1    **Title:** Empirical abundance distributions are more uneven than expected given their statistical baseline

2    **Running title:** SADs deviate from statistical baselines

3    **Author names and affiliations:**

4    Renata M. Diaz*[1], Hao Ye[2], S. K. Morgan Ernest[3]

5    [1] School of Natural Resources and Environment, University of Florida, Gainesville, Florida, USA.
6    renata.diaz@weecology.org; *corresponding author

7    [2] Health Science Center Libraries, University of Florida, Gainesville, Florida, USA. haoye@ufl.edu

8    [3] Department of Wildlife Ecology and Conservation, University of Florida, Gainesville, Florida, USA.
9    skmorgane@ufl.edu

10    **Address for correspondence:**

11    Renata M. Diaz
12    110 Newins-Ziegler Hall
13    PO Box 110430,
14    Gainesville, FL 32611-0430
15    renata.diaz@weecology.org
16    Phone: (352) 846-0643
17    Fax: (352) 392-6984

18    **Statement of authorship:** RMD and SKME conceived the analysis; HY devised the algorithm to sample the
19    feasible set, reviewed the coded implementation, and wrote the explanatory vignette; RMD conducted
20    the analyses and wrote the first draft of the manuscript; all authors contributed substantively to revisions.

21    **Data accessibility statement**: All data used are available publicly via Zenodo and figshare. Upon
22    publication, all code and data will be archived and made publicly available via Zenodo.

23    **Keywords:** Species abundance distributions; feasible set; combinatorics; macroecology; constraints

24    **Conflict of interests:** The authors declare no conflicts of interest.

25    **Type of article**: Reviews and Syntheses

26    **Word counts:**

27    Abstract: 147
28    Main text: 4202
29    No text boxes

30    **Number of references:** 25

31    **Number of figures, tables, and text boxes:** 4 figures; 1 table; 0 text boxes

32

33    Abstract

34    The prevalence of the species abundance distribution's hollow-curve shape across many communities is

35    frequently assumed to reflect ecological processes structuring communities. However, this hollow curve

36    can also emerge as a statistical phenomenon of dividing a particular number of individuals into a given

37    number of species. While the hollow curve may be a statistical artefact, ecological processes may induce

38    subtle deviations between empirical species abundance distributions and their statistically most probable

39    forms. Examining ~22,000 communities, we found that empirical species abundance distributions are

40    more skewed and uneven than their statistical baselines. However, small communities – with few species

41    or individuals – exhibit poorly-resolved statistical baselines, thereby reducing our capacity to detect

42    deviations. The extraordinarily skewed and uneven nature of empirical species abundance distributions

43    provides new avenues for testing ecological theory, while the issues posed by small communities

44    illustrate the limitations of statistical baselines for studying ecological patterns in small samples.

45    Introduction

46    The distribution of how the total number of individuals in a community are divided among the species in

47    that community, or the species abundance distribution (SAD), is one of the few ecological patterns whose

48    shape is so consistent that it is often considered an ecological law (Lawton 1999, McGill et al. 2007).

49    Across varied ecosystems and taxa, the species abundance distribution is dominated by a few very

50    abundant species and a larger number of increasingly rare species, generating a distinctive hollow- or J-

51    shaped curve (Fisher et al. 1943). Community ecologists have used the SAD to test numerous theories

52    intended to determine which biological processes are most important for structuring assemblages of

53    species, by comparing theoretical predictions for the SAD to observed SADs (McGill et al. 2007). However,

54    this approach has proven inconclusive because many theories predict similar shapes for the SAD (McGill

55    et al. 2007), and even experimental manipulations generate little variation in the shape of the SAD (Supp

56    and Ernest 2014). After decades of attention to its consistent pattern and potential as a theoretical

57    benchmark, the utility of the SAD for assessing the processes structuring ecological communities remains

58    unclear.

