the
186 highest Jaccard similarity values between the estimated and observed RB FN
187 (Fig. 3.a), particularly in high-dimensional environmental space. When the
188 virtual FN was defined as CH, MVE algorithm got the highest Jaccard similarity
189 values in the low dimension ($e=2$), KDE was the winner in the middle
190 dimension ($e = 4$), and RB won in the high dimension ($e=6, 8$, Fig. 3.a). When
191 the virtual FN was defined as MVE, we failed to replicate the results of Blonder
192 and colleagues, who found that MVE consistently overestimated niche

193 volumes (Fig. 4c in Blonder *et al.*, 2014, our figure Fig. S4.c).

194 Overall, method performance varied as a function of the 'true shape' of the
195 virtual niche. That is to say, the RB method performed best when the true
196 shape was RB, and so forth. In general, CH tended to underestimate true
197 niche volumes. Similarly, MVE and RB underestimated true niche volumes
198 using small sample sizes, but overestimated niche volumes using larger
199 sample sizes. KDE tended to underestimate volumes of niches in high
200 dimensionality and overestimate volumes of niches in low dimensionality (Fig.
201 S4).

202 All methods performed well in terms of specificity and sensitivity using
203 large sample sizes ($m = 100, 1000$). Results for smaller sample sizes ($m = 10$),
204 however, were more variable. When considering sensitivity, KDE performed
205 well, as this method tends to generate broader niche estimates (Fig. S5).
206 Broader niche estimates, however, will generate more opportunities for type I
207 error, resulting in lower specificity values. Indeed, the KDE method performed
208 worst in terms of specificity for small sample sizes, whereas the CH method
209 performed best. Overall, the CH method performed well in terms of specificity,
210 but performed poorly when considering sensitivity (Fig. S6). As dimensionality
211 increased, the KDE method exhibited decreased sensitivity but increased
212 specificity, and underestimated the true volume of the niche. In other words,
213 estimates were constrained severely in high dimensions.

214

215 **Discussion**

216 Our results suggest that accuracy on niche estimations depends on the
217 research question and particularities of the data. A complex algorithm, such as
218 KDE, may function best when the goal is to fit models tightly to available data
219 and avoid environmental interpolation across 'holes' in environmental space.
220 These are often desirable features when exploring the occupied area or RN of
221 a species, or the distribution of non-living organisms (e.g., when mapping
222 potential wildfires). KDE, however, is sensitive to both sample size and
223 environmental dimensionality. Contrary to the claims of Blonder and
224 colleagues, KDE may overestimate niche volumes in low dimensions and
225 constrict niche volume estimates in high dimensions. We found that as
226 dimensionality increases, specificity increases as sensitivity decreases (Drake
227 2015; Figs. S5, S6).

228 The MVE algorithm performed best when the target shape is ellipsoid in
229 nature, which is often hypothesized to be the true shape of species' FNs
230 (Hutchinson, 1957; Maguire, 1973; Brown, 1984; Drake, 2015). The CH
231 method tended to generate narrow niche estimates relative to the KDE method,
232 as reflected in the specificity and sensitivity values. The CH algorithm may be
233 suitable when the goal is to estimate suitable environmental conditions
234 allowing environmental interpolation, but avoiding prediction of suitable
235 conditions in novel environments.

236 The analyses conducted herein support the idea that there is often not a

237 single 'best' algorithm or method that fits with all ecological applications and
238 data configurations for estimating species' niches (Guillera-Arroita *et al.*, 2015;
239 Qiao *et al.*, 2015). As is now common practice in phylogenetics, we propose
240 that the best niche model should be selected from a variety of model
241 hypotheses, based on its fit to the nature of the data and the specific research
242 question (Diniz-Filho *et al.*, 2015).

243

244 **Abbreviations**

245 **CH** Convex-hull.

246 **E** The environmental variables forming an environmental space.

247 **e** The number of environmental variables used to create the
248 environmental space.

249 **FN** Fundamental niche.

250 **KDE** Multivariate kernel density method.

251 **N** 'True' virtual niche.

252 **n** Estimation of the 'true' fundamental niche from a sample of
253 occurrences.

254 **MVE** Minimum-volume ellipsoid.

255 **m** Number of occurrences collected from the environmental space used
256 to estimate **N**, the true virtual niche.

257 **RN** Realized niche.

258 **RB** Range-box.

259 **SDM** Species distribution modeling.

