

1 Stronger Evidence for Trait–Environment Association by Pre-processing of  
2 Abundance Tables

3  
4 Cajo J. F. ter Braak<sup>1</sup>

5 <sup>1</sup>Mathematical & Statistical Methods group – Biometris, Wageningen University &  
6 Research, Wageningen, the Netherlands

7 Email: cajo.terbraak@wur.nl

8 ORCID: <http://orcid.org/0000-0002-0414-8745>

9

10 Version 1 : February, 14, 2026; CC-BY

11 Version 1.2: February, 24, 2026; CC-BY

12

13 **Abstract**

14

15 Understanding trait–environment relationships is central to predicting community  
16 responses to environmental change, yet statistical evidence for such relationships is  
17 often weak in observational datasets. Here, I introduce an  $N_2$ -processing method for  
18 abundance tables, grounded in the observation that the precision of  
19 community-weighted means (CWMs) and species niche centroids (SNCs) is  
20 proportional to Hill’s effective number  $N_2$ . I refine this concept by defining  
21 informativeness, which down-weights ubiquitous species whose abundances provide  
22 little environmental information. The resulting species and site weights are  
23 implemented through iterative proportional fitting (IPF), which preserves the species–  
24 site interaction structure while aligning totals with informativeness and effective  
25 numbers.

26 The performance of  $N_2$ -processing was evaluated using simulations and 80 published  
27 trait–environment datasets from the CESTES global database, analyzed with  
28 fourth-corner correlation analysis (FC), RLQ, and double-constrained correspondence  
29 analysis (dc-CA). In simulations,  $N_2$ -processing substantially increased the power of  
30 the max test, especially when abundances were untransformed. In the CESTES  
31 analyses, the proportion of datasets exhibiting at least weak trait–environment  
32 association increased from 13% when using FC with multiple-testing adjustment to  
33 68% when using dc-CA with  $N_2$ -processing. The improvement was largest when  
34 associations were weak and diminished when strong prior abundance transformations  
35 (such as the logarithmic transformation) had already down-weighted dominant  
36 species.

37 A square-root abundance transformation followed by  $N_2$ -processing is recommended  
38 as an effective and robust approach to strengthening inference on trait–environment  
39 relationships. The  $N_2$ -processing is essential because it removes arbitrary, data-analytic  
40 choices and ensures that dc-CA operates under conditions that match the assumptions  
41 of the model it is testing.

42 Key words: Hill Number, weighted averaging, community weighted means, trait-  
43 environment association, fourth-corner correlation, RLQ, double-constrained  
44 correspondence analysis (dc-CA).

## 45 **1. Introduction**

46 Understanding how biodiversity and the ecosystem services it supports will change  
47 under global pressures is a central question in ecology. Predicting shifts in species  
48 assemblages along environmental gradients is a major part of this challenge, yet the  
49 large number of species involved makes a purely species-based, data-driven approach  
50 impractical. A more efficient strategy is to describe species through their functional  
51 traits, just as environments are described by sets of environmental variables.

52 Trait-based ecology links species traits, species abundances and environmental  
53 conditions, enabling the detection of trait–environment relationships that clarify  
54 underlying mechanisms (Garnier *et al.* 2015; Green *et al.* 2022; Westoby 2025).  
55 Trait–environment analyses integrate species abundance data, environmental  
56 variables, and species traits, each with its own structure and limitations (de Bello *et*  
57 *al.* 2025).

58 Careful pre-processing of these tables is essential for drawing reliable inferences  
59 about community responses to environmental filters (McCune 2015; Broennimann,  
60 Collart & Guisan 2026). Here the focus is on pre-processing the abundance table to  
61 strengthen evidence for trait–environment associations.

62 We consider pre-processing approaches tailored to the fourth-corner correlation  
63 analysis (FC) (Legendre, Galzin & Harmelin-Vivien 1997; Dray & Legendre 2008;  
64 ter Braak 2017), the three-table RLQ method (Dolédec *et al.* 1996; Dray *et al.* 2014),  
65 and double-constrained correspondence analysis (dc-CA) (ter Braak, Šmilauer & Dray  
66 2018; ter Braak & van Rossum 2025a). The latter extends community-weighted mean  
67 (CWM) regression (Peres-Neto, Dray & ter Braak 2017; Pinho *et al.* 2024) by  
68 allowing simultaneous inclusion of multiple traits and environmental variables.  
69 Although more advanced methods exist, these approaches, especially dc-CA, provide  
70 a practical and robust framework for testing trait–environment relationships in  
71 multivariate ecological data (ter Braak & van Rossum 2025a; ter Braak 2026).

72 Abundance is the weight in a CWM and is often left untransformed in view of  
73 Grime’s dominance hypothesis (Grime 1973), which suggests that dominant species  
74 largely determine ecosystem functioning and thus argues for using raw abundances.  
75 Nevertheless, in practice, transforming abundances to presence–absence is often  
76 considered as an alternative. This paper proposes a statistical perspective on  
77 abundance transformations aimed at strengthening statistical evidence for trait–  
78 environment association—an important goal, as such evidence is often weak in  
79 observational data (Anderegg 2023). To ensure that models generalize to unseen trait  
80 values and trait combinations, site-level CWM regression should be accompanied by  
81 species-level Species-Niche-Centroid (SNC) regression. For reliable inference, both  
82 regressions must be significant, as first elaborated for FC through the max test (ter  
83 Braak, Cormont & Dray 2012; ter Braak 2026).

84 Ter Braak (2019) proposed  $N_2$ -weighting of CWM and SNC regressions, with  $N_2$  the  
85 Hill diversity number of order 2 (Hill 1973a). Hill numbers of all orders have the  
86 interpretation as effective numbers of occurrences. Here it is shown that order 2 is  
87 optimal for use in weighted averaging—a key ingredient of FC, RLQ and dc-CA—as  
88 the precision of a weighted average (CWM or SNC) is proportional to  $N_2$ . However, a

89 species occurring with equal abundance in all sites has  $N_2$  equal to the number of sites,  
90 yet is not informative about the environment (Lepš & de Bello 2023). A second  
91 limitation is that Ter Braak's (2019)  $N_2$ -weighting is difficult to extend to dc-CA, as it  
92 would require separate species-level and site-level dc-CA analyses.

93 Both issues are resolved here. First, the effective number  $N_2$  is replaced by a measure  
94 of informativeness,  $N_2(1 - N_2/n)$ , with  $n$  the number of sites, so that a species  
95 uniformly distributed across sites is no longer informative. Second, the weighting  
96 scheme is reformulated as a transformation of the abundance table that preserves  
97 species–site interactions while modifying species and site totals to reflect  
98 informativeness and the effective number of informative species. These totals then act  
99 as optimal weights in FC, RLQ, and dc-CA. Technically, this is achieved via iterative  
100 proportional fitting (IPF).

101 The effectiveness of this approach was demonstrated in preliminary simulations in ter  
102 Braak (2026). Here the simulations are extended, the method is fully documented, and  
103 evaluated across 80 published trait–environment datasets from the CESTES database  
104 (Jeliaskov *et al.* 2020), by comparing FC, RLQ, and dc-CA with and without  
105  $N_2$ -processing and with and without prior abundance transformations. Pre-processing  
106 is considered beneficial when it increases statistical power in simulations and tends to  
107 strengthen evidence for trait–environment associations (i.e., lowers test  $P$ -values)  
108 across the CESTES datasets.

109 The paper is structured as follows. Section 2 describes the method and its evaluation,  
110 Section 3 presents the results, and the paper concludes with a discussion. Technical  
111 details and additional results are provided in three supplements.

## 112 **2. Theory and Methods**

### 113 **2.1 $N_2$ -processing**

114 Correspondence analysis (CA) originated in ecology under the name *reciprocal*  
115 *averaging* (Hill 1973b), because species scores in CA are weighted averages of latent  
116 site scores, and site scores are in turn weighted averages of species scores. Weighted  
117 averages are now most familiar to ecologists through the community-weighted mean  
118 trait (CWM). The CWM of a site is the abundance-weighted average of a trait (or trait  
119 composite) across species. Conversely, the species niche centroid (SNC) is the  
120 abundance-weighted average of an environmental variable (or gradient) across sites  
121 for a given species (Peres-Neto *et al.* 2017). Fourth-corner correlation analysis (FC)  
122 can be seen as a combination of CWM- and SNC-regression analysis (Peres-Neto *et*  
123 *al.* 2017; ter Braak, Peres-Neto & Dray 2018).

124 Under a null model where abundance is unrelated to the variable being averaged, the  
125 precision (inverse variance) of a weighted average is proportional to  $N_2$ , the Hill  
126 diversity of order 2 (Hill 1973a) (Supplement A).  $N_2$  represents the *effective number*  
127 of occurrences. For example, if abundances are 100, 1, and 1, the number of  
128 occurrences is three, but the effective number  $N_2 \approx 1.04$  because one weight  
129 dominates. Although the weighted average is computed over three values, only one  
130 effectively contributes to its precision. At the other extreme, a species occurring  
131 almost everywhere carries little environmental information (Lepš & de Bello 2023). If  
132 a species occurs equally in all sites, its weighted average equals the mean  
133 environment, making it uninformative. Thus, informativeness increases with effective

134 occurrence when rare but decreases again when the species becomes ubiquitous.  
135 Under the null model, species informativeness is:  
136  $\text{Informativeness} = N_2(1 - N_2/N)$ .  
137 where  $N$  is the total number of sites (Supplement A).  
138 In CA and its constrained forms—CCA and dc-CA (including FC)—species and site  
139 totals act as implicit weights. A simple and effective pre-processing approach is to  
140 divide species abundances by their total abundance and multiply by  $N_2(1 - N_2/N)$ ; so  
141 species totals reflect informativeness.  
142 This rationale also applies to sites: CWMs from species-poor sites are less precise  
143 than those from richer sites, suggesting weighting by the effective number of  
144 informative species. Down-weighting extremely species-rich sites ( $N_2$  close to the  
145 number of species) was considered but found to offer little advantage in the data sets  
146 and simulations presented here. Dividing site abundances by site totals and  
147 multiplying by  $N_2$  adjusts site weights, but then species totals no longer reflect  
148 informativeness. This catch-22 is resolved by iterative proportional fitting (IPF)  
149 (Fienberg 1970; Idel 2016): repeatedly applying both adjustments until species and  
150 site totals stabilize (Supplement B). Convergence is typically fast, yielding species  
151 and site totals approximately proportional to informativeness and effective numbers  
152 while preserving the interaction structure of the table. Informativeness and effective  
153 numbers are recomputed at each iteration.  
154 This  $N_2$ -processing procedure is implemented in R (R Core Team 2025) as the  
155 function *ipf2N2* in the package *douconca* (ter Braak & van Rossum 2025b).

## 156 **2.2 Evaluation**

157 The effectiveness of  $N_2$ -processing was evaluated using fourth-corner correlation  
158 analysis, RLQ, and dc-CA on both simulated datasets (ter Braak 2019; ter Braak  
159 2026) and on 80 CESTES datasets (Jeliazkov et al. 2020), summarized in Supplement  
160 C. Statistical evidence for trait–environment association was assessed through  $P$ -  
161 values (Muff *et al.* 2021) obtained from the permutational max test (ter Braak et al.  
162 2012; ter Braak 2026) with 999 random permutations.  $N_2$ -processing was applied with  
163 and without abundance transformations (*e.g.* logarithmic) to determine whether it can  
164 replace such transformations and whether it strengthens evidence for trait–  
165 environment association in datasets with less extreme abundance variation.

166 The simulated data sets were used to study power and Type I error rates. These were  
167 generated from generalized linear mixed models with random site-dependent trait and  
168 species-dependent environmental effects, known as GLMM3 (Niku *et al.* 2021), fitted  
169 to the Aravo and Revisit data in ter Braak (2019) and then modified by varying  
170 parameters, notably the strength of the trait–environment association(ter Braak 2019).

171 The 80 CESTES data sets were analyzed with and without pre-processing of  
172 abundance tables using FC, RLQ, and dc-CA. FC used the minimum  $P$ -value of all  
173 pairwise fourth-correlations between traits and environmental and was assessed with  
174 and without Bonferroni adjustment for multiple testing (FCadj and FC\_un).  
175 Differences between pre-processing and analysis methods were statistically evaluated  
176 using linear regression of evidence, expressed as  $-\log(P)$ , on the factors *method of*  
177 *analysis*,  *$N_2$ -processing method* (yes, no) and *type of abundance transformation*, with

178 covariate factor *CESTES data set*. Method of analysis consisted of four levels (FC\_un,  
 179 FCadj, RLQ and dc-CA). The types of abundance transformation were also quantified  
 180 as ‘none’ =0, ‘square-root’= 1, ‘log’ = 2 and ‘P/A’ =4, to more concisely describe the  
 181 interaction between  $N_2$ -processing method (yes, no) and type of abundance  
 182 transformation.

183 The logarithmic transformation of abundance was  $\log(y/s+1)$ , with  $y$  abundance and  $s$   
 184 the smallest non-zero value, usually  $s = 1$  for count data. Non-negative right-skewed  
 185 traits and environmental variables were log-transformed. When the number of traits or  
 186 environmental variables was too high relative to species or sites (fewer than 10  
 187 residual degrees of freedom), the number of variables used in dc-CA was reduced by  
 188 selecting variables via CCA. This selection avoided collinearity and discarded  
 189 variables least related to the pre-processed abundance table.

190 In this paper, FC and RLQ were fitted to data using the R packages *ade4* v1.7-23  
 191 (Thioulouse *et al.* 2018) and dc-CA using *douconca* v1.2.5 (ter Braak & van Rossum  
 192 2025a). In RLQ, the Hill-Smith option was used to deal with categorical trait and  
 193 environmental variables; in dc-CA, abundances were *not* divided by the site totals.

### 194 3. Results

195  $N_2$ -processing of the abundance table increased the power to detect trait–environment  
 196 associations in both Aravo- and Revisit-simulations (Fig. 1). The increase was largest  
 197 without any prior transformation and smallest for prior transformation to presence-  
 198 absence data. With  $N_2$ -processing, all prior transformations led to a similar power.  
 199 Without  $N_2$ -processing, there was slight inflation of the Type I error rate when the  
 200 variance of the site-dependent trait effect was increased relative to its fitted value  
 201 (Fig. 1). The simulations contained a single trait and single environmental variable  
 202 only, and therefore apply to both dc-CA and RLQ, as well as fourth-corner correlation  
 203 analysis.

204 The median number of IPF iterations was 76, the maximum 8015. Only 0.3% of the  
 205 data sets and their transformations required more than 1000 iterations.

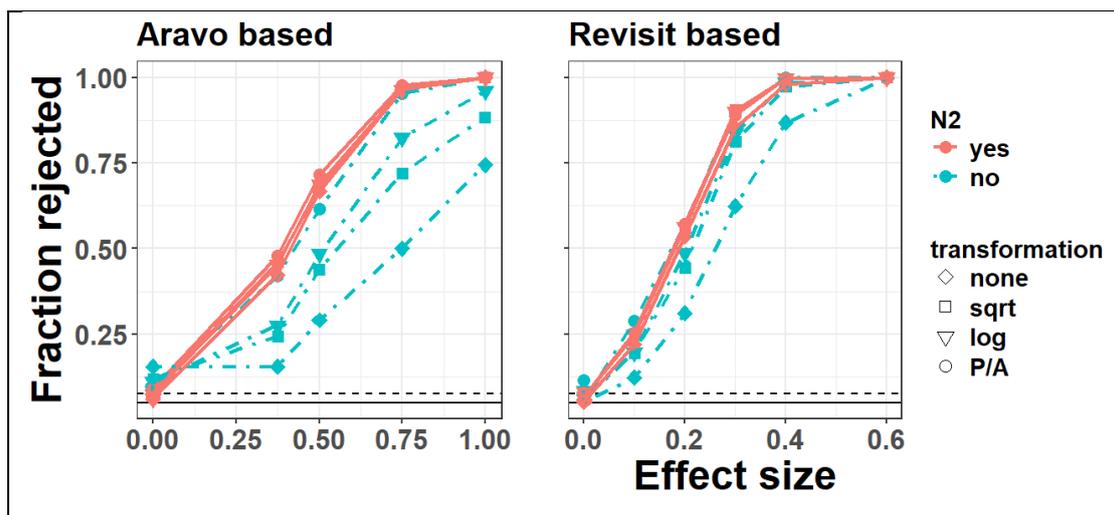


FIG. 1. Rejection rates of pre-processing methods against the size of the trait-environment association (Effect size), based on the simulations of GLMM3 models for the Aravo and Revisit data in ter Braak (2019). Pre-precession methods are transformation of abundance (none, sqrt: square-root, log: logarithmic and P/A: transformation to presence/absence) with and without subsequent  $N_2$ -processing

( $N_2$ ). Additional simulations with increased trait variance were carried out, as in ter Braak (2019), for effect size zero, and the lines start at this scenario so as to emphasize potential Type I error rate inflation. The horizontal solid line is at the nominal significance threshold; rates above the dashed line (at 0.078) are significantly greater than 0.05.

