

1 Running head: Bias in CWM analysis

2 **BIAS IN COMMUNITY-WEIGHTED MEAN ANALYSIS RELATING SPECIES ATTRIBUTES**  
3 **TO SAMPLE ATTRIBUTES: JUSTIFICATION AND REMEDY**

4

5 David Zelený

6

7 Institute of Ecology and Evolutionary Biology, National Taiwan University, No. 1, Sec. 4,

8 Roosevelt Rd., Taipei 10617, Taiwan, and

9 Department of Botany and Zoology, Faculty of Science, Masaryk University, Kotlářská 2, Brno

10 61137, Czech Republic

11

12 email: [zeleny.david@gmail.com](mailto:zeleny.david@gmail.com)

13

14 **# Words: 10 718**

15 **# References: 50**

16

## 17 **Abstract**

18 A common way to analyse relationship between matrix of species attributes (like functional traits  
19 of indicator values) and sample attributes (e.g. environmental variables) via the matrix of species  
20 composition is by calculating community-weighted mean of species attributes (CWM) and  
21 relating it to sample attributes by correlation, regression, ANOVA or other method. This  
22 *weighted-mean approach* is used in number of ecological fields (e.g. functional and vegetation  
23 ecology, biogeography, hydrobiology or paleolimnology), and represents an alternative to other  
24 methods relating species and sample attributes via species composition matrix (like the fourth-  
25 corner problem and RLQ analysis).

26 Here, I point out two important problems of weighted-mean approach: 1) in certain cases,  
27 which I discuss in detail, the method yields highly biased results in terms of both effect size and  
28 significance of the relationship between CWM and sample attributes, and 2) this bias is  
29 contingent upon beta diversity of species composition matrix. CWM values calculated from  
30 samples of communities sharing some species are not independent from each other and this lack  
31 of independence influences the number of effective degrees of freedom. This is usually lower  
32 than actual number of samples entering the analysis, and the difference further increases with  
33 decreasing compositional heterogeneity of the dataset. Discrepancy between number of effective  
34 degrees of freedom and number of samples in analysis turns into biased effect sizes and inflated  
35 Type I error rate in case that significance of the relationship is tested by standard tests, a problem  
36 which is analogous to analysis of two spatially autocorrelated variables.

37 Consequences of the bias is that reported results of studies using rather homogeneous  
38 (although not necessarily small) compositional datasets may be overly optimistic, and results of

39 studies based on datasets differing by their compositional heterogeneity are not directly  
40 comparable. I describe the reason for this bias and suggest guidelines how to decide in which  
41 situations the bias is actually a problem for interpretation of results. I also introduce analytical  
42 solution accounting for the bias, test its validity on simulated data and compare it with an  
43 alternative approach based on the *fourth-corner* approach.

44

## 45 **Introduction**

46 *Weighted-mean approach* is a method to analyse link between species attributes and sample  
47 attributes by calculating community-weighted means of species attributes (CWM), which can be  
48 directly related to sample attributes by correlation, regression, ANOVA or other methods (Fig. 1).  
49 *Species attributes* are species properties (traits), behaviour (species ecological optima) or  
50 phylogenetic age, while *sample attributes* are characteristics of community samples measured in  
51 the field (environmental variables) or derived from matrix of species composition (species  
52 richness or positions of samples in ordination diagrams).

53         Weighted-mean approach is used in wide range of ecological fields. In functional ecology,  
54 testing the effect of environmental variables on changes in CWM is one of the approaches  
55 demonstrating the effect of environmental filtering on trait-mediated community assembly (Diaz  
56 *et al.* 1998; Shipley 2010; Laliberté *et al.* 2012). Similarly, CWM are used to predict changes in  
57 ecosystem properties, such as biomass production or nutrient cycling (Garnier *et al.* 2004; Vile *et*  
58 *al.* 2006), or ecosystem services like fodder production or maintenance of soil fertility (Diaz *et al.*  
59 2007). In biogeography, grid-based means of species properties (like animal body size, range  
60 size or evolutionary age) are linked to macroclimate or diversity (Blackburn & Hawkins 2004,

61 Hawkins & Diniz-Filho 2006, Hawkins et al. 2014). Vegetation ecologists use species indicator  
62 values (e.g. those of Ellenberg et al. 1992 or Landolt 1977) to estimate habitat conditions from  
63 calculated mean species indicator values of vegetation samples and relate them to soil, light or  
64 climatic variables (Schaffers & Sýkora 2000, Wamelink et al. 2002, 2005). In hydrobiology,  
65 reliability of saprobic index of Sládeček (1973) based on weighted mean of diatom indicator  
66 values, or similar indices (e.g. trophic diatom index, Kelly & Whitton 1995) is evaluated by  
67 relating them to measured water quality parameters. Similarly, in paleoecology the method used  
68 to reconstruct acidification of lakes from fossil diatom assemblages preserved in lake sediments  
69 is based on weighted means of diatom optima along pH gradient (ter Braak & Barendregt 1986)  
70 and as one of the transfer functions, (e.g. Birks *et al.* 1990) is considered to be a tool which have  
71 “revolutionised paleolimnology” (Juggins 2013). Other, more specific examples include relating  
72 community specialization index (mean of species specialization values weighted by their  
73 dominance in community) to environmental variables (Clavero & Brotons 2010, Fajmonová et al.  
74 2013), or attempts to verify whether plant biomass can be estimated from tabulated plant heights  
75 and species composition as mean of species heights weighted by their cover in a plot (Axmanová  
76 *et al.* 2012).