59    Accumulating evidence suggests that statistical constraints may actually generate the most striking

60    feature of the species abundance distribution – the hollow curve (White et al. 2012, Locey and White

61    2013, Frank 2019). Indeed, power-law or log-series distributions (i.e. hollow curves) emerge from many

62    distinct generative processes for abundance distributions generally (Frank 2009, Frank 2019). In the

63    specific case of the SAD, approaches using both statistical mechanics (i.e. the Maximum Entropy Theory

64    of Ecology (METE); Harte et al. 2008, Harte 2011) and combinatorics (i.e. 'the feasible set'; Locey and

65    White 2013) generate realistic hollow curves via the random division of the total number of individuals in

66    a community, $N$, into the total number of species present, $S$.  Given that these stochastically generated

67    SADs are excellent empirical fits to those seen in nature (Harte 2011, White et al. 2012, Locey and White

68    2013), failing to account for these considerations may have led us to focus on a distinctive but relatively

69    uninformative aspect of the SAD.

70    If SADs are statistically inclined to be hollow curves, it is no surprise that we have struggled to use the

71    hollow curve to definitively identify the ecological processes that shape SADs. However, this does not

72    necessarily mean that the SAD cannot be biologically informative. Biological factors may introduce subtle,

73    but meaningful, deviations between the shapes of observed SADs and the shapes of the SADs expected

74    due to the mathematical constraints imposed by $S$ and $N$, which we hereafter refer to as the "statistical

75    baseline" (Locey and White 2013, Harte and Newman 2014). Thus, it may be possible to use these

76    *deviations* to detect strong ecological processes or evaluate theories (Harte and Newman 2014, Xiao et al.

77    2016). If the vast majority of mathematically achievable SADs for a community share a similar shape, an

78    empirically observed SAD that deviates even slightly from this statistical baseline is unlikely to have

79    emerged at random (Locey and White 2013). Such a deviation is therefore likely to be the signature of a

80    non-random – i.e., biological – process operating on the relative abundances of species (Harte and

81    Newman 2014). It may be possible to evaluate ecological processes and theories based on how well they

82    predict these *deviations* between observed SADs and their statistical baselines.

83    Successfully interpreting SADs in this fashion depends on our capacity to detect and quantify deviations

84    between empirical observations and statistical baselines, which requires metrics and computational

85    approaches that allow us to quantify and interpret whatever deviations may exist. Here, we build upon

86    the combinatoric approach developed by Locey and White (2013) to define and explore the statistical

87    baselines for SADs. For a given $N$ (total number of individuals) and $S$ (total number of species), there exists

88    a finite (but potentially very large) set of possible distributions of individuals into species. Collectively, this

89    set of possible SADs is referred to as the feasible set, with each possible SAD constituting a single element

90    of the set. If an observed SAD is drawn at random from the set of mathematically possible SADs, it is likely

91    to have a shape similar to the shapes most common in the feasible set. The feasible set can therefore be

92    used as a statistical baseline for assessing whether observed SADs deviate from what is likely to occur

93    merely due to mathematical constraints, and to explore how the statistical baseline varies over across

94    ranges of values of S and N (Locey and White 2013).

95    The specificity, or vagueness, of the expectations derived from the statistical baseline is critically

96    important for disentangling the aspects of the SAD that are likely generated by statistical constraints from

97    those generated by other processes. If the vast majority of mathematically possible SADs are similar in

98    shape – generating a very specific, narrowly defined statistical baseline – then even small deviations

99    between an observed SAD and this baseline can signal the operation of ecological processes. However, if

100   many different shapes occur with more even frequency in the feasible set, the statistical baseline is less

101   specific and less well defined, and our sensitivity for distinguishing biological signal from statistical

102   constraints is greatly reduced. In general, a poorly defined statistical baseline is more likely to occur when

103   the size of the community, in terms of $S$ and $N$, is small, because in such cases there may be too few

104   possible SADs in the feasible set for a particular shape to emerge as the most common, and therefore

105   most likely, shape. When this occurs, we have reduced confidence that even an observation that deviates

106   from the statistical baseline did not emerge at random from the relatively restricted pool of possible

107   outcomes (Jaynes 1957). This general concern has been acknowledged in efforts to comparing ecological

108   observations to statistical baselines (Harte 2011, White et al. 2012, Locey and White 2013) but there has

109   not yet been a quantification of these effects for the SAD or an identification of the range of community

110   sizes most strongly affected. Because ecologists study the SAD for communities varying in size from the

111   very small – $S$ and $N$ < 5 – to the enormous – $S$ and $N$ >> 1000 – identifying the community sizes for which

112   we can and cannot confidently detect deviations from the statistical baseline is necessary to appropriately

113   contextualize our interpretations.