260

261 **Acknowledgements**

262 We thank editors Richard Field and Antoine Guisan and two anonymous
263 reviewers for invaluable suggestions that improved this manuscript. HQ was
264 supported by National Natural Science Foundation of China (A New Method to
265 Predict the Species Distributions, 31100390). LEE was supported by the
266 Minnesota Environment and Natural Resources Trust Fund, the Minnesota
267 Aquatic Invasive Species Research Center, and the Clean Water Land and
268 Legacy. JS was partially supported by NSF grant 1208472

269

270 **References**

- 271 Araújo, M.B. & Peterson, A.T. (2012) Uses and misuses of bioclimatic
272 envelope modeling. *Ecology*, **93**, 1527-1539.
- 273 Austin, M.P., Cunningham, R.B. & Fleming, P.M. (1984) New approaches to
274 direct gradient analysis using environmental scalars and statistical
275 curve-fitting procedures. *Plant Ecology*, **55**, 11-27.
- 276 Birch, L.C. (1953) Experimental background to the study of the distribution
277 and abundance of insects: III. The relation between innate capacity for

278 increase and survival of different species of beetles living together on the
279 same food. *Evolution*, **7**, 136-144.

280 Blonder, B., Lamanna, C., Violle, C. & Enquist, B.J. (2014) The
281 n-dimensional hypervolume. *Global Ecology and Biogeography*, **23**,
282 595-609.

283 Brown, J.H. (1984) On the relationship between abundance and distribution
284 of species. *American Naturalist*, **124**, 255-279.

285 Colwell, R.K. & Rangel, T.F. (2009) Hutchinson's duality: The once and
286 future niche. *Proceedings of the National Academy of Sciences USA*, **106**,
287 19651-19658.

288 Diniz-Filho, J.A.F., Rodrigues, H., Telles, M.P.D.C., Oliveira, G.D., Terribile,
289 L.C., Soares, T.N. & Nabout, J.C. (2015) Correlation between genetic
290 diversity and environmental suitability: taking uncertainty from ecological
291 niche models into account. *Molecular Ecology Resources*, **15**, 1059-1066.

292 Drake, J.M. (2015) Range bagging: A new method for ecological niche
293 modelling from presence-only data. *Journal of The Royal Society Interface*,
294 **12**, 10.1098/rsif.2015.0086.

295 Fielding, A. & Bell, J. (1997) A review of methods for the assessment of
296 prediction errors in conservation presence/absence models. *Environmental*
297 *Conservation*, **24**, 38-49.

298 Franklin, J. (2005) The elements of statistical learning: data mining,
299 inference and prediction. *The Mathematical Intelligencer*, **27**: 83-85.

300 Godsoe, W. (2010) I can't define the niche but I know it when I see it: A
301 formal link between statistical theory and the ecological niche. *Oikos*, **119**,
302 53-60.

303 Godsoe, W. (2014) Inferring the similarity of species distributions using
304 species' distribution models. *Ecography*, **37**, 130-136.

305 Guillera-Aroita, G., Lahoz-Monfort, J.J., Elith, J., Gordon, A., Kujala, H.,
306 Lentini, P.E., McCarthy, M.A., Tingley, R. & Wintle, B.A. (2015) Is my
307 species distribution model fit for purpose? Matching data and models to
308 applications. *Global Ecology and Biogeography*, **24**, 276-292.

309 Hastie, T.R. Tibshirani, & Friedman J. (2009). *The Elements of Statistical*
310 *Learning. Data Mining, Inference and Prediction*. Second Edition. Springer
311 *New York*.

312 Silverman, B.W. (1986). *Density Estimation for Statistics and Data Analysis*.
313 Chapman and Hall. *London*.

314 Hutchinson, G.E. (1957) Concluding remarks. *Cold Spring Harbor*
315 *Symposia on Quantitative Biology*, **22**, 415-427.