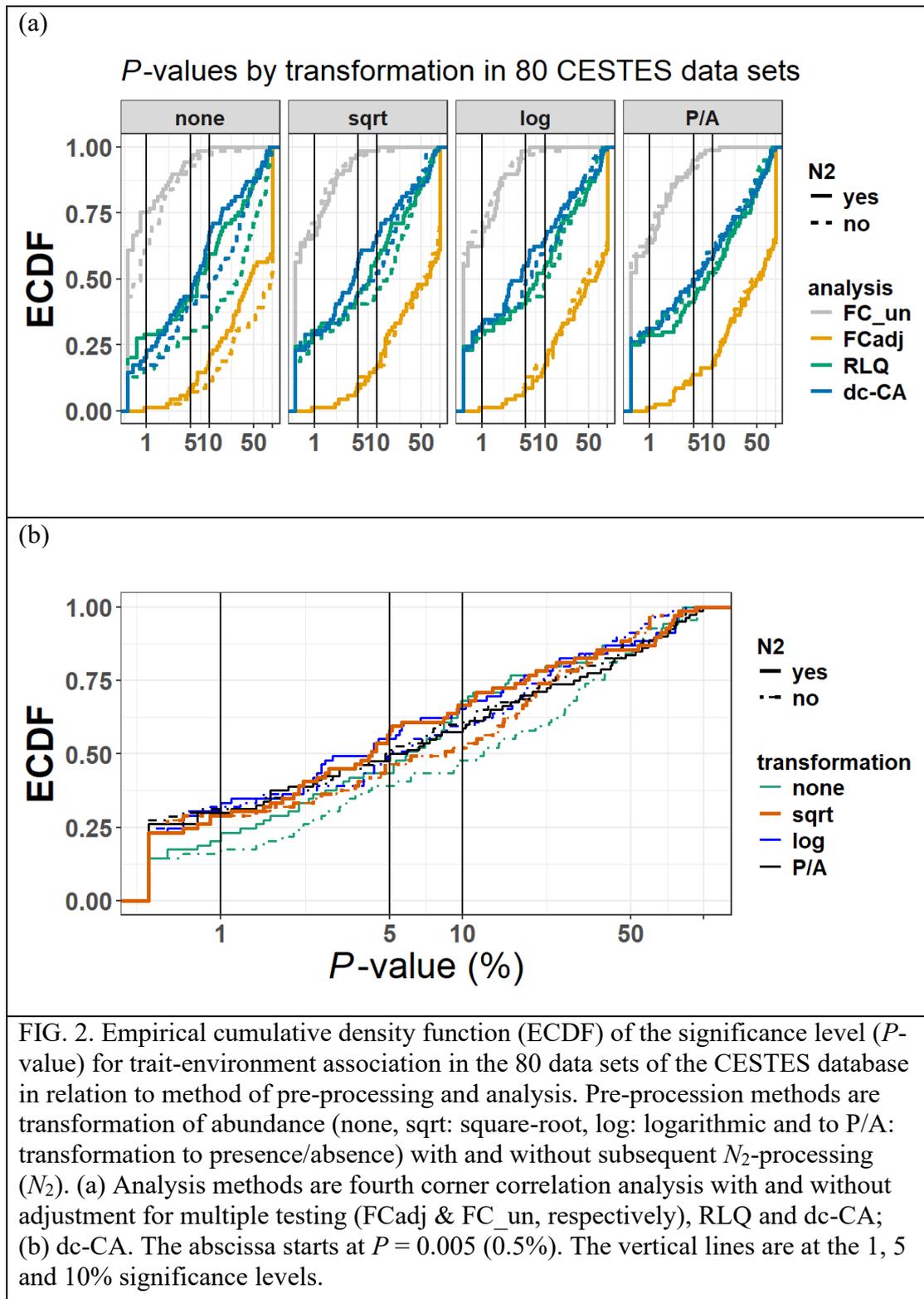
206

207 Evidence for trait-environment association in the 80 CESTES data sets (Fig. 2a)  
208 depended more on the analysis method than on the method of pre-processing  
209 (abundance transformation with and without subsequent  $N_2$ -processing) and whether  
210 FC was carried out with or without adjustment for multiple testing. Without  
211 adjustment, all data sets had at least one pair of trait and environmental variable that  
212 showed at least weak evidence ( $P < 0.10$ ) for trait-environment association,  
213 irrespective of the type of pre-processing, as the grey lines in all panels of Figure 2a  
214 reach 1 at the vertical line for the 10% significance level. There was one exception, a  
215 data set that showed weak significance with  $N_2$ -processing, but no evidence without  
216 pre-processing. However, adjustment for multiple testing is called for as the data sets  
217 typically contained many traits and environmental variables. The adjusted minimum  
218 P-value was below 0.10 only in 16-22% of the 80 data sets, as the yellow lines cross  
219 the vertical line for the 10% significance level at this range of values.

220 Without  $N_2$ -processing, the percentage of data sets showing at least weak trait-  
221 environment association ( $P < 0.10$ ) was much higher using RLQ or dc-CA: between  
222 38 and 51% for RLQ and 50-60 for dc-CA. With  $N_2$ -processing these percentages  
223 further increased to 52-60% for RLQ and 57-66% for dc-CA, with prior square-root  
224 and log transformation yielding the highest percentages at the 5% and 1% significance  
225 levels. The benefit of  $N_2$ -processing was most pronounced for untransformed data  
226 and disappeared once the data were converted to presence-absence (Fig. 2).

227 The percentage of data sets showing strong trait-environment association ( $P < 0.01$ )  
228 was between 18 and 29% for RLQ and 21-35% for dc-CA without  $N_2$ -processing.  
229 With  $N_2$ -processing, these percentages were similar for both RLQ and dc-CA, except  
230 for 'no prior transformation' where  $N_2$ -processing increased the percentage data sets  
231 with strong evidence in RLQ from 18 to 29% and in dc-CA from 21% to 24% (Fig.  
232 2a).

233 Within-data-set comparisons using linear regression indicated that the main difference  
234 in  $-\log(P)$  was between no processing and any pre-processing method; differences  
235 among the pre-processing methods themselves were minimal, in RLQ, dc-CA and  
236 fourth-corner correlation. Overall, dc-CA showed stronger evidence than RLQ ( $P <$   
237 0.01), while fourth-corner analysis with Bonferroni adjustment was much worse than  
238 both dc-CA and RLQ. There was no statistical evidence for interaction between  
239 *method of analysis* and *method of pre-processing* ( $P = 0.68$ ).



240

241 Focusing on dc-CA (Fig. 2b), there was strong evidence for effects of both abundance  
 242 transformation on  $-\log(P)$ , with the log transformation performing best and the  
 243 square-root transformation as close runner-up, and  $N_2$ -processing ( $P < 0.001$ ). A  
 244 significant interaction ( $P < 0.01$ ) indicated that the effectiveness of  $N_2$ -processing  
 245 diminished as the strength of the abundance transformation increased from none to

246 square-root, to log, with minor positive or even a small adverse effect after  
247 transformation to presence–absence (Fig. 2b; Supplement C).  
248 The median number of IPF iterations was 76, the maximum 5646 and 2% of the data  
249 sets and their transformations required more than 1000 iterations.  
250 Supplement C presents results for  $N_2$ -variants, explores additional power  
251 transformations, and further dissects the CESTES database by type of abundance  
252 measurement, ecosystem, and taxonomic group.

## 253 4. Discussion

254 As shown here, the precision of both CWM and SNC is proportional to the Hill's  
255 effective number  $N_2$ , at least when  $N_2$  is far below its maximum value. This motivates  
256 the use  $N_2$  or its modification *informativeness*, as weights in statistical analyses based  
257 on CWMs and SNCs, such as FC, RLQ and dc-CA. However, weighting species by  
258 their informativeness alters the site totals, and thus the implied site weights. This  
259 catch-22 is resolved by iterative proportional fitting (IPF) which conserves the  
260 species-site interaction structure of the table while adjusting the relative species and  
261 site totals to match, as closely as possible, their target informativeness and effective  
262 numbers. The species and site totals of the transformed table represent  
263 informativeness and effective number of informative species, respectively.  $N_2$ -  
264 processing replaces ad-hoc preprocessing choices with a principled, model-based one.

265 The proposed method increased the power of the max test for detecting trait–  
266 environment associations in simulations, and strengthened statistical evidence for  
267 such associations in the 80 datasets of the CESTES global database of traits,  
268 environments, and abundances. Its effect diminished when extreme abundances are  
269 down-weighted, as occurs with log-transformations and, even more strongly, with the  
270 transformation to presence–absence.

271 A key limitation of drawing general conclusions from analyses of simulated and real  
272 datasets, as compared with theory, is that it is virtually impossible to encompass all  
273 relevant scenarios and data set properties. Having said this, based on the analyses of  
274 the simulated and real datasets, a square-root transformation followed by  
275  $N_2$ -processing is recommended, offering high effectiveness while avoiding the  
276 arbitrary zero replacement required for log transforms.  $N_2$ -processing is essential; the  
277 square-root step can be replaced by other mild power transforms without loss of  
278 efficiency.  $N_2$ -processing is essential because it provides the correct model-based  
279 variance scaling, reduces the influence of other arbitrary data-analytic choices and  
280 maximizes the statistical power and validity of dc-CA across datasets.

281 Lower P-values in the analysis of the CESTES datasets are desirable, provided that  
282 genuine trait–environment associations are present in each dataset. This assumption is  
283 reasonable because every dataset in the CESTES compilation originates from a  
284 published, peer-reviewed study, implying that the underlying trait–environment  
285 relationships were considered scientifically credible by the original authors. Indeed,  
286 each dataset contained at least one trait–environment pair showing weak statistical  
287 evidence ( $P < 0.10$ ) for association.

288 However, without prior indication—ideally via pre-registration—of the primary trait–  
289 environment pair, multiple-testing correction is required. After such adjustment, only  
290 13% of the datasets retained weak to strong evidence for trait–environment  
291 association. Given the typically large number of traits and environmental variables, a

292 multivariate approach is more appropriate than conducting many univariate tests.  
293 Methods such as RLQ and dc-CA are specifically designed to address this issue.  
294 When dc-CA was combined with  $N_2$ -processing, the proportion of datasets showing at  
295 least weak evidence for trait–environment association increased markedly, from 13%  
296 using FC to 68% using dc-CA with  $N_2$ -processing, demonstrating the substantial gain  
297 in sensitivity provided by the proposed method.

298 In practice, species abundance is often more strongly related to environmental  
299 variables than to traits. Sampling designs typically represent one or two dominant  
300 environmental gradients, while species occurrence within the chosen taxonomic group  
301 cannot be controlled (Gobbi *et al.* 2022). Hence, CWM regression along these  
302 gradients often shows stronger evidence for trait–environment associations than SNC  
303 regression. Note that a species-level analysis such as SNC-regression is missing in  
304 many legacy trait-environment analyses.

305 Whereas a presence–absence table records occurrences, the entries in the  
306  $N_2$ -transformed table can be interpreted as *effective occurrences*. In contrast to  
307 presence–absence data, the transformed table retains quantitative information while  
308 down-weighting extreme abundances. Predicted effective occurrences can therefore  
309 be viewed as analogous to presence probabilities. An advantage of  $N_2$ -processing over  
310 conversion to presence–absence is that it preserves quantitative information while  
311 down-weighting common, uninformative species.

312 The improvement in statistical evidence for trait–environment association is greatest  
313 when the underlying associations are weak—arguably the most common situation in  
314 ecological datasets. In such cases, strengthening the analysis through dc-CA  
315 combined with  $N_2$ -processing becomes particularly important, as it can reveal subtle  
316 relationships that might otherwise remain undetected.

317 **Supplement A:  $N_2$  and a species’ informativeness**

318 **Supplement B: Iterative proportional fitting (IPF)**

319 **Supplement C: Real and simulated data sets with additional results**

320

321

322 **Acknowledgements**

323 I thank xxxxx for their valuable comments that helped to improve the manuscript. By  
324 prompting Microsoft Copilot (Version January 6,  
325 2026, <https://www.microsoft.com/copilot>) with ‘Improve: [my English text]’, I  
326 iteratively refined the language and style of sentences and paragraphs in this paper.  
327 All views, insights, and errors are mine.

328 **Data availability**

329 No new data have been generated. The R code is available at  
330 <https://doi.org/10.6084/m9.figshare.31398048>

331 **CRedit authorship contribution statement**

332 CB is the only author: Writing – review and editing, Writing – original draft,  
333 Visualization, Validation, Supervision, Software, Methodology, Formal analysis,  
334 Conceptualization.

335 **Competing interests and Funding**

336 The author declares no conflicts of interest. No funds, grants, or other support was  
337 received.

338

339 **5. References**

- 340 Anderegg, L.D.L. (2023) Why can't we predict traits from the environment? *New*  
341 *Phytologist*, **237**, 1998-2004. <https://doi.org/10.1111/nph.18586>
- 342 Broennimann, O., Collart, F. & Guisan, A. (2026) The ecospat R Package: A  
343 Collection of Pre-, Core-, and Post-Modeling Tools to Investigate Species  
344 Niches and Distributions. *R Coding for Ecology* (ed. D. Rocchini), pp. 67-98.  
345 Springer Nature Switzerland, Cham. [https://doi.org/10.1007/978-3-031-99665-](https://doi.org/10.1007/978-3-031-99665-8_3)  
346 [8\\_3](https://doi.org/10.1007/978-3-031-99665-8_3)
- 347 de Bello, F., Fischer, F.M., Puy, J., Shipley, B., Verdú, M., Götzenberger, L., Lavorel,  
348 S., Moretti, M., Wright, I.J., Berg, M.P., Carmona, C.P., Cornelissen, J.H.C.,  
349 Dias, A.T.C., Gibb, H., Lepš, J., Madin, J.S., Majeková, M., Pausas, J.G.,  
350 Segrestin, J., Sobral, M., Zanne, A.E. & Garnier, E. (2025) Raunkiæran  
351 shortfalls: Challenges and perspectives in trait-based ecology. *Ecological*  
352 *Monographs*, **95**, e70018.  
353 <https://esajournals.onlinelibrary.wiley.com/doi/abs/10.1002/ecm.70018>
- 354 Dolédec, S., Chessel, D., ter Braak, C.J.F. & Champely, S. (1996) Matching species  
355 traits to environmental variables: a new three-table ordination method.  
356 *Environmental and Ecological Statistics*, **3**, 143-166.  
357 <http://dx.doi.org/10.1007/BF02427859>
- 358 Dray, S., Choler, P., Dolédec, S., Peres-Neto, P.R., Thuiller, W., Pavoine, S. & ter  
359 Braak, C.J.F. (2014) Combining the fourth-corner and the RLQ methods for  
360 assessing trait responses to environmental variation. *Ecology*, **95**, 14-21.  
361 <http://dx.doi.org/10.1890/13-0196.1>
- 362 Dray, S. & Legendre, P. (2008) Testing the species traits environment relationships:  
363 The fourth-corner problem revisited. *Ecology*, **89**, 3400-3412.  
364 <https://doi.org/10.1890/08-0349.1>
- 365 Fienberg, S.E. (1970) An Iterative Procedure for Estimation in Contingency Tables.  
366 *The Annals of Mathematical Statistics*, **41.3**, 907-917.  
367 <https://doi.org/10.1214/aoms/1177696968>
- 368 Garnier, E., Navas, M.-L., Grigulis, K., Garnier, E., Navas, M.-L. & Grigulis, K.  
369 (2015) Trait-based ecology: definitions, methods, and a conceptual  
370 framework. *Plant Functional Diversity: Organism traits, community structure,*  
371 *and ecosystem properties*, pp. 0. Oxford University Press.  
372 <https://doi.org/10.1093/acprof:oso/9780198757368.003.0002>
- 373 Gobbi, M., Corlatti, L., Caccianiga, M., ter Braak, C.J.F. & Pedrotti, L. (2022) Hay  
374 meadows' overriding effect shapes ground beetle functional diversity in  
375 mountainous landscapes. *Ecosphere*, **e4193**. <https://doi.org/10.1002/ecs2.4193>
- 376 Green, S.J., Brookson, C.B., Hardy, N.A. & Crowder, L.B. (2022) Trait-based  
377 approaches to global change ecology: moving from description to prediction.  
378 *Proceedings of the Royal Society B: Biological Sciences*, **289**, 20220071.  
379 <https://doi.org/10.1098/rspb.2022.0071>
- 380 Grime, J.P. (1973) Competitive Exclusion in Herbaceous Vegetation. *Nature*, **242**,  
381 344-347. <https://doi.org/10.1038/242344a0>