77       Important thing to note is that although weighted-mean approach is technically relating  
78 two sets of variables (CWM and sample attributes), three matrices are in fact involved in the  
79 computation background (notation here follows RLQ analysis of Dolédec *et al.* 1996): matrix of  
80 *sample attributes* **R** with *m* sample attributes of *n* samples ( $n \times m$ ), matrix of *species composition*  
81 **L** with abundances (or presences-absences) of *p* species in *n* samples ( $n \times p$ ) and matrix of  
82 *species attributes* **Q** with *s* species attributes for *p* species ( $s \times p$ ). Weighted-mean approach is

83 just one of possible options how to tackle the problem of relating species attributes (**Q**) with  
84 sample attributes (**R**) via matrix of species composition (**L**): it combines **Q** with **L** into matrix of  
85 weighted-means **M** and relates it to **R** (Fig. 1). Alternative solution is to combine matrix of  
86 sample attributes **R** with species composition **L** by calculating weighted-mean of sample  
87 attributes (optima of individual species along given sample attribute or species centroids) and  
88 relate these values to species attributes **Q** (e.g. ter Braak & Looman 1986). Third option is to use  
89 methods suitable for simultaneously handling all three matrices (**R**, **L** and **Q**), such as the *fourth-*  
90 *corner approach* (Legendre *et al.* 1997), related ordination method called RLQ analysis  
91 (Dolédec *et al.* 1996) and other alternatives (Jamil *et al.* 2013, Brown *et al.* 2014).

92 In weighted-mean approach, relationship between CWM and sample attributes is in most  
93 cases tested by standard parametric or permutation test. However, not all types of ecological  
94 questions, which are usually solved by weighted-mean approach, should actually be tested by  
95 standard tests. In certain situations and types of null hypotheses, weighted-mean approach  
96 combined with standard tests generates biased results, which are more optimistic than would be  
97 actually warranted by analysed data. This bias includes unreliable estimates of effect size (e.g.  
98 correlation coefficients in case of correlation or  $r^2$  in case of linear regression) and inflated Type  
99 I error rate, leading to more frequent rejection of the null hypothesis than would be expected.  
100 The key point before applying the weighted-mean approach is to explicitly decide, based on  
101 critical inspection of the context of the study question and tested null hypothesis, what is actually  
102 the relationship between species attributes or sample attributes and species composition, and  
103 which of these relationships is actually fixed and which is random (more on the terms “fixed”  
104 and “random” below). Inspiration for this issue can be seen in application of the fourth-corner

105 approach (Legendre et al. 1997), for which Dray & Legendre (2008) demonstrated the problem  
106 of deciding the right permutation test (from five permutation models) to test the actual question  
107 in hand, with a risk of inflated Type I error rate in situation of wrong choice. For *weighted-mean*  
108 *approach*, this issue was shown by Zelený & Schaffers (2012) in a specific context of relating  
109 mean Ellenberg indicator values (species attributes) to sample attributes derived from species  
110 composition matrix (like ordination scores or species richness), and also by Peres-Neto et al.  
111 (2012) in the context of metacommunity phylogenetics. Both studies proposed numerical  
112 solutions: Zelený & Schaffers (2012) introduced modified permutation test, based on permuting  
113 species instead of sample attributes, and Peres-Neto et al. (2012) suggested to use *sequential test*  
114 (ter Braak et al. 2012) using the *fourth-corner* statistic (Legendre et al. 1997). Additionally,  
115 Šmilauer & Lepš (2014) touched on this issue in the context of CWM-RDA method (Kleyer et al.  
116 2012).

117 In this study, first I review categories of questions and null hypotheses which are  
118 commonly analysed by weighted-mean approach. Using simulated data, I show for which  
119 category there is a risk of biased results if tested by standard tests, and describe in detail what  
120 exactly causes this bias. Namely, I argue that the bias is caused by mismatch between number of  
121 samples in weighted-mean analysis and actual number of effective degrees of freedom, since  
122 community samples sharing some of the species with other samples do not count for the full  
123 degree of freedom in this analysis. I will also demonstrate that the amount of bias depends on the  
124 compositional heterogeneity (beta diversity) of the species composition matrix in a way that with  
125 increasing heterogeneity the bias decreases, which makes comparison of results between datasets  
126 of different compositional heterogeneity difficult. Note that for numerical simplicity, I ignore

127 intraspecific variation in species attributes. Finally, I will review methods available for solving  
128 the problem of inflated Type I error rate in weighted mean approach (namely *modified*  
129 *permutation test* of Zelený & Schaffers 2012 and solution based on combining *fourth-corner*  
130 *statistic* and *sequential permutation test* in weighted regression as introduced by Peres-Neto et al.  
131 2012) and introduce novel solution, here called *two-step permutation test*. Although the  
132 examples, ecological interpretations and reasoning used here are focused on relationship of  
133 species functional traits or species indicator values with sample attributes analysed by weighted-  
134 mean approach, the general context is valid also for other types of species and sample attributes  
135 linked by weighted-mean approach.