316 Hutchinson, G.E. (1978) *An introduction to Population Ecology*. Yale
317 University Press, New Haven.

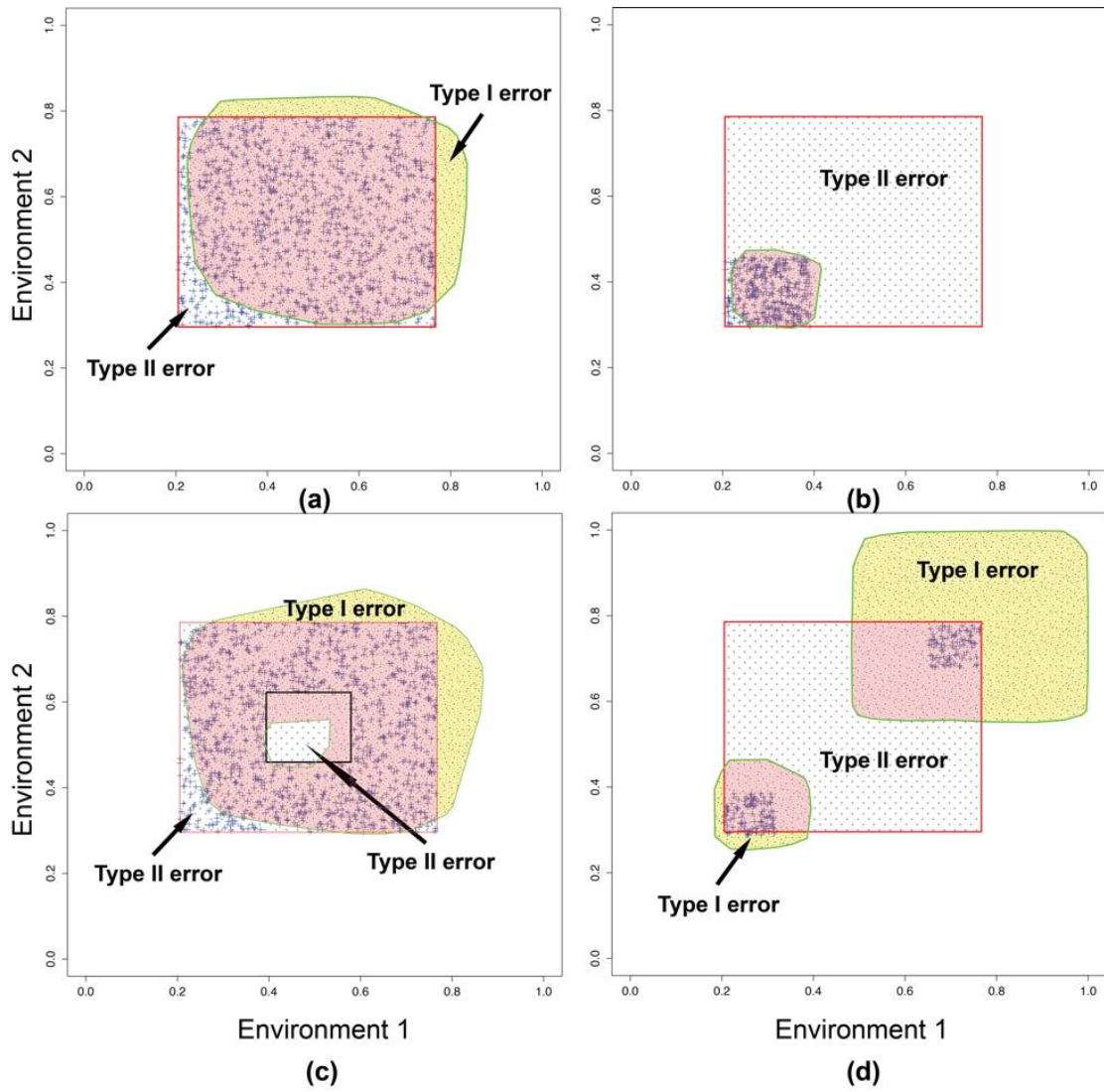
318 Jaccard, P. (1912) The distribution of the flora in the alpine zone. *New*
319 *Phytologist*, **11**, 37-50.

320 Maguire, B.Jr (1973) Niche response structure and the analytical potentials
321 of its relationship to the habitat. *American Naturalist*, **107**, 213-246.

322 Peterson, A.T., Soberón, J., Pearson, R.G., Anderson, R.P.,
323 Martínez-Meyer, E., Nakamura, M. & Araújo, M.B. (2011) *Ecological Niches*
324 *and Geographic Distributions*. Princeton University, Princeton and Oxford.
325 Qiao, H., Soberón, J. & Peterson, T.A. (2015) No silver bullets in correlative
326 ecological niche modeling: insights from testing among many potential
327 algorithms for niche estimation. *Methods in Ecology and Evolution*, **6**,
328 1126-1136.