- 382 Hill, M.O. (1973a) Diversity and evenness: a unifying notation and its consequences.  
 383 *Ecology*, **54**, 427-432. <https://doi.org/10.2307/1934352>
- 384 Hill, M.O. (1973b) Reciprocal averaging: an eigenvector method of ordination.  
 385 *Journal of Ecology*, **61**, 237-249. <https://doi.org/10.2307/2258931>
- 386 Idel, M. (2016) A review of matrix scaling and Sinkhorn's normal form for matrices  
 387 and positive maps. <https://arxiv.org/pdf/1609.06349.pdf>
- 388 Jeliaskov, A., Mijatovic, D., Chantepie, S., Andrew, N., Arlettaz, R., Barbaro, L.,  
 389 Barsoum, N., Bartonova, A., Belskaya, E., Bonada, N., Brind'Amour, A.,  
 390 Carvalho, R., Castro, H., Chmura, D., Choler, P., Chong-Seng, K., Cleary, D.,  
 391 Cormont, A., Cornwell, W., de Campos, R., de Voogd, N., Doledéc, S., Drew, J.,  
 392 Dziock, F., Eallonardo, A., Edgar, M.J., Farneda, F., Hernandez, D.F.,  
 393 Frenette-Dussault, C., Fried, G., Gallardo, B., Gibb, H., Gonçalves-Souza, T.,  
 394 Higuti, J., Humbert, J.-Y., Krasnov, B.R., Saux, E.L., Lindo, Z., Lopez-  
 395 Baucells, A., Lowe, E., Marteinsdottir, B., Martens, K., Meffert, P., Mellado-  
 396 Díaz, A., Menz, M.H.M., Meyer, C.F.J., Miranda, J.R., Mouillot, D., Ossola,  
 397 A., Pakeman, R., Pavoine, S., Pekin, B., Pino, J., Pocheville, A., Pomati, F.,  
 398 Poschlod, P., Prentice, H.C., Purschke, O., Ravel, V., Reitalu, T., Renema,  
 399 W., Ribera, I., Robinson, N., Robroek, B., Rocha, R., Shieh, S.-H., Spake, R.,  
 400 Staniaszek-Kik, M., Stanko, M., Tejerina-Garro, F.L., Braak, C.t., Urban,  
 401 M.C., Klink, R.v., Villéger, S., Wegman, R., Westgate, M.J., Wolff, J.,  
 402 Żarnowiec, J., Zolotarev, M. & Chase, J.M. (2020) A global database for  
 403 metacommunity ecology, integrating species, traits, environment and space.  
 404 *Scientific Data*, **7**, 6. <https://doi.org/10.1038/s41597-019-0344-7>
- 405 Legendre, P., Galzin, R.G. & Harmelin-Vivien, M.L. (1997) Relating behavior to  
 406 habitat: Solutions to the fourth-corner problem. *Ecology*, **78**, 547-562.  
 407 <https://doi.org/10.2307/2266029>
- 408 Lepš, J. & de Bello, F. (2023) Differences in trait–environment relationships:  
 409 Implications for community weighted means tests. *Journal of Ecology*, **111**,  
 410 2328-2341. <https://doi.org/10.1111/1365-2745.14172>
- 411 McCune, B. (2015) The front door to the fourth corner: variations on the sample unit  
 412 × trait matrix in community ecology. *Community Ecology*, **16**, 267-271.  
 413 <https://doi.org/10.1556/168.2015.16.2.14>
- 414 Muff, S., Nilsen, E.B., O'Hara, R.B. & Nater, C.R. (2021) Rewriting results sections  
 415 in the language of evidence. *Trends in Ecology & Evolution*.  
 416 <https://doi.org/10.1016/j.tree.2021.10.009>
- 417 Niku, J., Hui, F.K.C., Taskinen, S. & Warton, D.I. (2021) Analyzing environmental-  
 418 trait interactions in ecological communities with fourth-corner latent variable  
 419 models. *Environmetrics*, **n/a**, e2683. <https://doi.org/10.1002/env.2683>
- 420 Peres-Neto, P.R., Dray, S. & ter Braak, C.J.F. (2017) Linking trait variation to the  
 421 environment: critical issues with community-weighted mean correlation  
 422 resolved by the fourth-corner approach. *Ecography*, **40**, 806-816.  
 423 <http://dx.doi.org/10.1111/ecog.02302>
- 424 Pinho, B.X., Melo, F.P.L., ter Braak, C.J.F., Bauman, D., Maréchaux, I., Tabarelli,  
 425 M., Benchimol, M., Arroyo-Rodriguez, V., Santos, B.A., Hawes, J.E.,  
 426 Berenguer, E., Ferreira, J., Silveira, J.M., Peres, C.A., Rocha-Santos, L.,  
 427 Souza, F.C., Gonçalves-Souza, T., Mariano-Neto, E., Faria, D. & Barlow, J.  
 428 (2024) Winner–loser plant trait replacements in human-modified tropical  
 429 forests. *Nature Ecology & Evolution*. [https://doi.org/10.1038/s41559-024-](https://doi.org/10.1038/s41559-024-02592-5)  
 430 [02592-5](https://doi.org/10.1038/s41559-024-02592-5)

431 R Core Team (2025) *R: A language and environment for statistical computing*,  
432 *version 4.1*. R Foundation for Statistical Computing. [www.R-project.org](http://www.R-project.org),  
433 Vienna, Austria

434 ter Braak, C.J.F. (2017) Fourth-corner correlation is a score test statistic in a log-  
435 linear trait–environment model that is useful in permutation testing.  
436 *Environmental and Ecological Statistics*, **24**, 219-242.  
437 <http://dx.doi.org/10.1007/s10651-017-0368-0>

438 ter Braak, C.J.F. (2019) New robust weighted averaging- and model-based methods  
439 for assessing trait–environment relationships. *Methods in Ecology and*  
440 *Evolution*, **10**, 1962-1971. <https://doi.org/10.1111/2041-210X.13278>

441 ter Braak, C.J.F. (2026) Fourth-corner latent variable models overstate confidence in  
442 trait–environment relationships and what to use instead. *Environmental and*  
443 *Ecological Statistics*. <https://doi.org/10.1007/s10651-025-00696-0>;  
444 <https://edepot.wur.nl/707477>.

445 ter Braak, C.J.F., Cormont, A. & Dray, S. (2012) Improved testing of species traits–  
446 environment relationships in the fourth-corner problem. *Ecology*, **93**, 1525-  
447 1526. <https://doi.org/10.1890/12-0126.1>

448 ter Braak, C.J.F., Peres-Neto, P.R. & Dray, S. (2018) Simple parametric tests for  
449 trait–environment association. *Journal of Vegetation Science*, **29**, 801-811.  
450 <https://doi.org/10.1111/jvs.12666>

451 ter Braak, C.J.F., Šmilauer, P. & Dray, S. (2018) Algorithms and biplots for double  
452 constrained correspondence analysis. *Environmental and Ecological Statistics*,  
453 **25**, 171-197. <https://doi.org/10.1007/s10651-017-0395-x>

454 ter Braak, C.J.F. & van Rossum, B.-J. (2025a) Linking multivariate trait variation to  
455 the environment: the advantages of double constrained correspondence  
456 analysis with the R package douconca. *Ecological Informatics*, **88**, 103143.  
457 <https://doi.org/10.1016/j.ecoinf.2025.103143>

458 ter Braak, C.J.F. & van Rossum, B.-J. (2025b) R package douconca: Double  
459 constrained correspondence analysis for multi-trait multi-environment analysis  
460 v1.2.5. <https://doi.org/10.32614/CRAN.package.douconca>

461 Thioulouse, J., Dray, S., Dufour, A.-B., Siberchicot, A., Jombart, T. & Pavoine, S.  
462 (2018) *Multivariate Analysis of Ecological Data with ade4*. Springer New  
463 York, New York, NY.978-1-4939-8850-1

464 Westoby, M. (2025) Trait-based ecology, trait-free ecology, and in between. *The New*  
465 *Phytologist*, **245**, 33-39. <http://dx.doi.org/10.1111/nph.20197>

466

## 467 **6. Supplement A: $N_2$ and a species' informativeness**

468 This supplement shows the importance of Hill number of order 2,  $N_2$ , in weighted  
469 averaging of traits (CWM) and environmental variables (SNC). Here it shown that the  
470 variance of an SNC or CWM under a permutational null model is inversely  
471 proportional to  $N_2$ , at least when  $N_2$  is small compared to the number of units  
472 (number of species for CWM or number of sites for SNC).  $N_2$  can thus be interpreted  
473 the effective number of occurrences or effective number of species in a CWM- and  
474 SNC-analysis, respectively A measure of informativeness of species is also defined to  
475 cover the case of when its effective number of occurrences, its  $N_2$ , is not small, e.g.  
476 when a species occurs equally in nearly all sites.

477 The derivation in section 1.1 is also of some interest for winner-loser analysis  
478 (Tabarelli, Peres & Melo 2012; Chisté *et al.* 2016; Filgueiras *et al.* 2021; García  
479 Criado *et al.* 2023) as it derives the variance of an SNC under permutation of the  
480 environmental variables across sites.

### 481 **6.1 Mean and variance of SNC and CWM under a permutational null** 482 **model**

483 Being weighted averages, the SNC and the CWM can be expressed as a linear  
484 combination  $z = \sum_i p_i x_i$ , provided  $p_i$  is an abundance fraction such that  $\sum_i p_i = 1$ . In  
485 the case of an SNC,  $x_i$  is the value of a particular environmental variable or  
486 environmental gradient  $x$  at site  $i$ , and  $p_i$  is the abundance of a particular species at  
487 site  $i$  divided by the total abundance of the species. In the case of a CWM,  $x_i$  is the  
488 value of a particular trait  $x$  of species  $i$ , and  $p_i$  is the abundance of species  $i$  at a  
489 particular site divided by the total abundance at the site.

490 The expectation of  $z$  under random permutation of the values  $\{x_1, x_2, \dots, x_n\}$  is

$$491 E(z) = E(\sum_i p_i x_i) = \sum_i p_i E(x_i) = \sum_i p_i \bar{x} = \bar{x} \sum_i p_i = \bar{x}, \quad (\text{A1})$$

492 where  $\bar{x} = \frac{1}{n} \sum_i x_i$ , the average of the  $x$ -values, as  $\sum_i p_i = 1$  and  $E(x_i) = \bar{x}$  under  
493 random permutation. This formula shows that the expected value of SNC under  
494 random permutation is the same for all species, irrespective of their total or number of  
495 occurrences.

496 The general formula for the variance of a linear combination is (Rao 1973)

$$497 \text{var}(z) = \text{var}(\mathbf{p}'\mathbf{x}) = \mathbf{p}' \text{cov}(\mathbf{x}) \mathbf{p}, \quad (\text{A2})$$

498 where  $\mathbf{p}' = (p_1, p_2, \dots, p_n)$  and  $\text{cov}(\mathbf{x})$  is the covariance matrix of  $x_1, x_2, \dots, x_n$  under  
499 random permutation. The diagonal elements are each  $s_x^2$ , the variance of  $x$  with  
500 denominator  $n$  instead of  $n - 1$  (intuitively, the sample  $x_1, x_2, \dots, x_n$  is, in the  
501 permutations, the population of interest). The covariance between the elements  $x_i$  and  
502  $x_j$  ( $i \neq j$ ) is a constant  $c$ , say, because of the symmetry of permutations. The constant  
503 can be obtained by noting that the variance of the sum of all  $x$ -values is 0 under  
504 permutation and can be expressed as

$$505 \text{var}(x_1 + x_2 + \dots + x_n) = \sum_i \text{var}(x_i) + 2 \sum_{i < j} \text{cov}(x_i, x_j) = n s_x^2 + n(n - 1)c.$$

506 Solving for  $c$  gives  $c = -s_x^2/(n - 1)$ , which completes the specification of  $\text{cov}(\mathbf{x})$ .  
507 Insertion in the general formula (A2), yields

508  $var(z) = s_x^2 \sum_i \left( p_i^2 - \frac{1}{n-1} p_i \sum_{j \neq i} p_j \right),$  (A3)

509 which can be simplified to

510  $var(z) = \frac{s_x^2}{N_2} \left( 1 - \frac{N_2-1}{n-1} \right),$  (A4)

511 where  $N_2 = \frac{1}{\sum_i p_i^2}$ , the Hill number of order 2.

512 The first term in (A4) is variance divided by  $N_2$ , which resembles the usual formula  
 513 for the variance of a mean: the population variance divided by the number of units  
 514 drawn.  $N_2$  can thus be interpreted as the effective number of units in a CWM- or  
 515 SNC-analysis. When  $N_2$  is small compared to  $n$  (e.g. sites are relatively species-poor  
 516 or species are relatively rare in the data set) the second term is small. On neglecting  
 517 the second term, the precision, *i.e.* the inverse of the variance, is proportional to the  
 518 effective number  $N_2$ .

519 The variance  $s_x^2 \left( 1 - \frac{N_2-1}{n-1} \right)$  can be interpreted as the variance that is left on average  
 520 after simple random sampling of  $N_2$  distinct units without replacement from a  
 521 population of  $n$  units (Cochran 1977). Consider  $N_2 = 1$  and  $N_2 = n$ , then the formula  
 522 says that the variance under permutation is  $s_x^2$  and 0, respectively. In the first case a  
 523 single  $x$ -value is randomly selected so that the variance is that of  $x$  ( $s_x^2$ ), and in the  
 524 second case all values of  $x$  are selected with equal weight ( $p_i = 1/n$ ) so that the SNC  
 525 or CWM is equal to  $\bar{x}$  and is constant under permutation, so that the permutational  
 526 variance is 0, in agreement with the formula. A species for which  $N_2 = n$  is thus  
 527 uninformative. On replacing  $n-1$  by  $n$  in the term  $1 - \frac{N_2-1}{n-1}$ , the term is approximately  
 528 equal to  $1 - \frac{N_2}{n}$ , as  $n$  is large in practice. A simple measure of a species'  
 529 informativeness under the null model is therefore  $N_2(1 - N_2/n)$ .

530

## 531 6.2 Conclusion

532 Formula (A4) gives a definite interpretation of the Hill number  $N_2$  as the effective  
 533 number of occurrences or effective number of species for SNC and CWM,  
 534 respectively. It the Hill number of order 2 ( $N_2$ ) that plays a special role in weighted  
 535 averaging, instead of any order, as Lepš and de Bello (2023) appear to suggest.

536 The simple measure of the informativeness of a species under the null model is  
 537  $N_2(1 - N_2/n)$ .

538

## 539 References

540

- 541 Chisté, M.N., Mody, K., Gossner, M.M., Simons, N.K., Köhler, G., Weisser, W.W. &  
 542 Blüthgen, N. (2016) Losers, winners, and opportunists: How grassland land-  
 543 use intensity affects orthopteran communities. *Ecosphere*, **7**, e01545.  
 544 <https://esajournals.onlinelibrary.wiley.com/doi/abs/10.1002/ecs2.1545>  
 545 Cochran, W.G. (1977) *Sampling Techniques*, 3rd ed. John Wiley & Sons, New York.  
 546 978-0-471-16240-7

- 547 Filgueiras, B.K.C., Peres, C.A., Melo, F.P.L., Leal, I.R. & Tabarelli, M. (2021)  
 548 Winner–Loser Species Replacements in Human-Modified Landscapes. *Trends*  
 549 *in Ecology & Evolution*, **36**, 545-555.  
 550 <https://www.sciencedirect.com/science/article/pii/S0169534721000562>
- 551 García Criado, M., Myers-Smith, I.H., Bjorkman, A.D., Normand, S., Blach-  
 552 Overgaard, A., Thomas, H.J.D., Eskelinen, A., Happonen, K., Alatalo, J.M.,  
 553 Anadon-Rosell, A., Aubin, I., te Beest, M., Betway-May, K.R., Blok, D.,  
 554 Buras, A., Cerabolini, B.E.L., Christie, K., Cornelissen, J.H.C., Forbes, B.C.,  
 555 Frei, E.R., Grogan, P., Hermanutz, L., Hollister, R.D., Hudson, J., Iturrate-  
 556 Garcia, M., Kaarlejärvi, E., Kleyer, M., Lamarque, L.J., Lembrechts, J.J.,  
 557 Lévesque, E., Luoto, M., Macek, P., May, J.L., Prevéy, J.S., Schaepman-  
 558 Strub, G., Sheremetiev, S.N., Siegwart Collier, L., Soudzilovskaia, N.A.,  
 559 Trant, A., Venn, S.E. & Virkkala, A.-M. (2023) Plant traits poorly predict  
 560 winner and loser shrub species in a warming tundra biome. *Nature*  
 561 *Communications*, **14**, 3837. <https://doi.org/10.1038/s41467-023-39573-4>
- 562 Lepš, J. & de Bello, F. (2023) Differences in trait–environment relationships:  
 563 Implications for community weighted means tests. *Journal of Ecology*, **111**,  
 564 2328-2341. <https://doi.org/10.1111/1365-2745.14172>
- 565 Rao, C.R. (1973) *Linear statistical inference and its application*. 2nd ed. Wiley, New  
 566 York. 9780470316436
- 567 Tabarelli, M., Peres, C.A. & Melo, F.P.L. (2012) The ‘few winners and many losers’  
 568 paradigm revisited: Emerging prospects for tropical forest biodiversity.  
 569 *Biological Conservation*, **155**, 136-140.  
 570 <https://www.sciencedirect.com/science/article/pii/S0006320712002893>  
 571

## 572 **7. Supplement B: Iterative proportional fitting (IPF)**

573 This supplement described the  $N_2$ -processing method in detail. Section 7.1 describe  
 574 the standard iterative proportional fitting (IPF) procedure. In dc-CA, RLQ and FC, the  
 575 marginals of the obtained table are the implied weights. A disadvantage of the  
 576 standard procedure for use in dc-CA, RLQ and FC is, that the marginals of the  
 577 obtained table and thus the implied weights to not reflect the effective or informative  
 578 numbers of the obtained table. In order to obtain marginals that reflect effective or  
 579 informative numbers these numbers need to be recomputed (updated) in each step of  
 580 the IPF procedure. This is detailed in section 0, while the details of the stopping  
 581 criterion are described in section 1.3.

582 *Notation:*  $n_{ij}$  is the abundance of the  $j^{\text{th}}$  out of  $m$  species in the  $i^{\text{th}}$  out of  $n$  sites ( $n_{ij} \geq 0$ ).  
 583 The  $n \times m$  abundance table thus consist of non-negative entries denoted by  $n_{ij}$  and after  
 584 transformation or during or after IPF. A ‘+’ replacing an index denotes summation  
 585 over that index, e.g.  $y_{i+}$  and  $y_{+j}$  denote the total for the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column,  
 586 respectively.