136

### 137 **Types of species and sample attributes**

138 When thinking about possible alternative types of questions which are commonly being analysed  
139 using weighted-mean approach, it proves as useful to distinguish whether species and sample  
140 attributes are fixed or random. Terminology behind distinction into fixed and random is diverse  
141 (Gelman 2005); here I don't use this terms in the sense of ANOVA (fixed vs random terms),  
142 neither in the sense inferring their importance (fixed factors are important while random are  
143 nuisance factors). As *fixed attributes* I consider those which are specific for given dataset and  
144 related to given species or samples, and this specificity is acknowledged by the  
145 question/hypothesis being tested in a way that this link is deemed as given, and not further  
146 questioned or tested. *Random attributes* represent a subset of some larger pool of values, and  
147 their link to species composition is not acknowledged by the question being tested. In the narrow  
148 sense of permutation tests, fixed attributes should not be permuted among each other, while

149 random attributes can. For interpretation, effect of fixed attributes is limited only for given set of  
150 attribute values and in the context of community datasets included in the analysis and is not  
151 likely to be generalized beyond, while in case of random attributes interpretation is focused on  
152 the general effect of given attribute, not only the subset of attribute values used in the study.  
153 Species traits measured on individuals from plots of given community datasets can be considered  
154 as fixed, while species traits taken from large trait databases and measured often in completely  
155 different context should more likely be considered as random. Similarly, species richness or  
156 sample ordination scores derived from community matrix are more likely fixed sample attribute,  
157 while environmental variables measured in the field or derived from GIS layers may be  
158 considered as random. Indeed, this distinction is often dependent on author's view of the  
159 problem and on the theoretical context of the study, and the same variables can be seen as fixed  
160 or random in different context.

161         In the original description of the fourth-corner problem (Legendre et al. 1997), focused  
162 on linking fish behavioral and biological characteristics to environmental variables, both species  
163 and sample attributes have been considered as fixed, and random was matrix of species  
164 composition. Alternative permutation models were then used to test alternative hypotheses about  
165 the mechanisms assembling the community (like environmental control over individual species  
166 or species assemblages, lottery or random species attributes). In weighted-mean approach,  
167 decision about fixed or random nature of attributes directly influences the decision for  
168 meaningful way to test the relationship, and is therefore crucial for selecting correct statistical  
169 test. Important is to note that all hypotheses tested by weighted-mean approach make (implicit or  
170 explicit) assumption that either species or sample attributes are fixed, and their link to species



171 composition is *a priori* acknowledged and therefore not questioned (and not tested). In some  
172 cases, whether the species or sample attributes are better considered as fixed or random depends  
173 also on the need to generalize the results of the study. If the results should be used e.g. for local  
174 application (e.g. whether CWM of species height in given agricultural system can predict well  
175 the harvested biomass), species attributes can be seen as fixed, since the results will be used  
176 solely in the context of the studied system – if the same values are measured again for the same  
177 community, the results should be similar, but not applicable generally to other communities. If  
178 the aim is to generalize the results (e.g. to assess whether the species height itself, as tabulated in  
179 the floras, can be used as a tool to predict biomass yield), it is more reasonable to treat the  
180 species attributes as random and modify the analysis accordingly, so as even local study can  
181 contribute to more general description of such pattern.

182         Another useful distinction of attribute types is whether the sample or species attributes  
183 are internal or external. The main difference is that internal attributes are numerically derived  
184 from matrix of species composition, while external attributes are typically measured or estimated  
185 variables, not directly derived from species composition matrix. Internal species attributes are  
186 species optima calculated as weighted-means of sample attributes or as species scores on  
187 ordination axes, and similarly internal sample attributes are e.g. sample scores on ordination axes,  
188 species richness of individual samples or assignment of samples into groups based on  
189 compositional similarity (e.g. by numerical classification). External species attributes, in contrary,  
190 are measured traits or tabulated species indicator values, external sample attributes are measured  
191 or estimated environmental variables or experimental treatments. While the link of external  
192 species or sample attributes to species composition may be fixed or random and depends on the

193 context, internal attributes should always be considered as fixed, since they refer only to the  
194 context of the dataset from which they have been derived.

195

## 196 **Types of hypotheses tested by weighted-mean approach**

197 Considering the distinction between fixed and random (sample or species) attributes, questions  
198 and hypotheses commonly tested by weighted-mean approach fall into one of the three  
199 categories (see Table 1 for summary). *Category 1* assumes that while sample attributes are fixed,  
200 species attributes are random; *category 2* is opposite to the previous, with sample attributes  
201 considered random and species attributes fixed; and, finally, *category 3* assumes that both  
202 species and sample attributes are random. Below, I review in detail individual categories with  
203 examples of ecological questions/hypotheses for each of them.

### 204 *Category 1 – species attributes are random, sample attributes are fixed*

205 Hypotheses in this category explicitly acknowledge the link between sample attributes and  
206 species composition, or the link is implicit from the context or numerical background of the  
207 study, and they focus on testing the link of species attributes to species composition. The null  
208 hypothesis which is tested states that species attributes are not linked to species composition,  
209 while alternative hypothesis states that they are. Questions focused on relating CWM to internal  
210 sample attributes derived computationally from matrix of species composition fall into this  
211 category (e.g. relating mean Ellenberg indicator values to sample scores on unconstrained  
212 ordination, often used to interpret ecological meaning of ordination axes; Zelený & Schaffers  
213 2012). Also studies with external sample attributes considered to be fixed, like experimental

214 treatments, fall into this category in case that their effect on species composition is  
215 acknowledged, and the question is about how does species attributes response to it. Additional  
216 level of complexity is added in studies dealing with grid data with both CWM and internal  
217 sample attributes (e.g. species richness derived from community data) spatially autocorrelated  
218 due to spatial coherence of species distribution (B. Hawkins, *pers. comm.*). Zelený & Schaffers  
219 (2012) showed that standard parametric and permutation tests has inflated Type I error rate for  
220 this category, and as an alternative introduced *modified permutation test*, permuting species  
221 attributes instead of sample attributes as further discussed in this study.