329 **Biosketch**

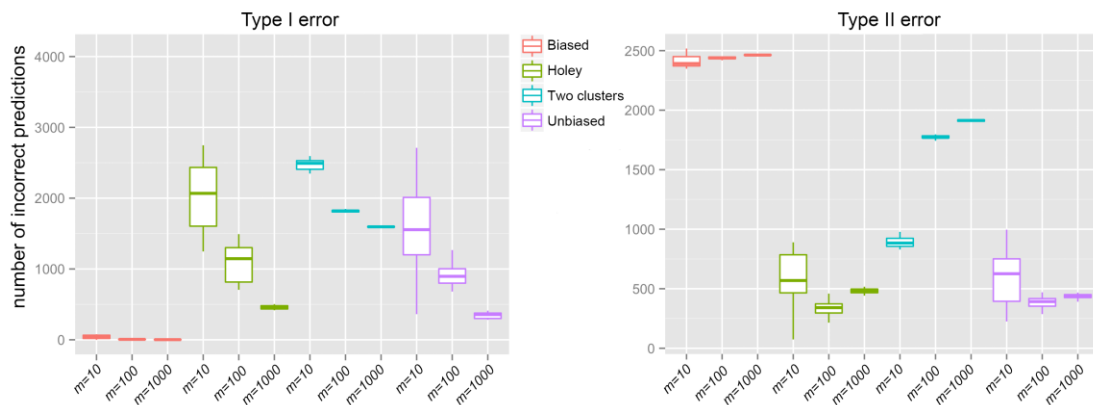
330 Research interests of the team include invasion ecology, virtual ecology, and
331 the evaluation of ecological niche modeling methods in ecology and
332 epidemiology.



334

335

336 **Figure 1. Type I and type II errors resulting from the KDE method using**
 337 **small sample sizes of $m = 1000$.** The red rectangle denotes a virtual
 338 fundamental niche (FN), while the blue points represent unbiased (a), biased
 339 (b), “holey”, as indicated by the black rectangle (c), and (d) two-clustered
 340 observations of the virtual FN. The green polygons are the estimated niche
 341 from the KDE method based on the blue observations. The overlap (pink) of
 342 the virtual FN and the estimated niche is the portion of virtual FN correctly
 343 predicted by the KDE method. The yellow area (b and d) outside of the virtual
 344 FN denotes type I error resulting from the KDE method. The white area with
 345 cross-shading denotes type II error resulting from the KDE method. Note that
 346 abundant occurrences reduce type I error at the cost of increased type II error.



347

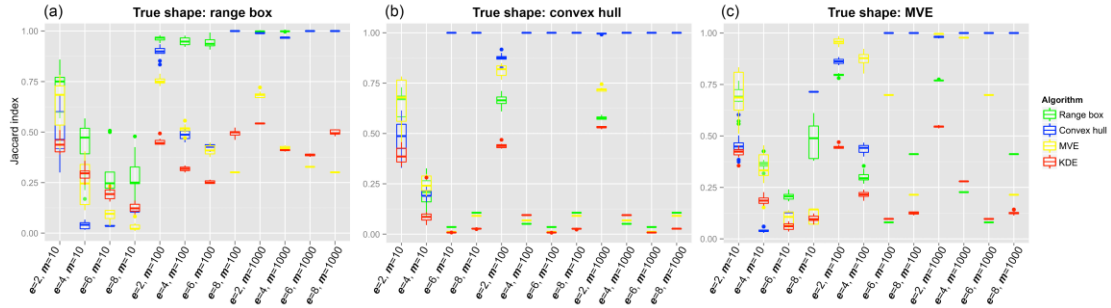
348 **Figure 2. Type I and II error for different sampling configurations**
 349 **estimated using the multivariate kernel density estimation (KDE) method.**

350 Left: type I error based on comparisons of the ‘true’ and estimated niche under
 351 unbiased (purple; Figs. 1.a, S1.a, S2.a), biased (red; Figs. 1.b, S1.b, S2.b),
 352 holey (green; Figs. 1.c, S1.c, S2.c), and two-clustered (blue; Figs. 1.d S1.d,
 353 S2.d) sampling configurations. Estimates are based on ten sampling replicates
 354 of 10, 100, and 1000 occurrences (m). Right: type II error from the same study

355 design. The y -axis indicates the number of false observations (left; type I error)

356 and the number of false negatives (right; type II error).

357



358

359 **Figure 3. Jaccard index for each modeling method based on the different**

360 **virtual fundamental niche shapes.** Fundamental niches (FN) were

361 represented as single hypercubes or range boxes (a), convex-hulls (b), and

362 ellipsoids (c). To estimate these 'true' FNs, we used four modeling methods:

363 range-box (RB; green), convex-hull (CH; blue), minimum-volume ellipsoid

364 (MVE; yellow), and multivariate kernel density estimation (KDE; red). Each

365 boxplot represents the Jaccard index of the niche according to 10 independent

366 subsamples of observations ($m = 10, 100, 1000$) collected randomly in a 2- to

367 8-dimensional dataset (e). Boxes closer to the top indicate better predictions (n)

368 in the form of high similarity or overlap between estimated (n) and 'true' virtual

369 fundamental niches (N).

370

371

372