### 587 **7.1 IPF of a two-way table**

588 The original Deming and Stephan IPF procedure (Fienberg 1970) aims to modify a  
 589 two-way to specified row and column marginals with entries  $r_i$  and  $k_j$ , respectively. If  
 590  $r_+ \neq k_+$ , calculate  $c = k_+ / r_+$  and rescale  $r_i \leftarrow c \times \text{old } r_i$  such that new  $r_+ = k_+$ .

591 Initialization:  $y_{ij} \leftarrow n_{ij}$

592 Iterate until to the values stabilize:  
 593 Step 1. For each column  $j$  execute:  $y_{ij} \leftarrow y_{ij} (k_j / y_{+j})$   
 594 Step 2. For each row  $i$  execute:  $y_{ij} \leftarrow y_{ij} (r_i / y_{i+})$   
 595 Step 3. Goto step 1 or stop if the values of the current iteration are sufficiently close to  
 596 those of the previous iteration.

597 This IPF procedure can be applied to an abundance table by first calculating the  
 598 effective numbers  $N_{2i} = \frac{1}{\sum_j p_{ij}^2}$ , with  $p_{ij} = \frac{y_{ij}}{y_{i+}}$  and  $\tilde{N}_{2j} = \frac{1}{\sum_i q_{ij}^2}$ , with  $q_{ij} = \frac{y_{ij}}{y_{+j}}$ , and  
 599 informative numbers  $I_i = N_{2i} (1 - N_{2i}/m)$  and  $\tilde{I}_j = \tilde{N}_{2j} (1 - \tilde{N}_{2j}/n)$

600 then setting the marginals to either the effective number or the informative numbers  
 601 ( $r_i = N_{2i}$  or  $r_i = I_i$  and  $k_j = \tilde{N}_{2j}$  or  $k_j = \tilde{I}_j$ ) with subsequent rescaling such that  $r_+ =$   
 602  $k_+$ .

603 Executing Step 1 only with  $k_j = \tilde{I}_j$  without further iteration is the processing method  
 604 *iter0* of Supplement C in the section  $N_2$ -variants. The row totals obtained after Step 1  
 605 in this special case do not have an interpretation in terms of effective or informative  
 606 numbers for sites.

607 The pre-processing method described so far can be obtained by the function *ipf2N2* in  
 608 the *douconca* package (ter Braak & van Rossum 2025) by setting the argument  
 609 *updateN2* to FALSE, while the method denoted *iter0* can be obtained by setting  
 610 *max\_iter* = 0.

611  
 612

613 Four variants of  $N_2$ -processing are:

- 614 1.  $N_2$  (in full:  $N_{2i}N_2$ ): default  $N_2$ -processing, *i.e.* IPF to species marginal that  
 615 represents informativeness ( $N_2(1 - N_2/N)$ ) and site marginal that represents effective  
 616 numbers of species, in the notation introduced in the following:  $N_{2i}N_{2i}$ ;
- 617 2.  $N_{2i}N_{2i}$ : IPF to marginals that represent informativeness of species and of sites;  
 618
- 619 3.  $N_2N_2$ : IPF to marginals that represent effective numbers ( $N_2$ ) of species and of sites;
- 620 4. *iter0*: species marginal scaled to informativeness of species without further scaling  
 621 of the site marginal

622

623 For presence–absence data, the  $N_2N_2$  variant yields no transformation because species  
 624 and site totals in a P/A table coincide with their  $N_2$ -values. The default variant  $N_2 \equiv$   
 625  $N_{2i}N_2$  similarly produces minimal adjustment when species totals are much smaller  
 626 than half the number of sites, in which case species informativeness is effectively  
 627 proportional to their occurrence frequencies.

628 There is no proof that IPF always returns a table with the desired marginals (Idel  
 629 2016).

## 630 7.2 IPF while updating $N_2$ marginals

631 The IPF procedure of the previous section leads to marginals that are proportional to  
632 the effective or informative numbers of the initial data table, but the precision of any  
633 CWM or SNC calculated during permutation is proportional to the effective numbers  
634 of the final data table, which differ from those of the initial table. For this reason, the  
635 effective and informative numbers are recalculated after each step.

636 The full procedure, with sequence as in the function *ipf2N2* in the *douconca* package,  
637 is:

638 Initialization:

639  $y_{ij} \leftarrow n_{ij}$

640 Set, for each row,  $r_i = N_{2i}$  or  $r_i = I_i$  and, for each column  $k_j = \tilde{N}_{2j}$  or  $k_j = \tilde{I}_j$  and for  
641 subsequent rescaling calculate  $c = k_+ / r_+$ .

642 Iterate:

643 Step 1. For each row  $i$  execute:  $y_{ij} \leftarrow y_{ij} (c r_i / y_{i+})$

644 Step 2: For each column update  $\tilde{N}_{2j}$  the effective number of the column of the current  
645 table  $\{y_{ij}\}$  and set  $k_j \leftarrow$  new  $\tilde{N}_{2j}$  or new  $\tilde{I}_j = \tilde{N}_{2j} (1 - \tilde{N}_{2j}/n)$ .

646 Step 3. For each column  $j$  execute:  $y_{ij} \leftarrow y_{ij} (k_j / y_{+j})$

647 Step 4: For each row update  $N_{2i}$  the effective number of the row of the current table  
648  $\{y_{ij}\}$  and set  $r_i \leftarrow$  new  $N_{2i}$  or new  $I_i = N_{2i} (1 - N_{2i}/m)$ .

649 Step 5: Calculate  $c \leftarrow k_+ / r_+$ .

650 Step 6: Stop if criterion is met (section 1.3), otherwise continue with Step 1.

651 Step 7: Goto step 1.

## 652 7.3 Stopping criterion

653 This section describes the detail of the stopping criterion used in Step 6 of the full IPF  
654 procedure of the previous section.

655 At step 6, each column has the intended sum (because of step 3) , but the row sum  
656 ( $y_{i+}$ ) may differ from the intended value ( $c r_i$ ). What counts in dc-CA are the ratios  $y_{i+}/$   
657  $y_{++}$ . The criterion is therefore based on the ratios of the row sums and the intended  
658 values. As full equality may be impossible to reach and the small row sums have little  
659 influence in dc-CA, the criterion is based on the range of the ratios  $\Delta_i = (c r_i + 1) /$   
660  $y_{i+} + 1$ ). The criterion value is the maximum of  $\max(\Delta_i)$  and  $1/\min(\Delta_i)$ . The iteration  
661 process stops when the criterion value does not improve since the previous iteration or  
662 when the maximum number of iterations has been reached, with a default value of  
663 10000.

## 664 7.4 R code of function IPF2N2 in douconca v1.2.5

665 The R code of function *ipf2N2* implementing the schematic code of sections 1-3. But  
666 first the R code for two helper functions, *fN2* for calculating the Hill number  $N_2$  and  
667 *fN2N\_N2* for calculating the effective number and informativeness (if  $N_{2N\_N2} =$

668 TRUE). The functions cater for the limiting cases. An example: When the  
 669 informativeness of a species is 0, while that species is the only one in a particular  
 670 site. In this particular case the  $N_2$  of the site is set to a small number, effectively 0. The  
 671 species and site are not removed as species and sites with very small marginal totals  
 672 do not influence dc-CA.

673

```
674 fn2 <- function(x ) {
675   sx<- sum(x)
676   sx <- if(is.na(sx) || sx < 1.0e-6) {
677     return(.Machine$double.eps)} else {sx}
678   x <- x / sx
679   return(1 / sum(x * x))
680 }
```

681

```
682 fn2N_N2 <- function(Y,
683   margin,
684   N2N_N2 = TRUE) {
685   N2 <- apply(X = Y, MARGIN = margin, FUN = fn2)
686   margin1 <- if (margin == 1) 2 else 1
687   if (N2N_N2) N2 <- N2 * pmax(1 - N2 / dim(Y)[margin1], .Machine$double.eps)
688 }
689 return(N2)
690 }
```

691

```
692 ipf2N2 <- function(Y,
693   max_iter = 10000,
694   updateN2 = TRUE,
695   N2N_N2_species = TRUE,
696   N2N_N2_sites = FALSE
697 ) {
698
699   Y <- as.matrix(Y)
700   rownames <- rownames(Y)
701   colnames <- colnames(Y)
702   R <- rowSums(Y)
703   K <- colSums(Y)
704
705   if (any(R == 0)) {
706     warning("Some sites do not have species.\n ",
707       paste(names(R)[which(R == 0)], collapse = " "), "\n")
708   }
709   if (any(K == 0)) {
710     warning("Some species are absent in every site.\n",
711       paste(names(K)[which(K == 0)], collapse = " "), "\n")
712   }
713   K[K < .Machine$double.eps] <- .Machine$double.eps
714   R[R < .Machine$double.eps] <- .Machine$double.eps
715   N2spp <- N2spp0 <- fn2N_N2(Y, 2, N2N_N2 = N2N_N2_species)
716   N2sites <- N2sites0 <- fn2N_N2(Y, 1, N2N_N2 = N2N_N2_sites)
717   N <- nrow(Y)
718   crit0 <- 1.0e10
719   crit1 <- crit0 - 1
720   iter <- 0
721   ratio <- 1
```

```

722 if(max_iter == 0){
723   Y <- Y %%% diag(N2spp / K)
724 } else {
725 while (crit1 < crit0 && iter < max_iter) {
726   iter <- iter + 1
727   crit0 <- crit1
728   Y0 <- Y
729   Y <- Y * ((N2sites / R) * ratio)
730   K <- colSums(Y)
731   if (updateN2) N2spp <- fN2N_N2(Y, 2, N2N_N2 = N2N_N2_species)
732   Y <- Y %%% diag(N2spp / K)
733   R <- rowSums(Y)
734   if (updateN2) N2sites <- fN2N_N2(Y, 1, N2N_N2 = N2N_N2_sites)
735   ratio <- sum(N2spp)/sum(N2sites)
736   mm <- range((ratio * N2sites+1)/(R+1))
737   #mm <- range(c(mm, (N2spp+1)/(K+1)))
738   crit1 <- max(c(1/mm[1], mm[2]))
739 }
740 Y <- Y0
741 if (updateN2) {
742   N2sites <- fN2N_N2(Y, 1, N2N_N2 = N2N_N2_sites)
743   N2spp <- fN2N_N2(Y, 2, N2N_N2 = N2N_N2_species)
744 }
745
746 if (iter == max_iter && max_iter > 0) {
747   warning(paste0("No convergence in ", max_iter, " iterations."), "\n")
748 }
749 if (any(Y<0)){
750   warning("some values in preprocessed Y negative")
751   Y[Y < 0] <- 0 # tiny non-negative values should not occur
752 }
753 }
754 R <- rowSums(Y)
755 attr(Y, which = "N2species_original") <- N2spp0
756 attr(Y, which = "N2sites_original") <- N2sites0
757 attr(Y, which = "N2species") <- N2spp
758 attr(Y, which = "N2sites") <- N2sites
759 attr(Y, which = "R/N2") <- (R / N2sites) / (sum(R) / sum(N2sites))
760 attr(Y, which = "iter") <- iter
761 attr(Y, which = "crit") <- crit0
762 colnames(Y) <- colnams
763 rownames(Y) <- rownams
764 if (fN2(R) / length(R) < 0.5) {
765   message("Warning: unbalanced site totals in return value:",
766           "N2 of row sums less than halve the number of rows.\n")
767 }
768 if (any(N2sites <= .Machine$double.eps)) {
769   warning("After processing: Some sites do not have species.\n ",
770           paste(rownames(Y)[which(N2sites <= .Machine$double.eps)],
771                 collapse = " "), "\n")
772 }
773 if (any(N2spp <= .Machine$double.eps)) {
774   warning("After processing: Some species are absent in every site.\n ",
775           paste(colnames(Y)[which(N2spp <= .Machine$double.eps)],
776                 collapse = " "), "\n")
777 }
778 return(Y)
779 }

```

780

## 781 7.5 Help text

ipf2N2 {douconca}

### 782 Iterative proportional fitting of an abundance table to Hill-N2 marginals

#### 783 Description

784 Function for pre-processing/transforming an abundance table by iterative proportional fitting,  
785 so that the transformed table has marginals proportional to  $N_2$  or  $N_2 (1-N_2/N)$  with  $N$  the  
786 number of elements in the margin. Hill-N2 is the effective number of species. It is of intrinsic  
787 interest in weighted averaging (CWM and SNC) as their variance is approximately inversely  
788 proportional to  $N_2$  (ter Braak 2019), and therefore of interest in [dc CA](#).

#### 789 Usage

```
790 ipf2N2(  
791   Y,  
792   max_iter = 10000,  
793   updateN2 = TRUE,  
794   N2N_N2_species = TRUE,  
795   N2N_N2_sites = FALSE  
796 )
```

### 797 Arguments

|                             |  |
|-----------------------------|--|
| <code>Y</code>              | abundance table (matrix or dataframe-like), ideally, with names for rows and columns.  |
| <code>max_iter</code>       | maximum number of iterative proportional fitting (ipf) iterations. If <code>max_iter = 0</code> , the columns are divided by their effective number or informativeness ( $N_2$ or $N_2 (1-N_2/N)$ ), depending on the setting of <code>N2N_N2_species</code> ) without further pre-processing and the row sums are then, with <code>N2N_N2_species = TRUE</code> , sums of informativeness instead of effective number of informative species.   |
| <code>updateN2</code>       | logical, default <code>TRUE</code> . If <code>FALSE</code> the marginal sums are proportional to the $N_2$ -marginals of the initial table, but the $N_2$ -marginals of the returned matrix may not be equal to their marginal sum. If <code>updateN2 = TRUE</code> and <code>N2N_N2_species = TRUE</code> (the default), the column marginals are $N_2 (N-N_2) / N$ with $N$ the number of sites. The row sums are then proportional to, what we term, the effective number of informative species. If <code>N2N_N2_species = FALSE</code> , the returned transformed table has $N_2$ columns marginals, <i>i.e.</i> <code>colSums(Y2) = N2species(Y2)</code> with <code>Y2</code> the return value of <code>ipf2N2</code> . If converged, $N_2$ row marginals are equal to the row sums, <i>i.e.</i> <code>rowSums(Y2) = approx. const*N2sites(Y2)</code> and <code>const</code> a constant. |
| <code>N2N_N2_species</code> | Set species marginal to the value of $N_2 (1-N_2/N)$ for each species. Default <code>TRUE</code> . If <code>FALSE</code> , the marginal is set to the $N_2$ value of each species.   |
| <code>N2N_N2_sites</code>   | Default <code>FALSE</code> sets the marginal proportional to the $N_2$ value of each site. If <code>TRUE</code> , the marginal is set to $N_2 (1-N_2/m)$ , with $m$ the number of species.   |

798

## 799 Details

800 Applying `ipf2N2` with `N2N_N2_species=FALSE` to a presence-absence data table returns  
801 the same table. However, a species that occurs everywhere (or in most of the sites) is not  
802 very informative. This is acknowledged with the default option `N2N_N2_species=TRUE`.  
803 Then, with `N2N_N2_species=TRUE`, species that occur in more than half the number of  
804 sites are down-weighted, so that the row sum is no longer equal to the richness of the site  
805 (the number of species), but proportional to the number of informative species. The returned  
806 matrix has the intended species marginal (column sums), by construction of the algorithm,  
807 even without convergence. On convergence, it has the intended site marginal (row sums).

## 808 Value

809 a matrix of the same order as the input  $Y$ , obtained after ipf to N2-marginals.

## 810 References

811 ter Braak, C.J.F. (2019). New robust weighted averaging- and model-based methods for  
812 assessing trait-environment relationships. *Methods in Ecology and Evolution*, 10 (11), 1962-  
813 1971. [doi:10.1111/2041-210X.13278](https://doi.org/10.1111/2041-210X.13278)

814 ter Braak, C.J.F. (2026). Fourth-corner latent variable models overstate confidence in trait-  
815 environment relationships and what to use instead *Environmental and Ecological*  
816 *Statistics*. [doi:10.1007/s10651-025-00696-0](https://doi.org/10.1007/s10651-025-00696-0)

817

818 The arguments `N2N_N2_species` and `N2N_N2_sites` determine the four variants of  
819  $N_2$ -processing are follows:

- 820 1.  $N_2$  (in full:  $N_{2i}N_2$ ): default  $N_2$ -processing, *i.e.* IPF to species marginal that  
821 represents informativeness ( $N_2(1 - N_2/N)$ ) and site marginal that represents effective  
822 numbers of species, in the notation introduced in the following:  $N_{2i}N_{2i}$ ; *i.e.*  
823 `N2N_N2_species = TRUE and N2N_N2_sites = FALSE`.
- 824 2.  $N_{2i}N_{2i}$ : IPF to marginals that represent informativeness of species and of sites,  
825 *i.e.* `N2N_N2_species = TRUE and N2N_N2_sites = TRUE`.
- 826 3.  $N_2N_2$ : IPF to marginals that represent effective numbers ( $N_2$ ) of species and of sites,  
827 *i.e.* `N2N_N2_species = FALSE and N2N_N2_sites = FALSE`.
- 828 4. `iter0`: species marginal scaled to informativeness of species without further scaling  
829 of the site marginal, *i.e.* `max_iter = 0`.

830

831 In the section on  $N_2$  – variants in Supplement C, `updateN2 = TRUE`.

## 832 7.6 References

833 Fienberg, S.E. (1970) An Iterative Procedure for Estimation in Contingency Tables.  
834 *The Annals of Mathematical Statistics*, **41.3**, 907–917.