222 *Category 2 – species attributes are fixed, sample attributes are random*

223 Hypotheses in the second category assume that the species attributes are linked to species  
224 composition, and the null hypothesis states that sample attributes are not linked to species  
225 composition, while alternative hypothesis states that they are. Example are trait-based studies  
226 asking whether species traits can explain effect of environmental filtering on species abundances  
227 in community; these studies operate with an assumption that species traits (as species attributes)  
228 are functional, i.e. they influences the abundances of species in community, and the question  
229 being evaluated is whether the sample attributes (environmental factor) acts as an environmental  
230 filter on species abundances. Also descriptive studies without ambitions to be more generalized  
231 fall into this category – e.g. relationship between CWM of species indicator values (e.g. mean  
232 Ellenberg indicator values) and measured environmental variables, if the interpretation is  
233 restricted only for the community dataset included in the study. Finally, studies using internal  
234 species attributes (derived from species composition, e.g. as weighted-mean of sample attributes  
235 or as scores on ordination axes) also belong to this category.

236 *Category 3 – both species and sample attributes are random*

237 This category of hypotheses includes mostly observational studies without prior knowledge or  
238 expectations about link between any of matrices. The null hypothesis states that there is no link  
239 between species and sample attributes via the matrix of species composition because either  
240 species attributes or sample attributes (or both) are not linked to species composition. To reject  
241 this null hypothesis means to prove that both species and sample attributes are actually linked to  
242 species composition. Empirical studies describing general relationship between sample attributes  
243 and species attributes without explicitly or implicitly acknowledging some underlying  
244 assumptions or mechanisms belong to this category. Examples are studies relating CWM of  
245 functional traits to environmental variables without clear assumption that traits are functional,  
246 allowing to question whether particular trait are actually linked to species composition or not. In  
247 case of studies with species indicator values, these include relating mean indicator values to  
248 environmental variables with aim to generalize the result also out of the studied dataset (e.g.  
249 answering the question whether Ellenberg indicator values for soil reaction *per se* are good  
250 predictors of measured soil pH, i.e. not only in the context of given community dataset).

251

252 **Illustration of the bias using simulated community data**

253 In the next section, I will use simulated community data to illustrate performance of standard  
254 parametric test, if this is used to test hypotheses from each category defined above. The benefit  
255 of simulated data is the possibility to keep certain parameters fixed and to manipulate only those  
256 parameters whose effect is studied – in this case fixed are numbers of samples in the dataset, and

257 manipulated are links between species or sample attributes and species composition, and also  
258 compositional heterogeneity of community data.

259 Each artificial community dataset includes the set of three matrices (sample attributes **R**,  
260 species composition **L**, and species attributes **Q**), with the link between species or sample  
261 attributes and species composition (or both) broken by permutation of attributes. This creates  
262 four scenarios (Fig. 2, identical with scenarios 1-4 of Dray & Legendre 2008): *scenario 1* - both  
263 sample and species attributes are linked to species composition, *scenario 2* - sample attributes  
264 are linked to species composition, but species attributes are not, *scenario 3* - species attributes  
265 are linked to species composition, but sample attributes are not, and *scenario 4* - none of species  
266 or sample attributes are linked to species composition. For hypotheses in category 1 defined  
267 above, the scenario 2 represents the null hypothesis, for category 2 scenario 3 is the null  
268 hypothesis, and for category 3 the scenarios 2, 3 and 4 represent alternative states of null  
269 hypothesis (Table 1). Scenario 1 represents the power test for all three categories (i.e. it measures  
270 probability of getting significant results if the alternative hypothesis is true). Additionally, I also  
271 examined how observed bias depends on compositional heterogeneity (beta diversity) of  
272 community matrix, which influences the number of effective degrees of freedom in analysis (as  
273 explained in detail further in the section *Justification of the bias*). Note that all analyses in this  
274 paper were conducted using R-project (v. 3.2.3, R Core Team 2015); complete R scripts are  
275 available in Appendix S5, and all functions have been wrapped into R-packages *weimea*  
276 (abbreviation for *weighted mean*; source code for v. 0.58 available as Appendix S6).

277 *Description of simulated data*

278 I created an algorithm generating community data which are structured by two virtual ecological  
279 gradients (by extending the original one-gradient algorithm of Fridley et al. 2007 based on  
280 concept of Minchin 1987). The first gradient has *constant* length for all generated datasets and  
281 serves as a surrogate for measured environmental variable; in the analysis, positions of samples  
282 along this gradient are used as *sample attributes*, while the optima of species along this gradient  
283 are used as *species attributes*. The length of the second gradient is *variable* and its increasing  
284 length increases the compositional heterogeneity of species composition matrix. Community  
285 samples based on this simulated community ecospace were created by randomly locating  
286 samples along the first gradient, and species composition for each sample was derived by  
287 random assignment of fixed number of individuals to species identities weighted by relative  
288 abundances of species with non-zero probability of occurrence at given location of the gradient  
289 (see Appendix S1 for further details). With short second gradient, the resulting simulated  
290 community dataset was compositionally relatively homogeneous, with samples located nearby  
291 along the first gradient (with similar value of sample attribute) sharing rather high proportion of  
292 species. Increasing the length of the second gradient increased the compositional heterogeneity  
293 of the dataset, since two nearby samples may have quite different species composition (Fig. S1 in  
294 Appendix S1). Note that although scenarios 1-4 are conceptually analogous to scenarios 1-4 in  
295 Dray & Legendre (2008) used in the context of the *fourth-corner approach*, the model  
296 generating simulated communities is different, since Dray & Legendre (2008) used one-gradient  
297 model, which generates rather homogeneous communities, while I used two-gradients model,  
298 generating set of communities of increasing compositional heterogeneity.