114   Here we use the feasible set to define statistical baselines for empirical SADs for 22,000 communities of

115   birds, mammals, trees, and miscellaneous other taxa. We then compare *observed* SADs to their

116    corresponding statistical baselines and evaluate 1) if the shapes of observed SADs consistently deviate

117    from their statistical baseline, 2) how the specificity of the statistical baseline varies over ranges of $S$ and

118    $N$, and 3) whether this variation appears to be associated with variation in our capacity to detect

119    deviations between observations and the corresponding baselines.

120    **Methods**

121    *Datasets*

122    We used a compilation of community abundance data for trees, birds, mammals, and miscellaneous

123    additional taxa (White et al. 2012, data from Baldridge 2015 available at

124    https://doi.org/10.6084/m9.figshare.95843.v4, Baldridge 2016, data from Baldridge 2016 available at

125    https://zenodo.org/record/166725).  This compilation consists of cleaned and summarized community

126    abundance data for trees obtained from the Forest Inventory and Analysis (Woudenberg et al 2010) and

127    Gentry transects (Phillips and Miller 2002), birds from the North American Breeding Bird Survey (Sauer et

128    al. 2013), mammals from the Mammal Community Abundance Database (Thibault et al. 2011), and a

129    variety of less commonly sampled taxa from the Miscellaneous Abundance Database (Baldridge 2015).

130    Because characterizing the random expectation of the SAD is computationally intractable for very large

131    communities, we filtered our datasets to remove communities with more than 40720 individuals, which

132    was the largest community we successfully analyzed. This resulted in the removal of 4 communities from

133    the Miscellaneous Abundance Database. We further filtered the FIA database. Of the 103,343

134    communities in FIA, 92,988 have fewer than 10 species. Rather than analyze all these small communities,

135    we randomly selected 10,000 small communities to include in the analysis. We also included all FIA

136    communities with more than 10 species, which added 10,355 FIA communities to the analysis and

137    resulted in a total of 20,355 FIA communities. Finally, for sites that had repeated sampling over time, we

138    followed White et al. (2012) and Baldridge (2016) and analyzed only a single, randomly selected, year of

139    data, because samples taken from a single community at different time points are likely to covary. It

140    should be noted that our analyses include data from the Mammal Community Database and

141    Miscellaneous Abundance Database that were collected over longer timescales and cannot be

142    disaggregated into finer units of time. We also removed from our analyses any communities with only

143    one species, or for which N = S or N = S + 1, because these communities have only one mathematically

144    possible SAD. Our final dataset consisted of ~22,000 communities with S and N ranging from 2 to 250 and

145    3 to 40714, respectively (see Figure S1 in Supporting Information). Details and code for the filtering

146    process can be found in Appendix S2.

147    *Generating the statistical baseline*

148    We use the concept of the "feasible set" to establish a statistical baseline for the SAD (Locey and White

149    2013). For a given number of individuals *N*, there is a finite number of unique ways to partition those

150    individuals into *S* species. The complete set of these unique partitions is the feasible set. Because, in this

151    approach, neither species nor individuals are distinguishable from each other, partitions are unique if and

152    only if they differ in the number of species that have a particular abundance (Locey and White, 2013).

153    Operationally, this means that for *S = 3* and *N = 9*, the species abundances *(1, 3, 5)* and *(2, 2, 5)* count as

154    distinct partitions, but *(1, 3, 5)* and *(3, 1, 5)* are only one element of the feasible set because they each

155    contain one species with an abundance 1, 3, and 5, respectively, and they differ only in the *order* of the

156    numbers. In the absence of justification for additional assumptions regarding the distinguishability of

157    species and/or individuals, we adopted this simple set of assumptions that has previously been shown to

158    generate realistic statistical baselines (Locey and White 2013).