835 <https://doi.org/10.1214/aoms/1177696968>

836 Idel, M. (2016) A review of matrix scaling and Sinkhorn’s normal form for matrices  
837 and positive maps. <https://arxiv.org/pdf/1609.06349.pdf>

838 ter Braak, C.J.F. & van Rossum, B.-J. (2025) R package douconca: Double  
839 constrained correspondence analysis for multi-trait multi-environment analysis  
840 v1.2.5. <https://doi.org/10.32614/CRAN.package.douconca>

## 841 **8. Supplement C: Real and simulated data sets with additional** 842 **results**

843

844 A key limitation of drawing general conclusions from analyses of simulated and real  
845 datasets, as compared with theory, is that it is virtually impossible to encompass all  
846 relevant scenarios. This is reflected below in Figs. S1–S4, where the advantage of  
847 updating  $N_2$  during IPF (Appendix B) is not evident in Fig. S1 and only debatable in  
848 Fig. S4.

### 849 **8.1 Real and simulated data sets**

850 Four sets of simulated datasets were used. The first three—Aravo, Revisit, and  
851 Beetle—were generated from generalized linear mixed models (GLMMs) fitted to  
852 their respective real datasets. The fourth set reproduces the simulations in Fig. 3 of  
853 Niku et al. (2021), using the same model specification as in Fig. 3 of ter Braak (2026).  
854 Among the four sets, only the Beetle-simulations include multiple traits and multiple  
855 environmental variables.

856 In all sets, the underlying model is a GLM with fixed main effects of the trait, the  
857 environmental variables, and their interaction. Each set also includes latent variables  
858 associated with the trait or the environmental variables, as in GLMM3 models  
859 (ter Braak 2026). In the fourth set, on a GLLVM2(2), the fixed effects are  
860 supplemented by an environment-related latent variable and two additional unrelated  
861 latent variables to capture residual variation beyond that implied by the response  
862 distribution. The response distributions, in order, are Poisson, beta-binomial, and  
863 negative binomial (for the third and fourth sets). The link function is log, except in the  
864 second set, which uses a logit link.

865

#### 866 **Aravo-simulations**

867

868 The Aravo-simulated data sets were taken from Appendix A7 in the supporting  
869 information of ter Braak (ter Braak 2019). In that Appendix, a general linear mixed  
870 model with random species-dependent environmental effect and site-dependent trait  
871 effects (*i.e.* a GLMM3 model, which is termed MLM3 in ter Braak 2019) was fitted  
872 using the *aravo* data set in the R package *ade4* (Dray & Dufour 2007) with the single  
873 trait SLA (Specific Leaf Area) and the single environmental variable Snow (mean  
874 snowmelt date in Julian day averaged over 1997-1999) and with Poissonian response  
875 distribution.. It has 82 plant species and 75 sites. The fitted model was modified to  
876 generate data with varying strength of the trait-environment relationship (250 data sets  
877 for each scenario). The model fitted to these data was:

```
878 library(lme4)
```

```
879 fpol2 <- y ~ poly(trait,2) + poly(env,2) + trait:env +  
880 (1+trait|site) + (1+env|species))
```

```
881 MLM3 <- glmer(fpol2, data=dat, family= poisson,  
882 nAGQ=0, control = glmerControl(calc.derivs=F))
```

883 The simulated data are generated using `simulate(MLM3)` after modification of the  
884 regression coefficient of the trait-environment interaction. The data are identical to  
885 those used in ter Braak (2019).

886  
887

### 888 Revisit-simulations

889

890 The Revisit-simulated data sets were taken from ter Braak (ter Braak 2019), who fitted  
891 a GLMM3 model to a subset of the Whittaker Siskiyou Mountains Revisit  
892 data(Damschen, Harrison & Grace 2010) with the single trait leaf carbon-to-nitrogen  
893 ratio (C:N) and the environmental variable Topographic Moisture Gradient (TMG).  
894 Briefly, the data consists of the cover abundance of 75 species in 52 sites, calculated  
895 from the number of 100 quadrat corners per site that each species intersected. The  
896 model was fitted with a beta-binomial response distribution, which was then modified  
897 to generate data with varying strength of the trait-environment relationship (250 data  
898 sets for each scenario). The model fitted to these data was:

899

```
900 library(glmTMB)
```

```
901 formula.MLM3 <- y ~ trait*env +(1+trait|site)+(1+env|species)
```

```
902 MLM3 <- glmTMB(formula.MLM3, family = betabinomial, data= Revisit)
```

903 The data are generated using `simulate(MLM3)` after modification of the regression  
904 coefficient of the trait-environment interaction. The data are identical to those used in  
905 ter Braak (2019).

906

907

### 908 Beetle-simulations

909

910 The Beetle-simulations are obtained by fitting a GLMM3 model to ground beetle data  
911 from Ribera et al. (2001) which contains the numbers of individuals of 68 ground beetle  
912 species in 87 sites. The traits and environmental variables are those used by Niku et al.  
913 (2021). The four traits are total length (LTL), maximum vaulting of the  
914 pronotum(LPH), overwintering (OVE) with two classes adults only, and adults and/or  
915 larvae and breeding season (BRE) with three classes (spring, summer and autumn or  
916 winter), while the four environmental variables are the management intensity index  
917 (Management), percentage moisture content (Moist), elevation, and pH, with full  
918 definitions in the original paper (Ribera et al. 2001). The GLMM3 model was fitted to  
919 the data with negative binomial response distribution, “NB2” in `glmTMB` v2.0.2  
920 (Brooks *et al.* 2017), and model formula:

```
921 fGLMM3.H1 <- y ~ Management+ Moist + Elevation+pH +  
922 LTL + LPH +OVE + BRE +  
923 (Management+ Moist + Elevation+pH):(LTL + LPH +OVE + BRE)+  
924 (1 + Management+ Moist + Elevation+pH |species) +  
925 (1 + LTL + LPH +OVE + BRE |site)
```

926 and call:

```
927 GLMM3.fit <- glmTMB(fGLMM3.H1, data=datalong, family= nbinom2(link =  
928 "log"), se = TRUE)
```

929

930 Data sets with varying strength of the trait-environment relationship (250 data sets for  
931 each scenario) were generated using `simulate(model)` after setting the 16 regression  
932 coefficients of the trait-environment interaction to the product of the fitted coefficients  
933 and effect size (effect size = 0, 0.3, 0.45, 0.6, 0.75, 0.9).  
934

#### 935 GLLVM2(2) simulations

936  
937 The set-up the simulations of the GLLVM2(2) model was proposed by Niku et al.  
938 (2021) and the actual set used here is the same as that in ter Braak (2026). It has negative  
939 binomially distributed abundance of 40 species in 70 sites. The GLLVM2(2) model is  
940 a generalized linear mixed model with random species-dependent environmental  
941 effects and two latent variables to mimic structured variation beyond the unstructured  
942 negative binomial variation. The model differs from the previous GLMM3 models by  
943 their latent variables and by *not* including random site-dependent trait effects. The  
944 effect size of the trait-environment interactions and the variance of the random species-  
945 dependent environmental effects is varied in this set.

946 The data are generated using `simulate(model)` after modification of the regression  
947 coefficient of the trait-environment interaction. The data are identical to those used in  
948 ter Braak (2026).

949

950

#### 951 CESTES data set

952

953

954 The CESTES data base is a set of 80 curated datasets from published trait-studies  
955 (Jeliazkov *et al.* 2020; Jeliazkov & Chase 2024). CESTES stands for “metaCommunity  
956 Ecology: Species, Traits, Environment and Space”. Each dataset includes four  
957 matrices: species community abundances or presences/absences across multiple sites,  
958 species trait information, environmental variables and spatial coordinates of the  
959 sampling sites. Here, the information about space is not used. With number between  
960 parentheses, the data sets are from terrestrial (55), fresh water (14) and marine  
961 ecosystems (11) all over the globe with many taxonomic group represented and sites  
962 varying in extent from small to large. Measures of abundance can be classified as  
963 number-of-individuals (45), presence-absence (11), abundance-index (9) and other  
964 (15). Table 1 provides a summary of numbers of units and variables in the CESTES  
965 data base.  
966

967 Table 1. Summary of numbers of sites, species and variables in the CESTES data base  
 968 (nTraitsI and nEnvI are the number of traits and environmental variables, respectively,  
 969 after transforming categorical variables to indicator variables, as required in dc-CA).

| 970 |          | Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|-----|----------|------|---------|--------|------|---------|------|
| 971 | nSpecies | 6    | 32      | 53     | 71.6 | 107     | 208  |
| 972 | nTraits  | 1    | 5       | 8      | 13.8 | 16      | 89   |
| 973 | nTraitsI | 1    | 8       | 13     | 17.4 | 20      | 89   |
| 974 | nSites   | 7    | 27      | 40     | 70.9 | 88      | 364  |
| 975 | nEnv     | 1    | 6       | 9      | 11.9 | 13      | 47   |
| 976 | nEnvI    | 1    | 7       | 13     | 15.1 | 20      | 63   |

977  
 978

## 979 References

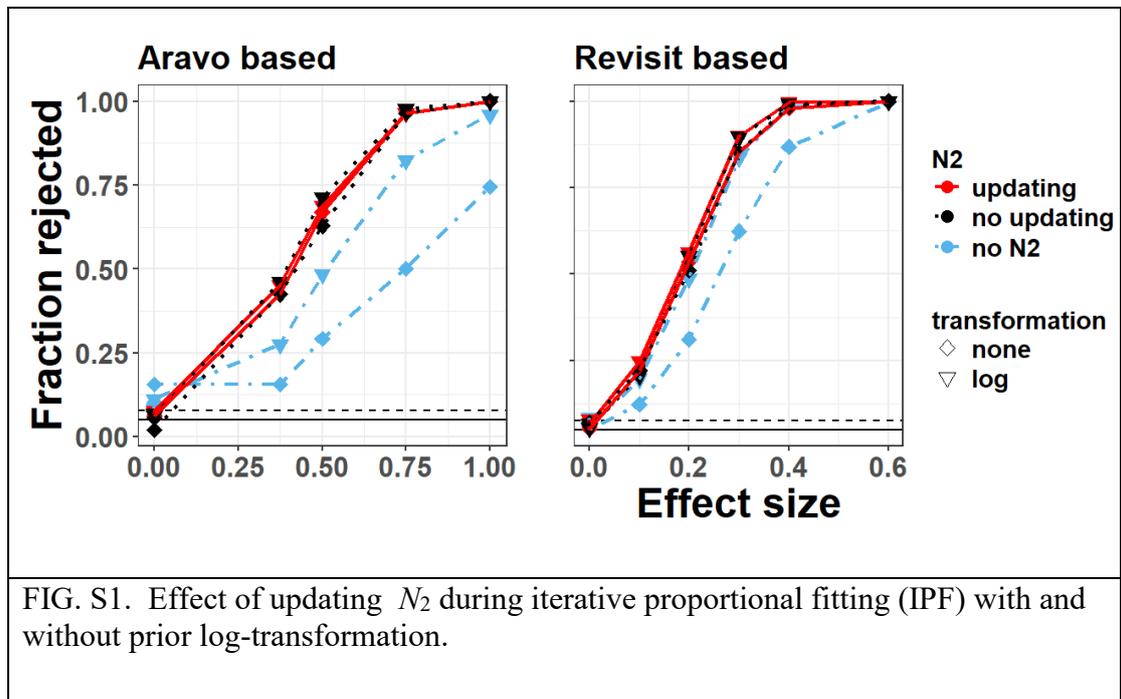
- 980 Brooks, M.E., Kristensen, K., van Benthem, K.J., Magnusson, A., Berg, C.W.,  
 981 Nielsen, A., Skaug, H.J., Maechler, M. & Bolker, B.M. (2017) glmmTMB  
 982 balances speed and flexibility among packages for zero-inflated generalized  
 983 linear mixed modeling. *The R Journal*, **9**, 378-400.  
 984 <https://doi.org/10.32614/RJ-2017-066>
- 985 Damschen, E.I., Harrison, S. & Grace, J.B. (2010) Climate change effects on an  
 986 endemic-rich edaphic flora: resurveying Robert H. Whittaker's Siskiyou sites  
 987 (Oregon, USA). *Ecology*, **91**, 3609-3619. <https://doi.org/10.1890/09-1057.1>
- 988 Dray, S. & Dufour, A.B. (2007) The ade4 package: implementing the duality diagram  
 989 for ecologists. *Journal of Statistical Software*, **22**, 1-20.  
 990 <https://doi.org/10.18637/jss.v022.i04>
- 991 Jeliaskov, A. & Chase, J.M. (2024) When Do Traits Tell More Than Species about a  
 992 Metacommunity? A Synthesis across Ecosystems and Scales. *The American*  
 993 *Naturalist*, **203**, E1-E18. <https://doi.org/10.1086/727471>
- 994 Jeliaskov, A., Mijatovic, D., Chantepie, S., Andrew, N., Arlettaz, R., Barbaro, L.,  
 995 Barsoum, N., Bartonova, A., Belskaya, E., Bonada, N., Brind'Amour, A.,  
 996 Carvalho, R., Castro, H., Chmura, D., Choler, P., Chong-Seng, K., Cleary, D.,  
 997 Cormont, A., Cornwell, W., de Campos, R., de Voogd, N., Doledéc, S., Drew, J.,  
 998 Dziock, F., Eallonardo, A., Edgar, M.J., Farneda, F., Hernandez, D.F.,  
 999 Frenette-Dussault, C., Fried, G., Gallardo, B., Gibb, H., Gonçalves-Souza, T.,  
 1000 Higtuti, J., Humbert, J.-Y., Krasnov, B.R., Saux, E.L., Lindo, Z., Lopez-  
 1001 Baucells, A., Lowe, E., Marteinsdottir, B., Martens, K., Meffert, P., Mellado-  
 1002 Díaz, A., Menz, M.H.M., Meyer, C.F.J., Miranda, J.R., Mouillot, D., Ossola,  
 1003 A., Pakeman, R., Pavoine, S., Pekin, B., Pino, J., Pocheville, A., Pomati, F.,  
 1004 Poschlod, P., Prentice, H.C., Purschke, O., Ravel, V., Reitalu, T., Renema,  
 1005 W., Ribera, I., Robinson, N., Robroek, B., Rocha, R., Shieh, S.-H., Spake, R.,  
 1006 Staniaszek-Kik, M., Stanko, M., Tejerina-Garro, F.L., Braak, C.t., Urban,  
 1007 M.C., Klink, R.v., Villéger, S., Wegman, R., Westgate, M.J., Wolff, J.,  
 1008 Żarnowiec, J., Zolotarev, M. & Chase, J.M. (2020) A global database for  
 1009 metacommunity ecology, integrating species, traits, environment and space.  
 1010 *Scientific Data*, **7**, 6. <https://doi.org/10.1038/s41597-019-0344-7>
- 1011 Niku, J., Hui, F.K.C., Taskinen, S. & Warton, D.I. (2021) Analyzing environmental-  
 1012 trait interactions in ecological communities with fourth-corner latent variable  
 1013 models. *Environmetrics*, **n/a**, e2683. <https://doi.org/10.1002/env.2683>
- 1014 Ribera, I., Dolédec, S., Downie, I.S. & Foster, G.N. (2001) Effect of Land  
 1015 Disturbance and Stress on Species Traits of Ground Beetle Assemblages.

1016 *Ecology*, **82**, 1112-1129. <https://doi.org/10.1890/0012->  
 1017 9658(2001)082[1112:EOLDAS]2.0.CO;2  
 1018 ter Braak, C.J.F. (2019) New robust weighted averaging- and model-methods for  
 1019 assessing trait–environment relationships. *Methods in Ecology and Evolution*,  
 1020 **10**, 1962-1971. <https://doi.org/10.1111/2041-210X.13278>  
 1021 ter Braak, C.J.F. (2026) Fourth-corner latent variable models overstate confidence in  
 1022 trait–environment relationships and what to use instead. *Environmental and*  
 1023 *Ecological Statistics*. <https://doi.org/10.1007/s10651-025-00696-0>;  
 1024 <https://edepot.wur.nl/707477>; <https://rdcu.be/eXukf>.  
 1025

## 1026 8.2 N2 updating

1027

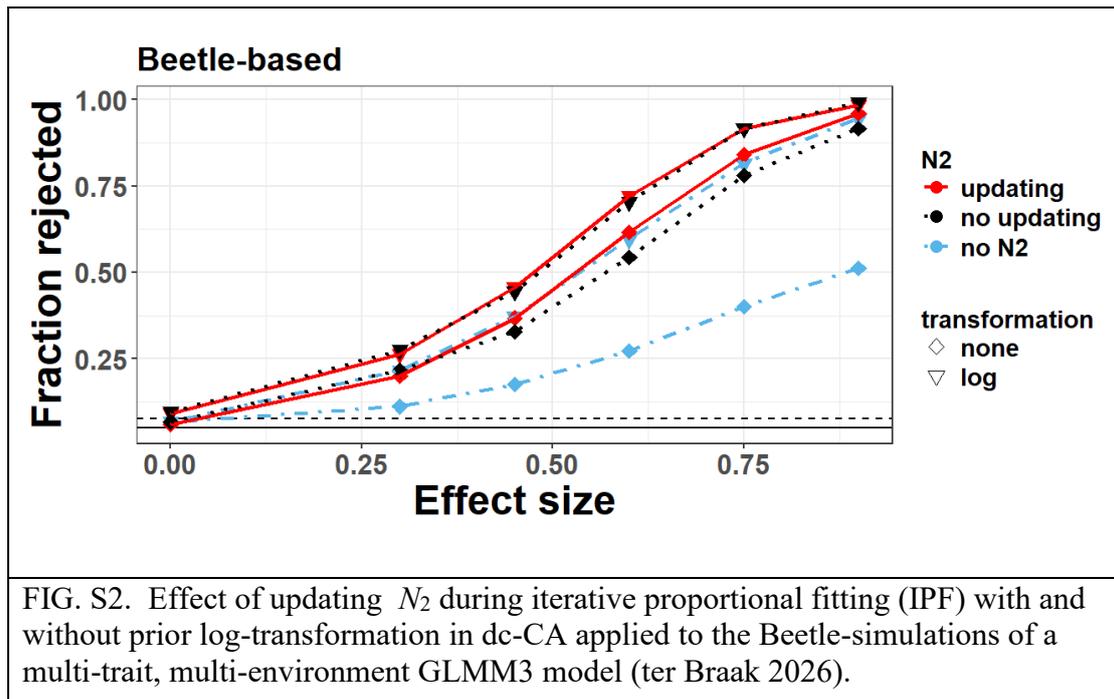
1028 Updating  $N_2$  (species informativeness and site effective number of informative  
 1029 species) during IPF (Supplement B) has minor effect on the power in the Aravo- and  
 1030 Revisit-simulations, surely when compared with no  $N_2$ -processing (Figure S1).