299 *Weighted-mean approach with standard parametric tests applied on simulated data*

300 Using the algorithm described above, for each of the four scenarios (1-4) I created ten levels of  
301 compositional heterogeneity, and for each combination of *scenario*  $\times$  *level of heterogeneity* I  
302 created 1000 datasets (4 scenarios  $\times$  10 levels of heterogeneity  $\times$  1000 replications = 40 000  
303 datasets). For each dataset I calculated CWM of species attributes, related it to sample attributes  
304 using Pearson's  $r$  correlation and tested its significance using parametric  $t$ -test (for additional  
305 results for least-square regression and  $r^2$  see Fig. S2 in Appendix S2). For each level of  
306 community heterogeneity in each scenario I counted proportion of correlations significant at  $\alpha =$   
307 0.05 (note that this proportion is identical to the proportion of significant regressions).

308 From the three scenarios with no direct link between species and sample attributes  
309 (scenarios 2, 3 and 4), analysis of data generated by scenario 2 (Fig. 3) reveals the bias – the  
310 correlation coefficient deviates from zero more than in other cases (Fig. 3), and the test of  
311 significance shows inflated Type I error rate (Fig. 4). This bias is decreasing with increasing  
312 heterogeneity of the species composition matrix (Fig. 3 & 4, Scenario 2): for the most  
313 homogeneous dataset (*number of communities* = 1), the range of Pearson's  $r$  correlation  
314 coefficients (expressed as 2.5% and 97.5% quantiles) is between -0.751 and 0.751 with 60% of  
315 correlations significant, while for the most heterogeneous dataset with high beta diversity  
316 (*number of communities* = 10) the range of Pearson's  $r$  values is between -0.381 and 0.354 with  
317 15% of correlations significant (compared to 2.5 and 97.5% quantile range values of  $r$  observed  
318 in scenarios 3 and 4 being in average between -0.278 and 0.281, with expected number of  
319 significant results being close to 5%). Similarly inflated are values of coefficient of  
320 determination ( $r^2$ ; Fig. S2, Scenario 2) calculated by least-square linear regression.

321 Dray & Legendre (2008) showed that the *fourth-corner approach*, if tested by the  
322 permutation test based on reshuffling sample attributes (or rows of species composition matrix,  
323 respectively, model 2 in their paper) also reveals inflated Type I error rate in case of scenario 2. I  
324 applied the fourth-corner analysis also on the simulated community data described above.  
325 Results show, in line with Dray & Legendre (2008), biased values of the fourth-corner statistic  
326 and inflated Type I error rate for the model 2 permutation test (Figs. S3 & S4, Appendix S2),  
327 with the bias (and inflated Type I error rate) decreasing with increasing compositional  
328 heterogeneity (Fig S3 & S4 with Scenario 2, Appendix S2).

329

### 330 **Justification of the bias**

331 Simulation study above showed that if hypotheses in category 1 or 3, for which scenario 2 is  
332 relevant, are tested by standard parametric or permutation tests, results are prone to biased  
333 estimates of model parameters and inflated Type I error rate, and the bias is contingent upon  
334 compositional heterogeneity of the community dataset. In this section, I explore the reasons for  
335 this bias, which will be used as a theoretical base for solution proposed in the next section, and  
336 also explain the link of bias to compositional heterogeneity. Note that only hypotheses in  
337 category 1 and 3 are influenced by this bias, since for the category 2 with fixed species and  
338 random sample attributes, scenario 2 is not relevant and the bias therefore does not occur.

339 A peculiar feature of CWM of species attributes is their “numerical burden”, namely that  
340 they are calculated from species attributes assigned to individual species and from species  
341 composition of individual samples, and therefore inherit part of information from both sources.  
342 The numerical difference between calculated CWM values of two community samples is indeed



343 constrained by a difference in species composition of these two samples - if they have identical  
344 species composition (or identical relative species abundances), their calculated weighted-means  
345 do not differ (because they cannot), and if their species composition differs only slightly, the  
346 difference in their weighted-mean values is likely to be small. This property of weighted-mean  
347 values, which are not independent from each other, has notable consequences for analysis with  
348 sample attributes, if these are themselves related to species composition. In fact the situation is  
349 analogous to the analysis of two spatially autocorrelated variables, just the autocorrelation is not  
350 happening in the geographical space, but in the compositional space (more about this analogy in  
351 the next section). Two values of CWM calculated from community samples sharing some  
352 species do not bring two independent degrees of freedom into analysis, because samples used for  
353 their calculation are not independent - difference in their CWM are predictable (to some extent)  
354 from difference in their species composition in a way that the more similar is species  
355 composition of two community samples, the more similar must be also their calculated CWM.  
356 This problem scales up to the dataset level: in case of two compositional datasets with the same  
357 number of samples used in weighted-mean approach, the dataset which is compositionally more  
358 homogeneous has lower number of effective degrees of freedom compared to the more  
359 heterogeneous one.

360         Although there are many ways how to quantify compositional heterogeneity of the dataset,  
361 promising is to use beta diversity measure based on Whittaker's index of association (Whittaker  
362 1952; in Legendre & Legendre 2012 as  $D_0$ ). This index can be numerically derived from  
363 differences in species composition between two calculated CWM values, and therefore quantifies  
364 the dissimilarity in species composition which is directly related to weighted-mean approach

365 analysis (Appendix S3). To obtain single value of beta diversity for given dataset, one may use  
366 beta diversity metric of Legendre & De Cáceres (2013) quantifying the variation in species  
367 composition, which can be calculated also from symmetric matrix of dissimilarities among all  
368 pairs of samples (here using Whittaker's index of association). Advantage of this beta diversity  
369 metric is that it is independent on the size of the dataset, and the underlying dissimilarity  
370 coefficient is directly related to weighted-mean approach.