159    While it is possible to list all possible partitions in the feasible set for small *S* and *N*, the size of the feasible

160    set increases rapidly with *S* and *N*. Therefore, characterizing the statistical properties of the feasible set

161    for large *S* and *N* can be computationally intensive. This renders it necessary to draw samples from the

162    feasible set, rather than enumerating all of its elements. Unbiased sampling of large feasible sets is a

163    nontrivial computational problem that has constrained previous efforts in this vein (Locey and White

164    2013). We developed an algorithm to efficiently and uniformly sample feasible sets even for large values

165    of $S$ and $N$. In brief, the algorithm takes a generative approach to sampling the feasible set. Individuals are

166    allocated one species at a time, beginning with the least abundant species. At each step, the number of

167    individuals to allocate for the current species is determined at random, with probability based on the

168    number of feasible sets with that specific abundance for that species, and conditional on the individuals

169    that have already been allocated. For example, if we have 3 species and 7 individuals, the least abundant

170    species can have an abundance of 1 or 2. Allocating 1 individual to the least abundant species allows for

171    the SADs (1, 1, 5), (1, 2, 4), and (1, 3, 3), and allocating 2 individuals to the least abundant species means

172    the only possible SAD is (2, 2, 3). We therefore allocate 1 individual with probability 3/4, and 2 individuals

173    with probability 1/4. If, at the first step, we allocated 1 individual to the least abundant species, the

174    second species can have an abundance of 1, 2 or 3 with equal probability, because there is exactly 1

175    possible SAD with each of these possible abundances for the first two species. This process is continued

176    until all individuals have been allocated. We implemented this algorithm in an R package, available on

177    GitHub at [www.github.com/diazrenata/feasiblesads](www.github.com/diazrenata/feasiblesads); details of the sampling methodology are available in

178    Appendix S3.

179    For every community in our database, we drew 4000 samples from the feasible set to characterize the

180    distribution of statistically probable shapes for the SAD. We filtered the 4000 samples to unique

181    elements. For small values of S and N, it can be impossible or highly improbable for the 4000 samples

182    from the feasible set to all be unique, but for large communities, all 4000 are usually unique. We refer to

183    this as the sampled feasible set.

184    *Comparing observed SADs to their baselines*

185  If all SADs in a feasible set are equally likely to occur, then an SAD with a particular $S$ and $N$ is likely to

186  have a shape similar to the shape that is most common among the SADs in the feasible set for the same $S$

187  and $N$; in contrast, strong processes may cause observed SADs to have shapes that deviate from this

188  statistical baseline (Locey and White 2013). We focus on two metrics to describe the shape of the SAD,

189  skewness and Simpson's evenness. Skewness measures the asymmetry of a distribution around its mean,

190  and Simpson's evenness is a commonly used metric in ecology for assessing how equitably abundance is

191  distributed across species. By calculating these metrics for each of the samples in the community's

192  sampled feasible set (see *Generating the statistical baseline*, above), we generated a distribution

193  describing the general shape (i.e. evenness or skewness) expected from the randomly sampled SADs.

194  Note that skewness, as implemented in the R package "e1071" (Meyer et al. 2019), always evaluates to 0

195  for distributions with only two species, and we therefore excluded those cases from analyses of skewness

196  (but included those communities for analyses using Simpson's evenness).

197  To assess whether the shape of an observed SAD was statistically unlikely, we calculated Simpson's

198  evenness and skewness for the observed SAD and compared these observed values to the distributions of

199  evenness and skewness obtained from that community's sampled feasible set. An observed SAD's

200  deviation from its feasible set was determined by computing the percentile rank of its skewness and

201  evenness relative to the sampled distributions for skewness and evenness, respectively. These percentile

202  ranks are then comparable across different community sizes, allowing broad-scale assessment across

203  wide ranges of $S$ and $N$. After aggregating across communities, if observed SADs reflect random draws

204  from their feasible sets, their percentile rank values should be uniformly distributed from 0 to 100.

205  However, if observed SADs are consistently more skewed or even than their feasible sets, the percentile

206  values will be disproportionately concentrated towards the extremes. Because an earlier survey in this

207  space (Locey and White 2013) found that the tendency is for empirical SADs to be more skewed and less

208  even than their feasible sets, we used one-tailed 95% confidence intervals and tested for unusually *high*

209     values for skewness and *low* values for evenness. This comparison is not meaningful if there are very few

210     unique values in the distributions of skewness and evenness, which can occur for small feasible sets. We

211     therefore excluded communities for which the distribution of skewness or evenness values from the

212     sampled feasible set had fewer than 20 unique values (in these cases, it is impossible for an observation

213     to fall above or below the 95th or 5th percentile, respectively). Our final aggregated analyses included

214     22,142 communities for evenness and 22,325 communities for skewness.