1031

1032 In the Beetle-simulations (Fig. S2), updating  $N_2$  led to higher power without prior  
 1033 transformation and, as in Fig. S1, to similar power after log-transformation.

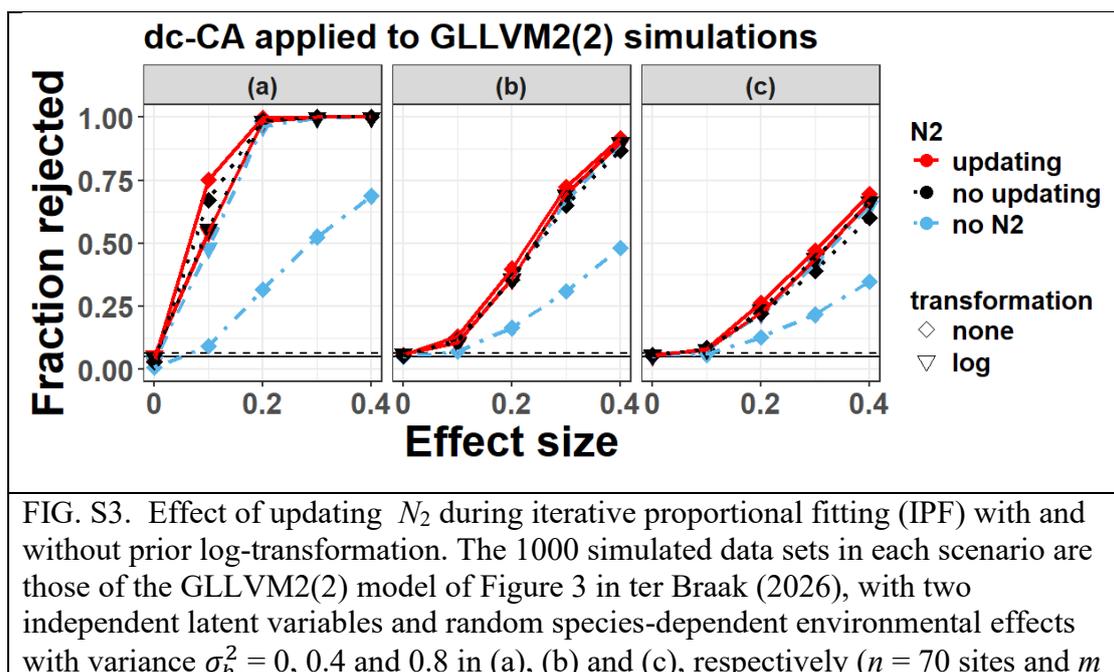
1034



1035

1036 The same pattern is visible in Figure S3 in the simulations of a GLLVM2(2) model in  
 1037 Figure 3 of both Niku et al. (2017) and ter Braak (2026): improvement without prior  
 1038 transformation but no effect of updating  $N_2$  after log-transformation. Fig. S3a show a  
 1039 relatively rare case in which  $N_2$ -processing after a log-transformation led to lower  
 1040 power than  $N_2$ -processing alone.

1041 The empirical power in Fig. S3 without pre-processing is much lower than in Fig. 3 of  
 1042 ter Braak (2026) because the latter used abundance data standardized by site totals,  
 1043 whereas Fig. S3 did not. Dividing by site totals is the default in *douconca* when  $N_2$ -  
 1044 processing is omitted, making dc-CA align more closely with CWM regressions,  
 1045 which generally avoid differential site weighting.  $N_2$ -processing replaces such ad-hoc  
 1046 choices with a principled one.



= 40 species). The scenarios were proposed in the Figures 1 and 3 of Niku et al. (2021).

1047

1048 Updating  $N_2$  during IPF yielded generally very similar  $P$ -values in dc-CA analysis of  
1049 the CESTES data sets (Figure S3).

1050

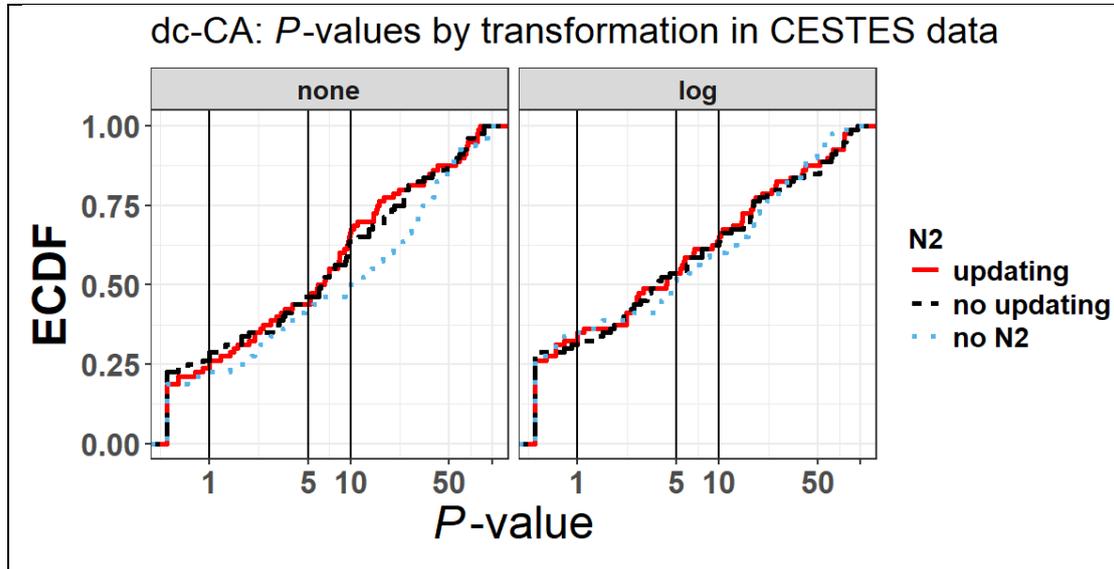


FIG. S3. Effect of  $N_2$  and updating  $N_2$  during iterative proportional fitting (IPF) with and without prior log-transformation: empirical cumulative density function (ECDF) of the significance level ( $P$ -value) using dc-CA for trait-environment association in the 80 data sets of the CESTES database in relation to updating  $N_2$  and prior transformation (none and log). The abscissa starts at  $P = 0.005$  (0.5%). The vertical lines are at the 1, 5 and 10% significance levels.

1051

### 1052 8.3 $N_2$ variants

1053

1054 The  $N_2$ -variant of the main text ( $N_2$  in Figure S4), which is the default in the *ipf2N2*  
1055 function in *douconca*, showed equal or more power in the Aravo- and Revisit-  
1056 simulations than the other three variants (Fig. S4, Supplement B). With and without  
1057 log-transformation, the  $N_2$ -variant in which both marginals are effective numbers  
1058 ( $N_2N_2$  in Figure S4) performs worse than the main text variant in the Aravo-  
1059 simulations, whereas the simplest variant which is without iteration (*iter0* in Figure  
1060 S1) and divides each species' abundance by the informativeness of the species only,  
1061 performed worse in the Revisit-simulations. The variant in which both marginals  
1062 reflect informativeness performed similar to the  $N_2$ -variant of the main text.

1063

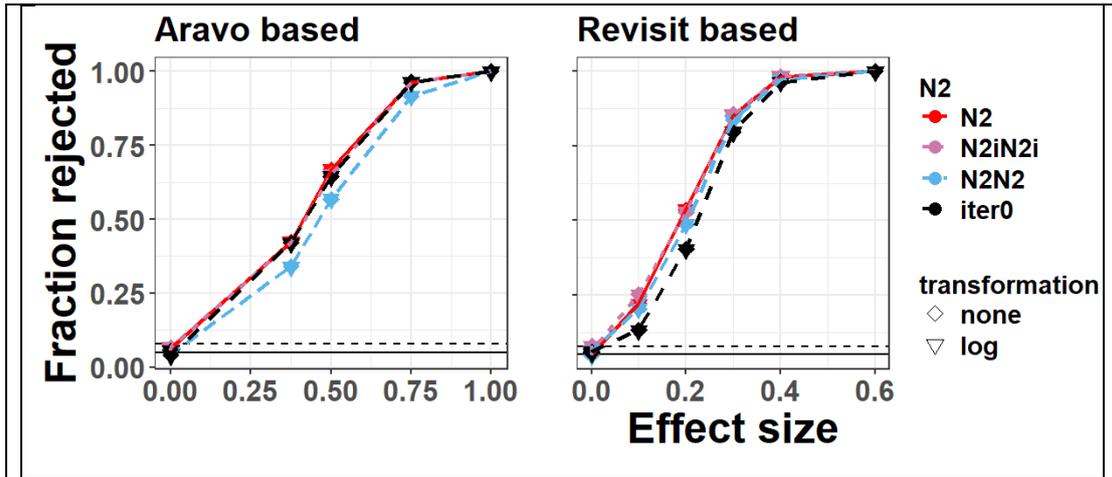


FIG. S4. Variants of  $N_2$ -processing with and without prior log-transformation:  $N_2$ : default  $N_2$ -processing, *i.e.* iterative proportional fitting (ipf) to species marginal that represents informativeness ( $N_2(1 - N_2/N)$ ) and site marginal that represents effective numbers of species, in the following notation  $N_2iN_2i$ ;  $N_2iN_2i$ : ipf to marginals that represent informativeness of species and of sites;  $N_2N_2$ , ipf to marginals that represent effective numbers ( $N_2$ ) of species and of sites;  $iter0$ : species marginal scaled to informativeness of species without further scaling of the site marginal.

1064

1065 The four variants of  $N_2$ -processing yielded similar  $P$ -value distributions in dc-CA  
 1066 analysis of the CESTES data sets with and without prior log-transformation (Figure  
 1067 S5). Without prior transformation, version 'iter0' did well below  $P = 0.07$ .

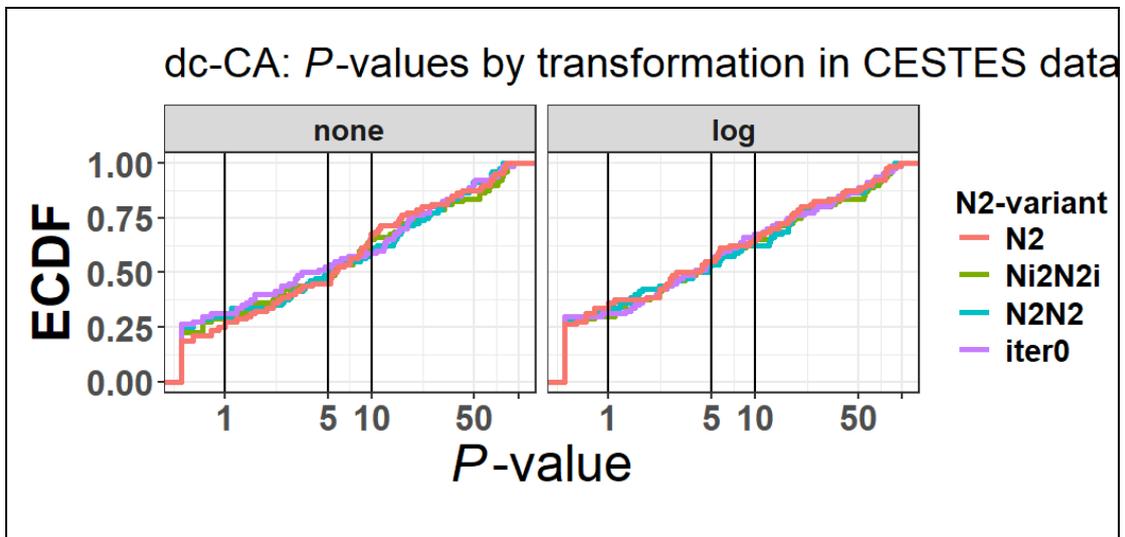


FIG. S5. Variants of  $N_2$ -processing with and without prior log-transformation: empirical cumulative density function (ECDF) of the significance level ( $P$ -value) using dc-CA for trait-environment association in the 80 data sets of the CESTES database in relation to method of  $N_2$ -variants (see legend Fig. S1) and prior transformation (none and log). The abscissa starts at  $P = 0.005$  (0.5%). The vertical lines are at the 1, 5 and 10% significance levels.

1068 **8.4 Prior power transformations**

1069 Figure S6 complements Figure 1 in the main text and shows the square-root  
 1070 transformation with other power transformations with and without subsequent  $N_2$ -  
 1071 processing. With  $N_2$ -processing, all power transformations gave similar the rejection  
 1072 rates in the Aravo- and Revisit-simulations.