371 Below, I will first illustrate what I mean by the differences in *effective number of degrees*  
372 *of freedom* in analysis. I use a simple example focused on relationship between community-  
373 weighted mean of traits with environmental variables, and compare two contrasting sampling  
374 designs differing substantially by number of effective degrees of freedom they bring to the  
375 analysis. Then, I elaborate in more detail the analogy to analysis of two spatially autocorrelated  
376 variables, since it offers deeper insight to the problem and inspiration for potential solutions.

377 *Two sampling designs: difference in effective number of degrees of freedom*

378 The following example will illustrate the mismatch between number of effective degrees of  
379 freedom and number of samples in analysis, using species functional traits. The weighted mean  
380 of species functional traits, if combined with environmental variables, is used to investigate  
381 whether the environment filters species into community via their trait properties. Assumption  
382 behind is that the traits are functional, meaning that they directly influence the probability of  
383 species occurrence in community under given environmental conditions, and this hypothesis  
384 therefore classifies into the category 2 described above. For example, plant species specific leaf  
385 area (SLA, Reich et al. 1992) is known to be related to the plant requirements for light, and  
386 shade tolerant species restricted to more shady habitats have larger thinner leaves with higher SLA

387 values (e.g., Lambers et al. 2008). One may therefore expect that light acts as an environmental  
388 filter for species entering given community, and this filtering is (at least partly) happening  
389 because of species specific SLA. To support this reasoning, let's conduct imaginary experiment:  
390 collect a dataset about composition of understory species in the forest and analyse it by the  
391 weighted-mean approach. Let's keep things simple in this example and assume that we will  
392 compare two forest types, one with open and one with closed canopy, to see whether the closed  
393 canopy filters the understory species with high SLA. When preparing the design of data  
394 collection, we have two options. The first is to choose two vegetation types (e.g. at the level of  
395 association) with contrasting canopy openness, search for them in the study area (this could be  
396 just one forest complex with mosaic of both vegetation types, or a larger region where these  
397 vegetation types occur), sample their species composition and measure the light in the canopy  
398 and SLA of individual species. The second option is to sample open- and closed-canopy forests  
399 without any restrictions about their species composition, possibly in wider area. Following the  
400 first sampling design we get dataset where samples of open-canopy forest are compositionally  
401 similar to each other, as well as samples made in closed-canopy forest (but the composition of  
402 open-canopy forest is different from the close-canopy one). In the second scenario we probably  
403 get samples which all are having rather different species composition.

404         A simplified example how the species composition of six samples collected by these two  
405 sampling designs would look likes in extreme cases is on Fig. 5. In the case of the first sampling  
406 design, three and three samples have exactly the same species composition, while in the case of  
407 the second design, species composition of each sample is distinctly different from the others (no  
408 species are shared among any pair of samples). In the first case, three and three samples have the

409 same weighted-mean values since their species composition is identical. In the second case, three  
410 and three samples have also the same weighted-mean values, not because of identical species  
411 composition, but because samples from the same canopy cover category (open or closed) are  
412 made to have similar distribution of trait values (although the species taxonomic identities are  
413 completely different).

414         Indeed, the real data would fall somewhere in between these two cases, but this example  
415 illustrates well the concept of effective degrees of freedom in weighted-mean analysis. The  
416 relationship of the light availability to weighted-mean of SLA is exactly the same in both cases  
417 (if analysed e.g. by one-way ANOVA in this situation, when environmental variable is  
418 categorical with two levels, open-canopy forest with more light and close-canopy forest with less  
419 light). The difference is in the number of effective samples, which each sampling design offers  
420 for answering our question whether light serves as an environmental filter of species occurrence  
421 in the community via species SLA. Three community samples with completely different species  
422 composition, yet similar CWM of SLA (low for open- and high for closed-canopy forest) offers  
423 considerably better information about interaction between and SLA *per se*, then do the three  
424 samples with identical species composition (and hence also identical low or high CWM of SLA).  
425 Also, what if our assumption is wrong and the SLA is in reality not a functional trait, and hence  
426 belongs to category 3 instead of category 2? Let's replace the real trait values by random one (by  
427 randomizing species attributes among species in table), calculate CWM and test their difference  
428 between open- and closed-canopy stands (and repeat this process 1000 times). In case of the first  
429 sampling design, 87% of tests detect significant difference (874 significant results of one-way  
430 ANOVA at  $P < 0.05$ ), while in case of the second sampling design the probability is near the

431 expected Type I error rate of 5% (54 significant results out of 1000). This agrees with the result  
432 of weighted-mean approach applied on simulated data saying that inflation of Type I error rate is  
433 high for homogeneous datasets (the case of the first sampling design) and decreases with  
434 increasing compositional heterogeneity (being effectively zero in dataset with all samples  
435 completely dissimilar as in case of the second sampling design).

436 *Analogy to the analysis of spatially autocorrelated variables*

437 The situation when *weighted-mean approach* is based on species composition data, in which  
438 some pairs of samples have the same or similar species composition and sample attributes are  
439 related to species composition, resembles an analysis of two spatially autocorrelated variables. In  
440 case of spatially autocorrelated variables, samples located more close to each other in  
441 geographical space have more similar values than expected if the values are randomly selected  
442 (Legendre & Legendre 2012). In case of weighted means, it is not the proximity in geographical  
443 space, but the proximity in compositional space, which reflects distances between samples  
444 expressed as their compositional dissimilarity.