215     *The narrowness of the expectation*

216     We also used the distributions of skewness and evenness from the sampled feasible set to describe the

217     relative specificity of the statistical baseline, in order to assess in what situations there could be

218     challenges in determining whether observed communities differ from their statistical baselines. We

219     quantified the narrowness of a distribution as the ratio of the range of values encompassed within a 95%

220     density interval relative to the full range of values in the distribution (Figure 1). This breadth index for the

221     statistical baseline ranges from 0 (a very narrow distribution and well-resolved baseline) to 1 (a very

222     broad distribution), and is comparable across feasible sets for varying combinations of $S$ and $N$. This

223     metric corresponds qualitatively to more computationally-intensive approaches to measuring the self-

224     similarity of the elements of feasible sets (see Appendix S4).

225     **Results**

226     *Observed SADs compared to their feasible sets*

227     For four of the five datasets we analyzed – BBS, Gentry, Mammal Communities, and Misc. Abund –

228     empirical SADs are highly skewed and highly uneven relative to their feasible sets, much more frequently

229     than would be expected by chance (Figure 2, Table 1). Combined across these four datasets, 16% of

230     observed SADs are more skewed than 95% of their feasible sets, and 31% are less even than 95% of their

231     feasible sets. For SADs randomly sampled from the feasible set, we would expect only 5% of observed

232    distributions to fall in these extremes. In contrast to the other datasets, the SADs from the FIA dataset

233    exhibit percentile scores that are more uniformly distributed: 5% of observations are more skewed than

234    95% of their feasible sets, and 9% of observations are less even than 95% of their feasible sets.

235    *The narrowness of the expectation*

236    The ability to detect deviations from the statistical baseline depends on the distribution of SADs in the

237    feasible set. Here, the statistical baseline for both skewness and evenness becomes more narrowly

238    defined as the size of the feasible set increases (Figure 1; Figure S5), making even small deviations in

239    skewness or evenness statistically meaningful and readily detectable. However, for communities with

240    relatively small feasible sets – fewer than approximately 1000 elements for skewness, and approximately

241    200 elements for evenness – the breadth index approaches 1, meaning that a 95% density interval of the

242    values in the distribution spans nearly the entire range of values (Figure S5). In particular, the FIA dataset

243    is dominated by small communities for which the breadth index is very high, reflecting relatively broad

244    and nonspecific statistical expectations for the shape of the SAD derived from the feasible set (Figure 3).

245    **Discussion**

246    We found widespread evidence that SADs for a range of real ecological communities are more skewed

247    and less even than expected given the distribution of shapes within their feasible sets. These deviations

248    may signal that ecological processes operate on top of statistical constraints, thereby driving the SAD

249    away from common shapes that would be observed in the absence of a dominating non-statistical

250    process. Our results suggest that the prevailing processes structuring these communities tend to be those

251    that cause abundance distributions to be more uneven – rather than those that produce more even

252    abundances across species. Ecological processes may lengthen the rare tail of the SAD, for example by

253    promoting the persistence of rare species at very low abundances (e.g. Yenni et al. 2012). Or, they could

254    drive abundant species to have larger populations that would be statistically expected, without driving

255    other species entirely to extinction (Chesson 2000). Although a disproportionate number of communities

256    deviated statistically from their feasible sets, there were also many communities for which we did not

257    detect deviations. In such cases, multiple ecological processes may operate simultaneously and with

258    countervailing impacts on abundance distributions, resulting in no dominating net effect on the shape of

259    the distribution beyond that imposed by fundamental constraints (Harte 2011; Harte and Newman 2014).

260    Going forward, testing whether ecological theories or common functional approximations (e.g. the log-

261    normal distribution) accurately predict this range of variation in deviations between observed SADs and

262    their statistical baselines may be much more fruitful than focusing only on the general form of the SAD

263    that emerges from statistical constraints (McGill et al. 2007; Locey and White 2013).

264    Unlike the other four datasets, communities in the FIA dataset showed weak or no evidence of deviations

265    from their feasible sets. These results may reflect statistical phenomena related to community size. The

266    FIA communities are by far the smallest across our datasets (Figure S1). Communities with small values of

267    S and N have smaller feasible sets; when there are relatively few possible SADs, the distributions of

268    evenness and skewness values from the feasible set are less narrowly peaked, meaning there is a weaker

269    statistical distinction between "common" and "extreme" shapes for the SAD (Figure 1). In fact, across the

270    range of community sizes present in our datasets, the feasible sets for small communities generally

271    generated broader distributions of evenness, and especially skewness, than those for large communities

272    (Figure S5).  For such communities, the deviations – or lack thereof – that we perceive are less

273    informative than for larger communities with more strongly defined statistical baselines (Jaynes 1957).