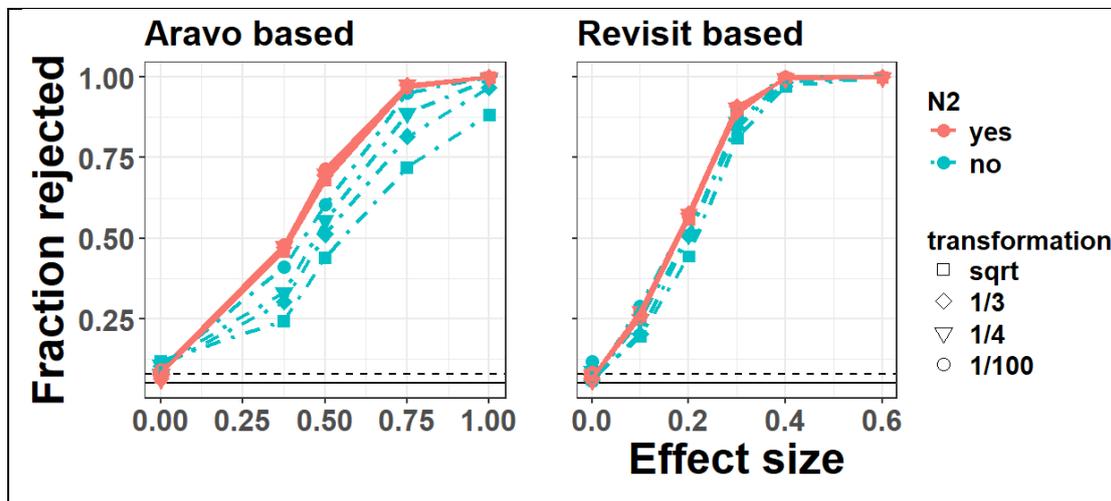


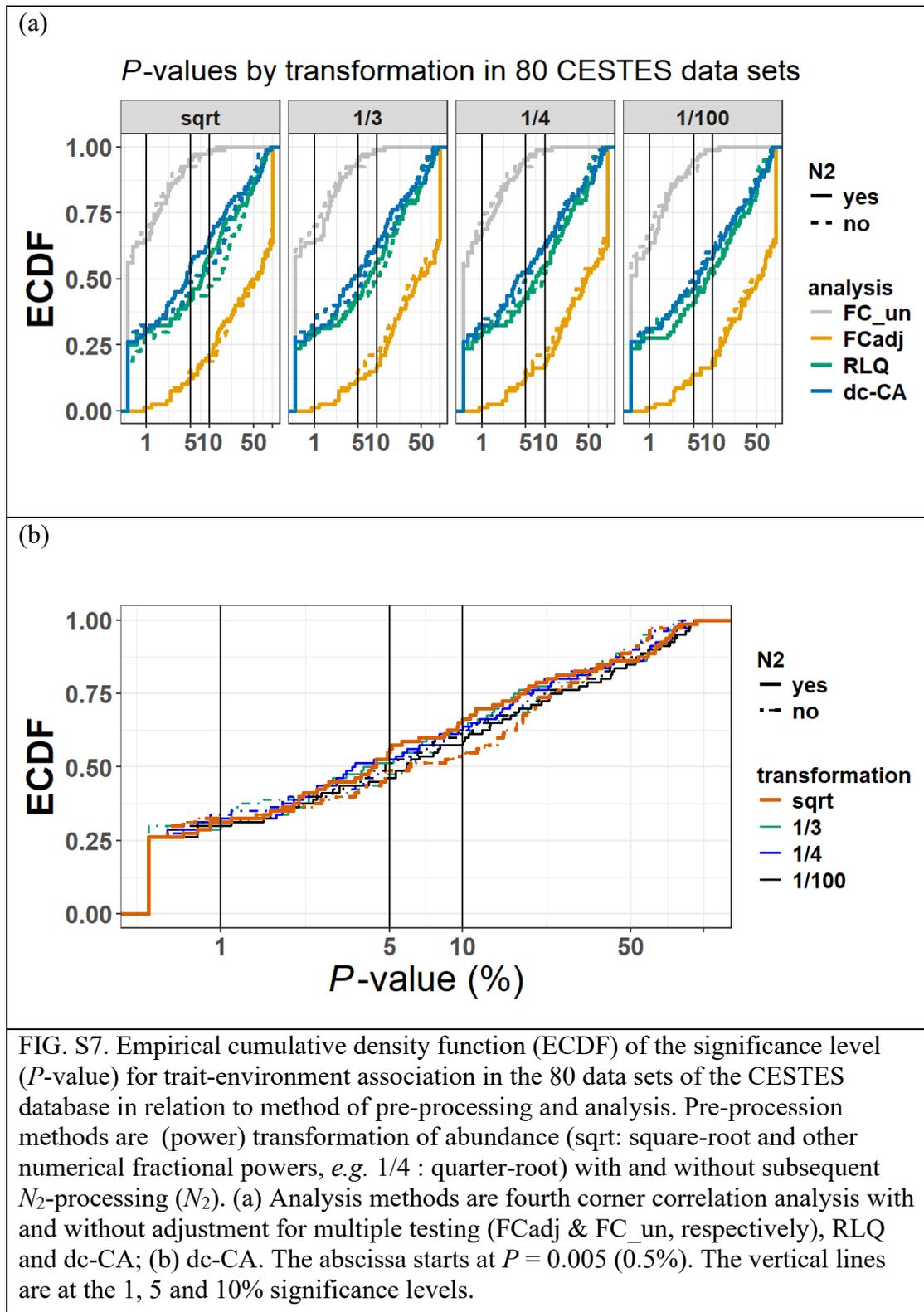
FIG. S6. Rejection rates of preprocessing methods against the size of the trait-environment association (Effect size), on the simulations of GLMM3 models for the Aravo and Revisit data in ter Braak (2019). Pre-processing methods are transformation of abundance (sqrt: square-root, and other fractional powers with and without subsequent  $N_2$ -processing ( $N_2$ ). Additional simulations with increased trait variance were carried out, as in ter Braak (2019), for effect size zero, and the lines start at this scenario so as to emphasize potential Type I error rate inflation. The horizontal solid line is at the nominal significance threshold; rates above the dashed line (at 0.078) are significantly greater than 0.05.

1073

1074 Figure S7 complements Figure 2 in the main text and shows the empirical  
 1075 distributions of the  $P$ -value obtained with a number of prior power transformations in  
 1076 the CESTES data base.  $N_2$ -processing generally increased the evidence but its effect  
 1077 diminished with the smaller powers. With cube-root there are slightly more data sets  
 1078 that give low  $P$ -values (below 0.025) without subsequent  $N_2$ -processing than with  
 1079 such processing (Fig. S7). The lines for power 1/100 are very similar to those for  
 1080 transformation to presence-absence data. The square-root transformation followed by  
 1081  $N_2$ -processing performs particularly well at the significance levels that matter most  
 1082 (5% and 10%) and shows no meaningful loss of power at the 1% level.

1083 .

1084



1085

1086 **8.5 Statistical analysis of the  $N_2$ -effect in the CESTES data base**

1087 This section focusses on the effect of prior-abundance transformation with or without  
 1088  $N_2$ -processing on *P*-values issued by FC, RLQ and dc-CA applied to the 80 CESTES

1089 data sets. The  $P$ -values are transformed to Evidence with Evidence =  $-\log(P)$  to give  
 1090 more emphasis on the small  $P$ -values. The analysis is by linear regression and Anova  
 1091 of the Evidence on to the factors *method of analysis*, *abundance transformation*, and  
 1092  $N_2$ -processing, with 3, 4 and 2 levels, respectively. The four levels of *method of*  
 1093 *analysis (mth)* are FC adjusted for multiple testing (FC-adj), RLQ and dc-CA. The  
 1094 three levels of *abundance transformation (pwr)* are none, square-root, log and  
 1095 transformation to presence absence. The two levels of  $N_2$ -processing ( $N_2$ ) are yes  
 1096 and no. As the Evidence varies with data set and is computed for each set, a paired  
 1097 analysis is asked for. This is achieved by adding the factor the *dataset* with 80 levels  
 1098 to each regression and Anova analysis. The Anova analyses use function LSD.test  
 1099 from the library agricolae (de Mendiburu 2023).

1100 There was no statistical evidence for interaction between *method of analysis* and  
 1101 *method of pre-processing*, here  $N_2 * pwr$ , ( $P = 0.68$ ) as judged by the comparison of  
 1102 the models

1103  $Evidence \sim dataset + N_2 * pwr * mth$

1104  $Evidence \sim dataset + N_2 * pwr + mth$ .

1105 The mean Evidence per *method of analysis* with grouping on least-significant  
 1106 difference (LSD) without adjustment at  $P = 0.05$  is

| ##   |          | Evidence | groups |
|------|----------|----------|--------|
| 1108 | ## dc-CA | 3.305665 | a      |
| 1109 | ## RLQ   | 3.025167 | b      |
| 1110 | ## FCadj | 1.149048 | c      |

1111 The LSD = 0.131. Note that an Evidence of 3 corresponds to an  $P$  value issued by the  
 1112 method of analysis of 0.05, indicating that the mean Evidence of RLQ and dc-CA is  
 1113 just below  $P = 0.05$ , the mean Evidence of FCadj of 1.149 corresponds to  $P = 0.32$ .

1114 The mean Evidence per *method of pre-processing* with grouping on least-significant  
 1115 difference (LSD) without adjustment at  $P = 0.05$ , with LSD = 0.213, is

| ##   |             | Evidence | groups |
|------|-------------|----------|--------|
| 1117 | ## yes.log  | 2.619430 | a      |
| 1118 | ## yes.sqrt | 2.616760 | a      |
| 1119 | ## no.P/A   | 2.578651 | a      |
| 1120 | ## no.log   | 2.574581 | a      |
| 1121 | ## yes.none | 2.535546 | a      |
| 1122 | ## yes.P/A  | 2.523273 | a      |
| 1123 | ## no.sqrt  | 2.464680 | a      |
| 1124 | ## no.none  | 2.033426 | b      |

1125

1126 The low values are caused by the low mean Evidence in the FCadj analysis. In the  
 1127 following, the  $P$ - values are restricted to those of dc-CA. The grouping for dc-CA  
 1128 only:

| ##   |             | Evidence | groups |
|------|-------------|----------|--------|
| 1129 | ## yes.log  | 3.505116 | a      |
| 1130 | ## yes.sqrt | 3.462850 | ab     |
| 1131 | ## no.log   | 3.408959 | ab     |
| 1132 | ## no.P/A   | 3.357034 | ab     |
| 1133 | ## yes.P/A  | 3.329233 | ab     |
| 1134 | ## no.sqrt  | 3.297212 | ab     |

```

1135 ## yes.none 3.225800      b
1136 ## no.none  2.859117      c

1137 with LSD= 0.270. No pre-processing at all is clearly worst. While N2-processing gives
1138 higher mean Evidence than no preprocessing (yes.none compared to no.none),
1139 additional pre-processing gives further improvement, with log-transformation
1140 performing best and the square-root as runner-up.

1141 Quantifying the abundance transformations (none=0, sqrt = 1, log = 2, P/A=4), the
1142 model with factors, Evidence ~ dataset + N2 * pwr is not statistically better
1143 fitting than the quantified model Evidence ~ dataset+ N2q+pwrq+ I(pwrq^2) +
1144 N2qpwrq (P = 0.97). The coefficient and R-squared- and F-statistics of the latter are:

1145 ## Coefficients:
1146 ##              Estimate Std. Error t value Pr(>|t|)
1147 ## (Intercept)      1.96258    0.31738   6.184 1.21e-09 ***
1148 ## datasetBarbaro2009a  0.96264    0.43391   2.219 0.026922 *
1149 ...
1150 ## datasetYates2014    -0.62001    0.43391  -1.429 0.153595
1151 ## N2q                 0.33811    0.11480   2.945 0.003363 **
1152 ## pwrq                0.49916    0.11168   4.469 9.51e-06 ***
1153 ## I(pwrq^2)          -0.11287    0.03430  -3.290 0.001064 **
1154 ## N2qpwrq            -0.12529    0.06136  -2.042 0.041643 *
1155 ## ---
1156 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
1157 ##
1158 ## Residual standard error: 0.8678 on 556 degrees of freedom
1159 ## Multiple R-squared:  0.8644, Adjusted R-squared:  0.8441
1160 ## F-statistic: 42.68 on 83 and 556 DF,  p-value: < 2.2e-16

1161 showing again that N2 has a significant effect (P < 0.01) on the Evidence without
1162 other preprocessing, but that this effect is diminishing with the stronger prior
1163 abundance transformation, as shown by the negative coefficient of the interaction
1164 N2qpwrq.

1165 With all power transformations tested the group means are (LSD = 0.234):

1166 ##              Evidence groups
1167 ## yes.log      3.505116      a
1168 ## yes.1/4     3.475380      a
1169 ## yes.sqrt    3.462850      a
1170 ## yes.1/3     3.459842     ab
1171 ## no.1/4      3.435104     ab
1172 ## no.1/3      3.410377     ab
1173 ## no.log      3.408959     ab
1174 ## no.1/100    3.360803     ab
1175 ## no.P/A      3.357034     ab
1176 ## no.1/1000   3.356174     ab
1177 ## yes.1/100   3.330693     ab
1178 ## yes.P/A     3.329233     ab
1179 ## yes.1/1000  3.309273     ab
1180 ## no.sqrt     3.297212     ab
1181 ## yes.none    3.225800      b
1182 ## no.none     2.859117      c

1183 demonstrating that N2-processing works and that the choice of prior abundance
1184 transformation (log, square-, cube- or quarter- root) is of minor importance.

```

1185 Discussion

1186

1187 While informative about which differences have no statistical meaning, the regression  
1188 models are, of course, simplistic as they focuss on the mean evidence expressed as -  
1189  $\log(P)$ , i.e. on the geometric mean  $P$ , only. Quantile regression might also be of  
1190 interest.

1191 Reference

1192 de Mendiburu F. (2023). agricolae: Statistical Procedures for Agricultural Research, R  
1193 package version 1.3-7. <https://doi.org/10.32614/CRAN.package.agricolae>.

1194 **8.6 Subdividing the CESTES data base by type**

1195 Here in CESTEST data base is subdivided by response type, ecosystem type, and  
1196 taxonomic group to examine whether the preferred transformation (prior-abundance  
1197 transformation with or without  $N_2$ -processing) varies across these categories. It should  
1198 be noted that several of these subsets contain relatively few datasets, and the  
1199 corresponding results should therefore be interpreted with caution and not  
1200 over-generalized.

1201

1202 Type of abundance measurement

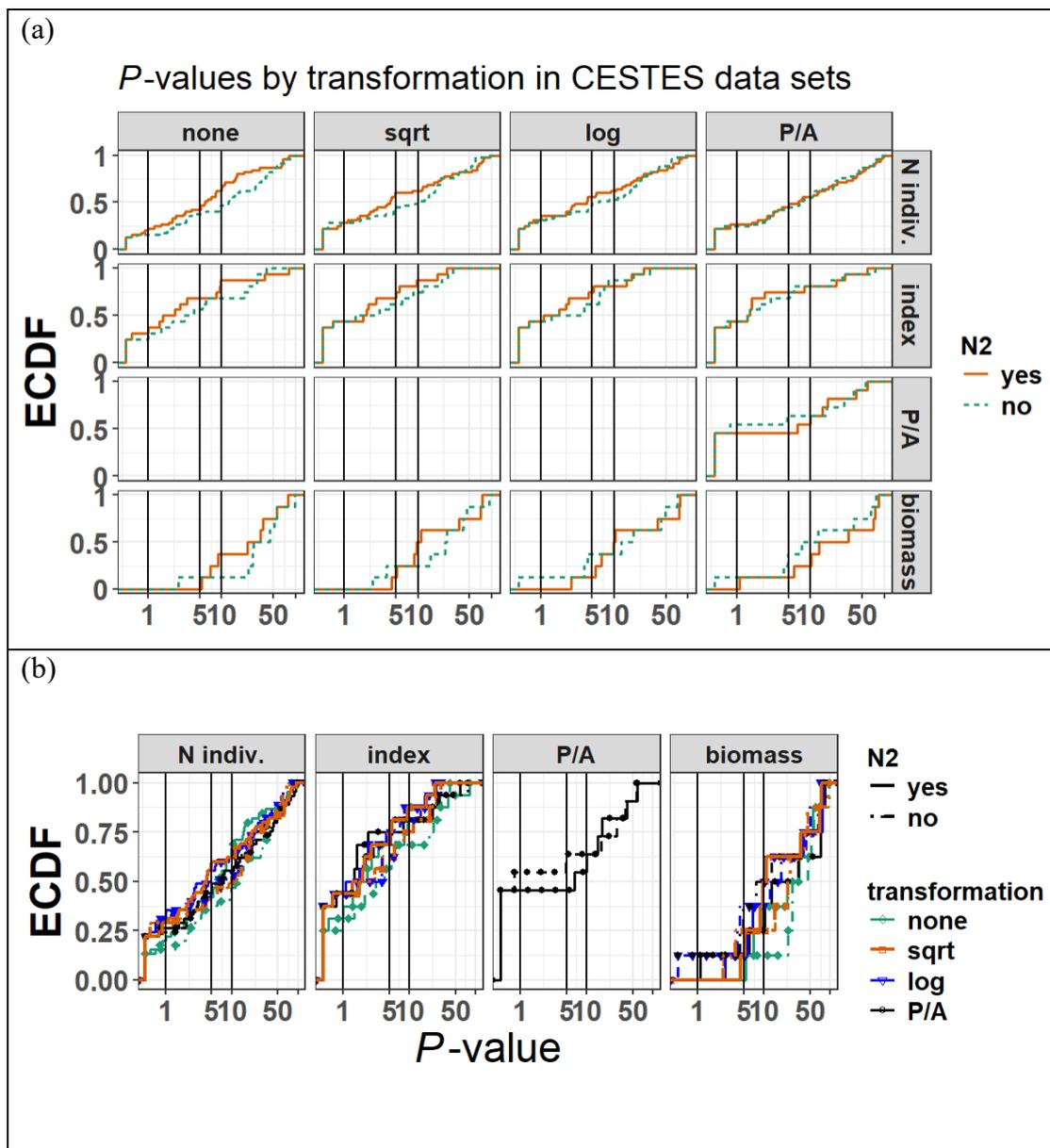
1203 Perhaps the best transformation (prior abundance transformation with or without  $N_2$ -  
1204 processing) differs between the type of response, *i.e.* how abundance was measured.  
1205 Fig. S8a focusses on the effect of  $N_2$ -processing per types of prior transformation and  
1206 prior transformation. Fig. S8b directly compares the different prior abundance  
1207 transformations, whereas Fig. S8c compares types of abundance measurement for  
1208 square-root transformed abundance with and without  $N_2$ -processing.

1209 In data sets in which the abundance is expressed as (average) number of individuals or  
1210 an index,  $N_2$ -processing yielded generally more evidence for trait-environment  
1211 association (lower  $P$ -values) with diminishing effect with stronger prior abundance  
1212 transformation (from none to P/A, Fig. S8a).

1213

1214 Pre-processing resulted in generally lower  $P$  values (more evidence for trait-  
 1215 environment association), with prior square-root and log-transformation being best,  
 1216 irrespective of the type of response (Fig. S8b), with the exception of response coded  
 1217 as presence-absence.  $N_2$ -processing on the raw abundance, when recorded as number  
 1218 of individuals showed stronger evidence (lower  $P$  values) for  $P \geq 0.10$  than with prior  
 1219 abundance transformation. For data sets that recorded presence-absence only all pre-  
 1220 processing methods show the same EDCF, except that  $N_2$ -processing showed lower  $P$   
 1221 values for  $P > 0.10$ . Apparently, only in data sets with weak to very weak evidence  
 1222 for trait-environment association contain species that are very common, so that  
 1223 informativeness works better than simply the number of occurrences. Studies using  
 1224 indices of abundance showed generally higher evidence than other measures of  
 1225 abundance. The fraction of data sets with  $P < 0.02$  was highest in data sets using  
 1226 presence-absence.