445 It has been shown that when analysing two positively *spatially autocorrelated* variables,  
446 spatial autocorrelation biases the results of statistical tests, inflating the type I error rate and thus  
447 resulting into too optimistic results (Legendre 1993). The problem is not with autocorrelation of  
448 individual variables themselves, but with spatially autocorrelated residuals when analysing their  
449 relationship (e.g. by linear regression). From the point of view of the degrees of freedom,  
450 samples located nearby in geographical space are not statistically independent, behaving to a  
451 certain degree as pseudoreplications (Legendre & Legendre 2012). A new observation does not  
452 bring completely new information, because its value can be partly derived from the value of a

453 nearby site, and the effective number of samples (i.e., effective number of degrees of freedom) is  
454 lower than the real number of samples. Since for standard parametric tests the number of degrees  
455 of freedom is important for choosing the correct statistical distribution, appropriate for a given  
456 sample size, disparity between the real number and effective number of samples leads to  
457 selection of narrower confidence intervals and hence a higher probability of obtaining significant  
458 results (Bivand 1980; Legendre 1993).

459         The same reasoning applies also for analysis of two *compositionally autocorrelated*  
460 variables. Two weighted-mean values calculated from two samples with similar species  
461 composition do not bring two full degrees of freedom to the analysis, as would be case of two  
462 weighted-mean values calculated from samples with distinctively different species composition.  
463 In a simple example with two sampling designs above, if we want to know whether species SLA  
464 really increases with decreasing light in the understory, we would learn more about this from two  
465 samples in the shaded understory which have different species composition and yet both have  
466 high CWM of SLA, than from two samples which both have similar species composition (and  
467 hence also similar CWM of SLA).

468         If sample attributes are not related to species composition, than the problem with  
469 effective number of degrees of freedom is not present; although weighted-means are still  
470 compositionally autocorrelated, sample attributes are not – in case of spatially autocorrelated  
471 variables this is analogous to situation when one variable is spatially autocorrelated, but the other  
472 is not, in which case the bias caused by autocorrelation doesn't appear. The situation which  
473 requires attention due to potential bias is therefore limited to cases when sample attributes are  
474 linked to species composition (i.e. they are fixed). This is the case of all internal sample

475 attributes derived from matrix of species composition, since they are linked to matrix of species  
476 composition (fixed) due to their numerical origin, and also the case of some of external sample  
477 attributes, if these are considered to be fixed (for examples see section *Types of species and*  
478 *sample attributes*).

479

## 480 **Proposed solutions**

481 Analogy between the bias in weighted-mean approach to the bias in analysis of spatially  
482 autocorrelated variables suggests potential toolbox for solving the problem. A simple option  
483 would be to stratify the dataset to reduce redundancy in species composition among samples, i.e.  
484 from pairs of samples with similar species composition remove one of them. Although methods  
485 for stratification based on species composition are available (e.g. Lengyel et al. 2011), it  
486 potentially results into throwing out a large number of expensive data. Alternative option would  
487 be to apply some correction for effective degrees of freedom in analysis, analogously to  
488 Dutilleul's method introduced for estimating effective number of samples in case of  
489 autocorrelated variables (Dutilleul 1993). The option I will further investigate here is based on  
490 comparison of results obtained by weighted-mean approach with those generated by a null model.

### 491 *Modified permutation test: comparison with the results of a null model*

492 Comparison with results of a null model is an analogy to testing the relationship between  
493 autocorrelated variables using toroidal shift, when one variable is permuted in a way that it  
494 preserves the original degree of spatial autocorrelation (Fortin & Dale 2005). Alternatively, one  
495 can generate random variables with the same degree of spatial autocorrelation as of the original

496 variable (Deblauwe et al. 2012). In case of compositionally autocorrelated variables used for  
497 weighted-mean analysis, such variables can be generated for weighted-mean values, by  
498 calculating weighted mean from randomized (or randomly generated) species attributes. Such  
499 weighted-mean of randomized species attributes inherits the same level of compositional  
500 autocorrelation as have the weighted-mean values of the real species attributes, because they are  
501 calculated by the same algorithm from the same species composition matrix. One can generate  
502 the null distribution of a test statistic (like  $t$ -value for correlation or  $F$ -value for regression) for  
503 each weighted-mean of randomized species attributes related to original sample attributes, and  
504 compare the observed statistic (relating the weighted-mean of real species attributes to sample  
505 attributes) to this null distribution. This is identical with the modified permutation test,  
506 introduced to test the relationship between weighted mean of species attributes and sample  
507 attributes by Zelený & Schaffers (2012) in case of relating mean Ellenberg indicator values with  
508 variables derived from ordination/classification based on the same species composition dataset.

509 To illustrate behaviour of the modified permutation test, I used the set of artificial  
510 community data as above, calculated the correlation between weighted-mean of species attributes  
511 and sample attributes for all four scenarios in communities of increasing heterogeneity, and  
512 tested the significance of this correlation using modified permutation test. Results show that in  
513 contrast to standard permutation test, inflated Type I error rate in case of the scenario 2  
514 disappears (Fig 6b). At the same time, in case of scenario 3 (species composition related to  
515 species attributes, but not to sample attributes) the test is overly conservative for homogeneous  
516 datasets. Additional power analysis (Appendix S4) reveals that the modified permutation test  
517 loses the power with decreasing sample size and mainly with decreasing number of species



518 which are being permuted (Fig S5 in Appendix S4). Modified permutation test seems therefore  
519 suitable for testing hypotheses in the category 1, which assume that species attributes are random,  
520 while sample attributes are fixed (linked to species composition) and for which scenario 2 is  
521 relevant for testing the null hypothesis. It is, however, not optimal for hypotheses in the category  
522 3, which assume that both species and sample attributes are not fixed (not linked to species  
523 composition), since in scenario 3, which is also relevant as an alternative null hypothesis for this  
524 category, the results are overly conservative (although only for the most homogeneous dataset,  
525 Fig. 6).