274    If the lack of discernable deviations from the feasible set for the FIA communities is indeed a byproduct of

275    their generally small size, then we would expect similarly-sized communities from other datasets to have

276    similar results. We identified 371 communities from other datasets with $S$ and $N$ matching communities

277    from FIA. We found no difference in the distribution of percentile scores between communities from FIA

278    and communities from other datasets (Figure 4; Table S6), confirmed via Kolmogorov-Smirnov tests (for

279     evenness, *D* = 0.04 and *p* = 0.87; for skewness, *D* = 0.07 and *p* = 0.37). Although 371 communities

280     constitute a small sample relative to the 20,355 FIA communities we analyzed, these results point to

281     community size, and not attributes specific to FIA, as a likely explanation for the weak evidence for

282     deviations across the full FIA dataset.

283     If this is indeed the case, it means that small-community considerations may affect our capacity to

284     meaningfully distinguish signal from randomness using this approach. FIA communities, with their broad

285     distributions of shape metrics and overall lack of detectable signal, have on the order of 10 species and

286     50-100 individuals. While these values do not constitute hard thresholds, they may indicate a general

287     range of values below which we have relatively diminished power to detect deviations from the statistical

288     baseline represented by the feasible set. To meaningfully draw inferences from deviations in these small

289     communities, we will likely need more sensitive metrics (than skewness and evenness), and/or theories

290     that generate more specific predictions for the SAD. In the absence of such, we may stand to learn the

291     most by focusing on SADs from relatively large communities.

292     It is also important to recognize that there are multiple plausible approaches to defining a statistical

293     baseline for the SAD, of which we have taken only one (Haegeman and Loreau 2008, Locey and White

294     2013). Our approach follows Locey and White (2013) and reflects the random partitioning of individuals

295     into species, with the resulting distributions considered unique if the species' abundance values are

296     unique, regardless of the order in which the values occur. The philosophy behind the feasible set reflects

297     a longstanding approach in the study of abundance distributions: to focus on the shape of the distribution

298     without regard to species' identities (McGill 2007). To include differences in *order* in the statistical

299     baseline would imply that identifying *which* species contain the most or least individuals is important for

300     evaluating theory. Other formulations for the statistical baseline may be equally valid and generate

301     different statistical expectations, including forms that approximate exponential, Poisson, or log-series

302    distributions (Harte et al. 2008, Favretti 2018). Comparing the results that emerge from different

303    baselines will be an important next step towards reinvigorating the use of the SAD as a diagnostic tool.

304    Our study demonstrates the utility, and the potential challenges, of applying tools from the study of

305    complex systems and statistical mechanics to the study of ecological communities (Haegeman and Loreau

306    2008, Harte 2011, White et al. 2012, Harte and Newman 2014). While concepts such as maximum

307    entropy and the feasible set are promising new horizons for macroecology, the small size of some

308    ecological communities may present difficulties that do not occur as often in the domains for which these

309    tools were originally developed (Jaynes 1957, Haegeman and Loreau 2008). When the observed numbers

310    of species and individuals are too small to generate highly resolved statistical baselines, these approaches

311    will be less informative than we might hope – as appears to be the case for the smallest communities in

312    our analysis. In larger communities, where mathematical constraints have stronger effects on the general

313    form of the SAD, our results show that these constraints alone do not fully account for the extremely

314    uneven SADs we often observe in nature – leaving an important role for ecological processes. This ability

315    to detect and diagnose the specific ways in which empirical SADs deviate from randomness can generate

316    new avenues for understanding how and when biological drivers affect the SAD. There are, of course, still

317    many elements to be improved in our ability to distinguish biological signal from randomness, including

318    assessing alternative statistical baselines and calibrating our expected power to detect deviations,

319    especially for small communities. Indeed, more sensitive metrics could also enable identification of

320    processes that operate through time (note that, in this analysis, we sampled only one time point for each

321    community). Continuing to explore and account for the interplay between statistical constraint and

322    biological process constitutes a promising and profound new approach to our understanding of this

323    familiar, yet surprisingly mysterious, ecological pattern.