1227



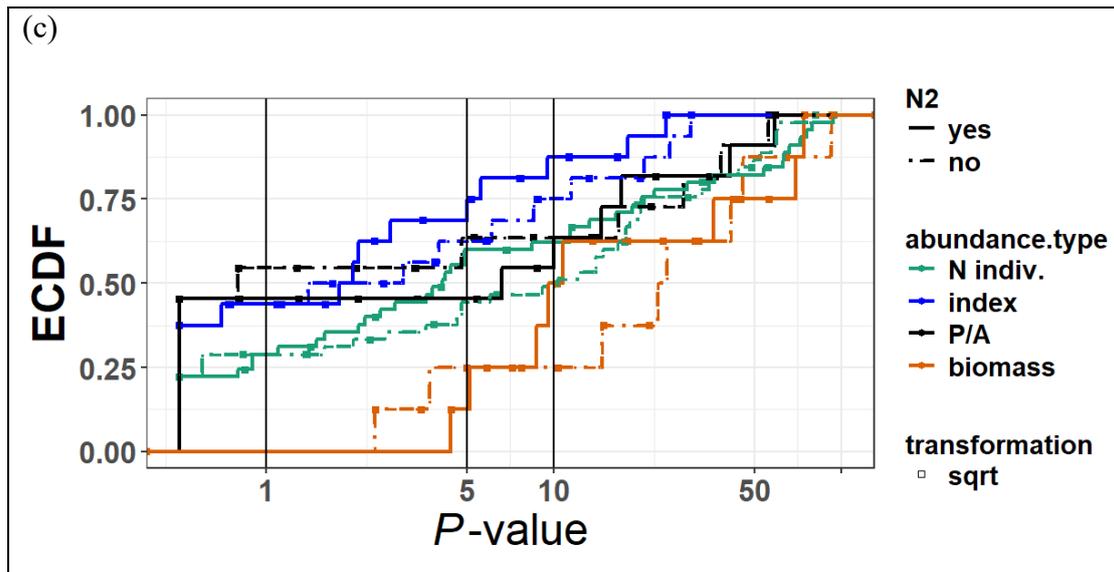


FIG. S8. Empirical cumulative density function (ECDF) of the dc-CA significance ( $P$ -value) for trait-environment association in the data sets of the CESTES database in relation to methods of pre-processing and type of abundance. Pre-processing methods are transformation of abundance (none, sqrt: square-root, log: logarithmic and to P/A: transformation to presence/absence) with and without subsequent  $N_2$ -processing ( $N_2$ ). Types of abundance are number of individuals ( $N=45$  datasets) and abundance index ( $N=16$ ), presence-absence ( $N=11$ ) and biomass ( $N=8$ ). (a) Effect of  $N_2$ -processing on  $P$ -values by prior transformation and type of abundance. (b) A  $P$ -value comparison of prior transformations per type. (c) A  $P$ -value comparison of types with square-root abundance transformation. The abscissa starts at  $P=0.005$  (0.5%). The vertical lines are at the 1, 5 and 10% significance levels.

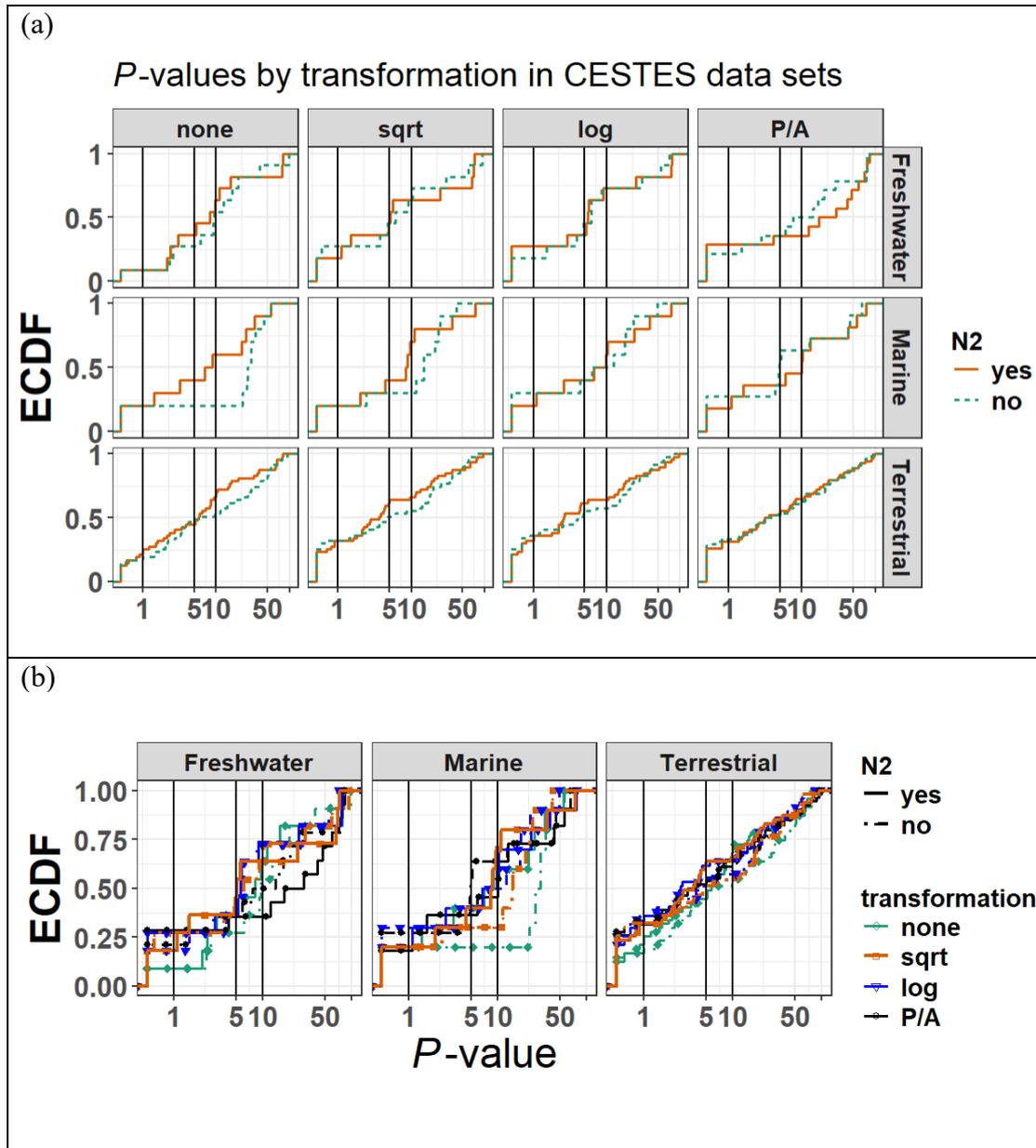
1228

1229

1230 Type of ecosystem  
 1231 Perhaps the best transformation (prior abundance transformation with or without  $N_2$ -  
 1232 processing) differs between the type of ecosystem. Fig. S9a focusses on the effect of  
 1233  $N_2$ -processing per type of ecosystem and prior transformation. Fig. S9b directly  
 1234 compares the different prior abundance transformations, whereas Fig. S9c compares  
 1235 ecosystem types for square-root transformed abundance with and without  $N_2$ -  
 1236 processing.

1237 The main patterns are similar to those in the main text with an effect of  $N_2$ -processing  
 1238 that diminishes with stronger prior abundance transformation and good performance  
 1239 of square-root and log-transformation, with little improvement due to  $N_2$ -processing  
 1240 (Fig. 9a). The effect of  $N_2$ -processing was largest in Marine data sets, but slightly  
 1241 negative with transformation to presence-absence (Fig. 9a). The fraction of data sets  
 1242 showing strong evidence ( $P < 0.01$ ) for trait-environment association was highest in  
 1243 Terrestrial studies (Fig. 9c).

1244



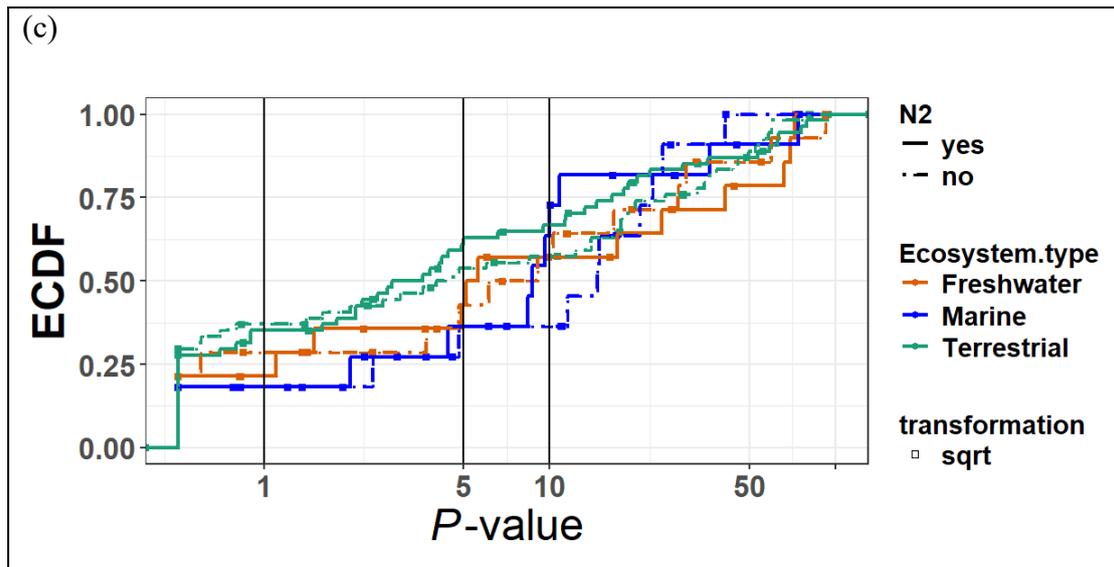


FIG. S9. Empirical cumulative density function (ECDF) of the dc-CA significance ( $P$ -value) for trait-environment association in the data sets of the CESTES database in relation to methods of pre-processing and type of ecosystem. Pre-processing methods are transformation of abundance (none, sqrt: square-root, log: logarithmic and to P/A: transformation to presence/absence) with and without subsequent  $N_2$ -processing ( $N_2$ ). Types of ecosystems are Freshwater ( $N=14$  datasets) and marine ( $N=11$ ) and terrestrial ( $N=54$ ). (a) Effect of  $N_2$ -processing on  $P$ -values by prior transformation and type of ecosystem. (b) A  $P$ -value comparison of prior transformations per type. (c) A  $P$ -value comparison of types. The abscissa starts at  $P = 0.005$  (0.5%). The vertical lines are at the 1, 5 and 10% significance levels.

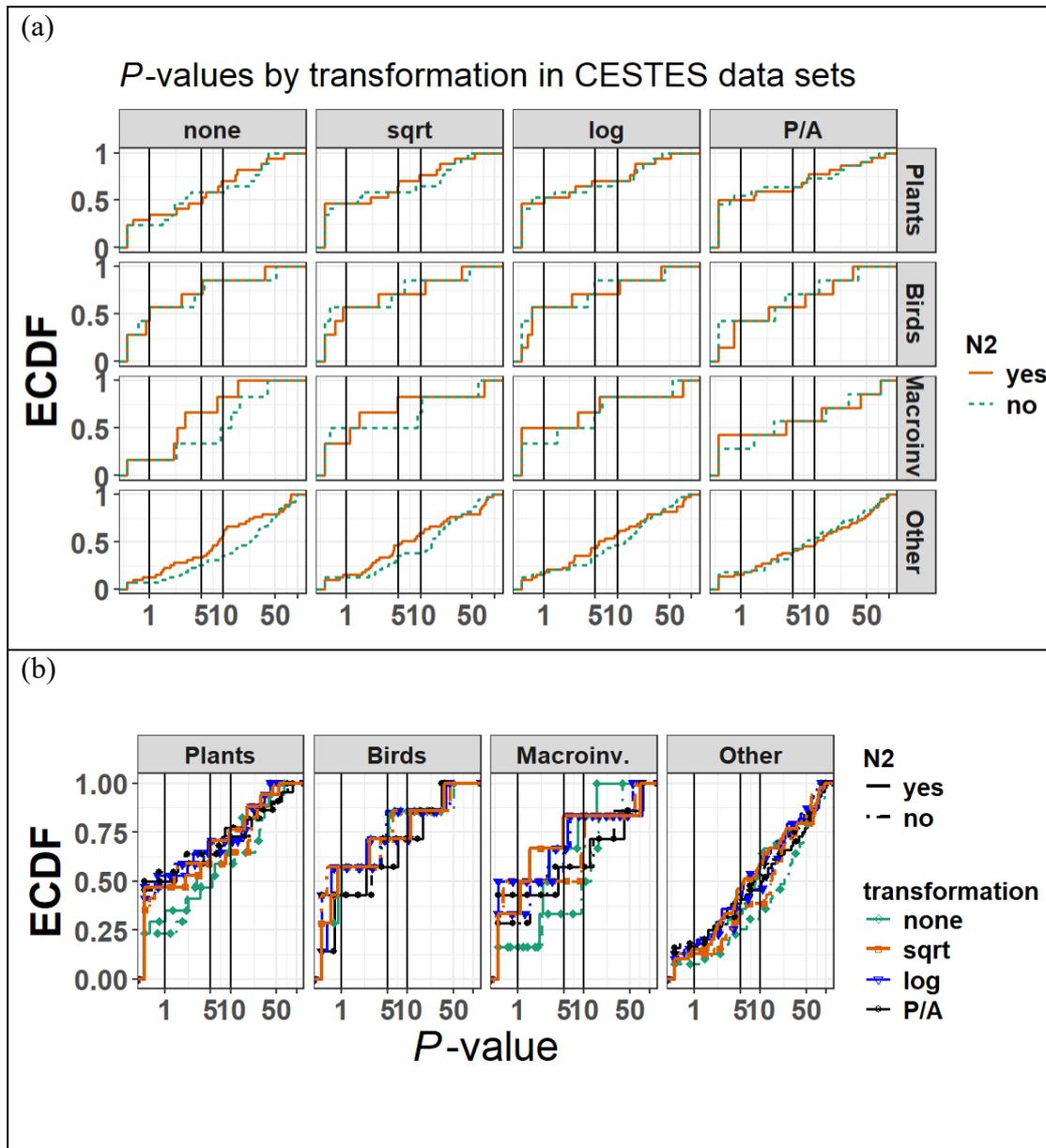
1245

1246

1247 Taxonomic group  
 1248 Perhaps the best transformation (prior abundance transformation with or without  $N_2$ -  
 1249 processing) differs between taxonomic groups. Fig. S10a focusses on the effect of  $N_2$ -  
 1250 processing per taxonomic group and prior transformation. Fig. S10b directly  
 1251 compares the different prior abundance transformations, whereas Fig. S10c compares  
 1252 taxonomic groups for square-root transformed abundance with and without  $N_2$ -  
 1253 processing.

1254 In the Plant and Bird data sets the effect of  $N_2$ -processing is small (Fig. 10a). The  
 1255 evidence for trait-environment association is higher in the Plant and Birds data sets  
 1256 than in the data sets on Macroinvertebrates, with least evidence in ‘Other’ consisting  
 1257 of 24 different taxonomic groups.

1258



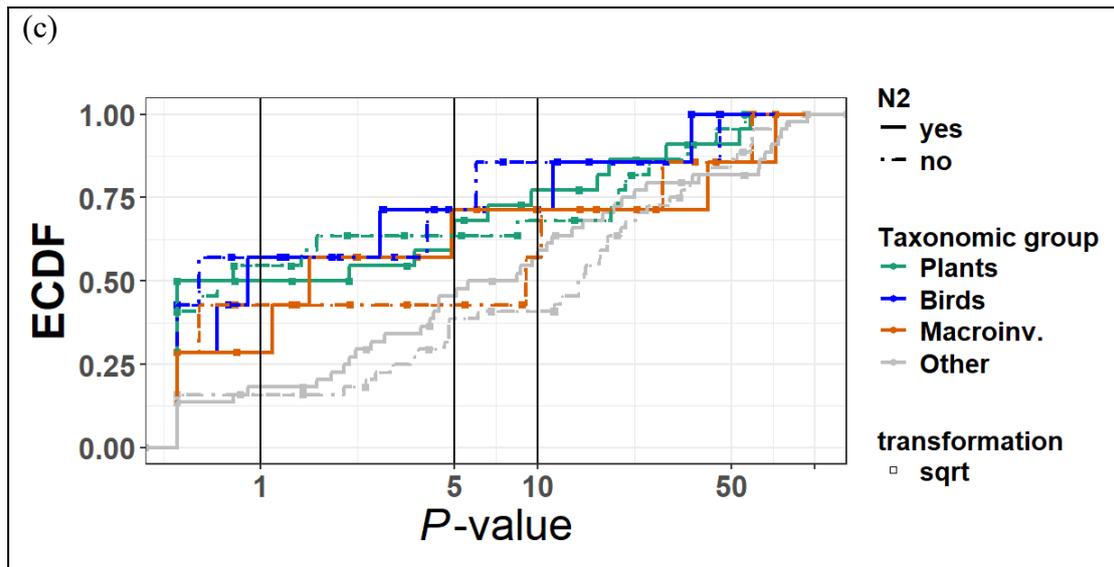


FIG. S10. Empirical cumulative density function (ECDF) of the dc-CA significance ( $P$ -value) for trait-environment association in the data sets of the CESTES database in relation to methods of pre-processing and taxonomic group. Pre-processing methods are transformation of abundance (none, sqrt: square-root, log: logarithmic and to P/A: transformation to presence/absence) with and without subsequent  $N_2$ -processing ( $N_2$ ). Taxonomic groups: Plants ( $N=22$  datasets) and Birds ( $N=7$ ), Macroinvertebrates ( $N=7$ ) and Other ( $N=44$ ). (a) Effect of  $N_2$ -processing on  $P$ -values by prior transformation and taxonomic group. (b) A  $P$ -value comparison of prior transformations per group. (c) A  $P$ -value comparison of taxonomic groups. The abscissa starts at  $P = 0.005$  (0.5%). The vertical lines are at the 1, 5 and 10% significance levels.

1259

1260

1261