324

330

331

332    References

333    Baldridge, E. (2015). Miscellaneous Abundance Database. figshare. Available at:

334         https://doi.org/10.6084/m9.figshare.95843.v4

335    Baldridge, E., Harris, D.J., Xiao, X. & White, E.P. (2016). An extensive comparison of species-abundance

336         distribution models. *PeerJ*, 4, e2823.

337    Baldridge, E., Harris, D.J., Xiao, X. & White, E.P. (2016). Data from *An extensive comparison of species-*

338         *abundance distribution models*. Zenodo. Available at: https://zenodo.org/record/166725.

339    Chesson, P. (2000). Mechanisms of Maintenance of Species Diversity. *Annual Review of Ecology and*

340         *Systematics*, 31, 343–366.

341    Favretti, M. (2018). Remarks on the Maximum Entropy Principle with Application to the Maximum

342         Entropy Theory of Ecology. *Entropy*, 20, 11.

343    Fisher, R.A., Corbet, A.S. & Williams, C.B. (1943). The Relation Between the Number of Species and the

344         Number of Individuals in a Random Sample of an Animal Population. *Journal of Animal Ecology*, 12,

345         42–58.

346    Frank, S.A. (2009). The common patterns of nature. *Journal of Evolutionary Biology*, 22, 1563–1585.

347    Frank, S.A. (2019). The common patterns of abundance: the log series and Zipf's law. *F1000Res*, 8, 334.

348    Haegeman, B. & Loreau, M. (2008). Limitations of entropy maximization in ecology. *Oikos*, 117, 1700–

349         1710.

350    Harte, J. (2011). *Maximum Entropy and Ecology: A Theory of Abundance, Distribution, and Energetics*.

351         Oxford University Press.

352    Harte, J. & Newman, E.A. (2014). Maximum information entropy: a foundation for ecological theory.

353         *Trends in Ecology & Evolution*, 29, 384–389.

354    Harte, J., Zillio, T., Conlisk, E. & Smith, A.B. (2008). Maximum Entropy and the State-Variable Approach to

355         Macroecology. *Ecology*, 89, 2700–2711.

356    Jaynes, E.T. (1957). Information Theory and Statistical Mechanics. *Phys. Rev.*, 106, 620–630.

357    Lawton, J.H. (1999). Are There General Laws in Ecology? *Oikos*, 84, 177.

358    Locey, K.J. & White, E.P. (2013). How species richness and total abundance constrain the distribution of

359        abundance. *Ecology Letters*, 16, 1177–1185.

360    McGill, B.J., Etienne, R.S., Gray, J.S., Alonso, D., Anderson, M.J., Benecha, H.K., *et al.* (2007). Species

361        abundance distributions: moving beyond single prediction theories to integration within an

362        ecological framework. *Ecol Letters*, 10, 995–1015.

363    Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A. & Leisch, F. (2019). *e1071: Misc Functions of the*

364        *Department of Statistics, Probability Theory Group (Formerly: E1071),* TU Wien. R package version

365        1.7-4. https://CRAN.R-project.org/package=e1071

366    Phillips, O. & Miller, J.S. (2002). *Global patterns of plant diversity: Alwyn H. Gentry's forest transect data*

367        *set*. Missouri Botanical Press.

368    Sauer, J.R., Link, W.A., Fallon, J.E., Pardieck, K.L. & Ziolkowski, D.J. (2013). The North American Breeding

369        Bird Survey 1966–2011: Summary Analysis and Species Accounts. *North American Fauna*, 1–32.

370    Supp, S.R. & Ernest, S.K.M. (2014). Species-level and community-level responses to disturbance: a cross-

371        community analysis. *Ecology*, 95, 1717–1723.

372    Thibault, K.M., Supp, S.R., Giffin, M., White, E.P. & Ernest, S.K.M. (2011). Species composition and

373        abundance of mammalian communities. *Ecology*, 92, 2316–2316.

374    White, E.P., Thibault, K.M. & Xiao, X. (2012). Characterizing species abundance distributions across taxa

375        and ecosystems using a simple maximum entropy model. *Ecology*, 93, 1772–1778.

376    Woudenberg, S.W., Conkling, B.L., O'Connell, B.M., LaPoint, E.B., Turner, J.A. & Waddell, K.L. (2010). The

377        Forest Inventory and Analysis Database: Database description and users manual version 4.0 for

378        Phase 2. *Gen. Tech. Rep. RMRS-GTR-245. Fort Collins, CO: U.S. Department of Agriculture, Forest*

379        *Service, Rocky Mountain Research Station. 336 p.*, 245.

380    Xiao, X., O'Dwyer, J.P. & White, E.P. (2016). Comparing process-based and constraint-based approaches

381        for modeling macroecological patterns. *Ecology*, 97, 1228–1238.

382    Yenni, G., Adler, P.B. & Ernest, S.K.M. (2012). Strong self-limitation promotes the persistence of rare

383        species. *Ecology*, 93, 456–461.

384

385     **Figure legends**

386     Figure 1. Large feasible sets may allow better detection of deviations from the statistical baseline by

387     generating more specific, narrowly-defined baselines. We illustrate this phenomenon using 3 hypothetical

388     communities: a small community ($S = 4$, $N = 34$; top row), an intermediate community ($S = 7$, $N = 71$;

389     middle row), and a large community ($S= 44$, $N = 13360$; bottom row). The large communiity has

390     approximately 6.59e+70 possible SADs in its feasible set, while the intermediate community has 60,289

391     and the small community has only 297. For every SAD sampled from the feasible set (left column), we

392     calculate the skewness (color scale) and evenness (not shown). The distributions of these values (right

393     column) constitute the statistical baseline. We define a "breadth index" as the ratio of the range

394     encompassed in the one-tailed 95% density interval (distance between red lines, right), compared to the

395     full range of values for the statistic (distance between the maximum and minimum values). As $S$ and $N$

396     increase, the size of the feasible set increases, resulting in a narrower statistical baseline (smaller breadth

397     index) – thus enabling easier detection of deviations that may be the result of ecological processes

398     affecting the SAD.

399

400     Figure 2. Many ecological communities are more skewed (left) or uneven (right) than their statistical

401     baselines. Percentile ranks are calculated by comparing each community to its sampled feasible set, with

402     very high or very low percentile ranks reflecting extreme values relative to statistical baselines. The

403     vertical red line marks the 95[th] percentile for skewness and the 5[th] percentile for evenness. Species

404     abundance distributions that are sampled at random from the feasible set will produce percentile ranks

405     that are uniformly distributed from 0 to 100, with approximately 5% of values above or below the 95[th]

406     and 5[th] percentiles, respectively. In contrast, most datasets have more communities that are highly

407     skewed or uneven than would be expected by chance.

408

409     Figure 3. Breadth indices of skewness (left) and evenness (right) indicate varying ability to detect the

410     deviations between observations and the statistical baseline. The breadth index (see Figure 1) quantifies

411     how narrowly-defined the statistical baseline is; high values indicate broad, poorly-defined statistical

412     baselines that may impede our ability to confidently detect deviations between observations and what is

413     expected given the baseline. Most datasets contain a mixture of communities with broad and narrow

414     statistical baselines, but some – particularly the skewness baseline for the Forest Inventory and Analysis –

415     have consistently high breadth indices across all of their communities, suggesting that skewness is not an

416     effective metric for distinguishing empirical observations from the feasible set. In general, the breadth

417     index for evenness (right panels) indicates more narrow statistical baselines than those for skewness.

418

419     Figure 4. Small communities exhibit consistently broad statistical baselines (top), and consistently weak

420     evidence of deviations for observed SADs (bottom), regardless of the originating dataset. For a subset of

421     371 communities from the Forest Inventory and Analysis with communities in other datasets with

422     matching $S$ and $N$, we generate distributions of breadth indices for skewness and evenness (top) and

423     compute corresponding percentile ranks for the observed SADs (bottom). Visually, there is no difference

424     between FIA (left panels) and other datasets (right panels), when they are matched in $S$ and $N$. This is

425     confirmed by Kolmogorov-Smirnov tests for the breadth indices (for evenness, $D = 0.04$ and $p = 0.91$; for

426     skewness, $D = 0.03$ and $p > 0.99$) and percentile ranks (for evenness, $D = 0.04$ and $p = 0.87$; for skewness,

427     $D = 0.07$ and $p = 0.37$).

428

429

430