#### 526 *Use of the fourth-corner statistic and the sequential test*

527 Dray & Legendre (2008) noted that the fourth-corner statistic  $r$ , introduced by Legendre et al.  
528 (1997), is “equal to the slope of the linear model, weighted by total species abundances, with the  
529 niche centroids as the response variable and the species trait as the explanatory variable”. This  
530 analogy was further elaborated by Peres-Neto et al. (2012, Appendix A), who presented  
531 algorithm how to use the *fourth-corner* statistic in weighted-mean approach. In short, both  $\mathbf{R}$  and  
532  $\mathbf{Q}$  matrices are first centred by weighted mean of row sums of  $\mathbf{L}$  (in case of  $\mathbf{R}$ ) and column sums  
533 of  $\mathbf{L}$  (in case of  $\mathbf{Q}$ ), and rescaled; then, the fourth-corner  $r$  statistic is the slope of regression  
534 between weighted mean of standardized  $\mathbf{Q}$  and standardized  $\mathbf{R}$ , weighted by row sums of  $\mathbf{L}$ .  
535 Advantage of the fourth-corner statistic is an option to use *sequential permutation test* introduced  
536 by ter Braak et al. (2012), which gives unbiased test of significance for all scenarios (for  
537 application on the simulated community data used above, see Fig. S3 & S4 in Appendix S2).  
538 This sequential permutation test combines results based on permuting sample attributes (model 2)  
539 and species attributes (model 4); if the first one is significant, than the second test is done, and

540 overall significance of the result is equal to the higher of these two test's  $P$ -values. Disadvantage,  
541 on the other side, is the fact that the combination of fourth-corner statistic and sequential test  
542 applies only to the regression between standardized (centred and rescaled) species and sample  
543 attributes, weighted by row sums of species composition matrix ( $\mathbf{L}$ ), and (to my knowledge) it  
544 cannot be used to test correlation, non-weighted regression or ANOVA between non-centred and  
545 standardized CWM and sample attributes.

#### 546 *Two-step permutation test*

547 As an analogy to the sequential test used together with the fourth-corner statistic, here I introduce  
548 two-step permutation test, which gives unbiased results for relationship between CWM and  
549 sample attributes for range of statistical metrics ( $t$ -value for correlation and  $F$ -value for linear  
550 regression tested here). The test is based on combination of standard and modified permutation  
551 test; while both tests give unbiased results for scenario 4, standard test gives unbiased results also  
552 for scenario 3 (in which sample attributes are not related to species composition), while modified  
553 permutation test gives unbiased results for scenario 2 (where sample attributes are related to  
554 species composition). The idea behind the two-step permutation test is to first test whether  
555 sample attributes ( $\mathbf{R}$ ) are related to matrix of species composition ( $\mathbf{L}$ ), without considering (or  
556 even knowing) the values of species attributes ( $\mathbf{Q}$ ). This could be achieved e.g. by constrained  
557 ordination, when sample attributes are used as explanatory variables explaining variation in  
558 species composition. Here I introduce more general solution (called *LR permutation test* within  
559 this paper as a notice that relationship between  $\mathbf{R}$  and  $\mathbf{L}$  matrices is tested), which can be directly  
560 connected to particular test statistic (e.g.  $t$ -value for correlation). The LR permutation test  
561 consists of the following steps: (i) generate artificial set of species attributes as species centroids

562 calculated from real sample attributes, (ii) use these species attributes to calculate CWM, (iii)  
563 calculate observed test statistic for relationship between artificial CWM and real sample  
564 attributes, and (iv) test this relationship. The test is based on comparing observed values of the  
565 test statistic (calculated in step iii) with the null distribution of the test statistic, generated in the  
566 following way: 1) randomize sample attributes, 2) use them to calculate species attributes as  
567 species centroids from weighted-means of sample attributes, 3) use these species attributes to  
568 calculate CWM, and 4) relate these calculated CWM with randomized sample attributes from  
569 step 1) to obtain the expected test statistic; repeat steps 1) to 4)  $n$ -times ( $n$  = number of  
570 permutations). If this test is significant, it means that the sample attributes are related to matrix of  
571 species composition, and relationship of CWM with sample attributes is consequently tested by  
572 modified permutation test. If the test is not significant, standard permutation test is used. When  
573 applied on the set of artificial communities used above, this sequential test gives unbiased values  
574 of Type I error rate for all three scenarios (2, 3 and 4) and for all levels of compositional  
575 heterogeneity (Fig. 6).

576

## 577 **Discussion**

578 Main motivation of this study was to show that results of weighted-mean approach critically  
579 depend on the correct decision about the test used for statistical inference. To help in this  
580 decision process, I suggested that each hypothesis can be classified into one of the three  
581 categories, given the explicit (or implicit) assumptions about the role of species and sample  
582 attributes. For each category, I suggested optimal strategy for testing the significance of  
583 relationship between CWM and sample attributes. The decision about appropriate category is not

584 always straightforward, although the decision whether species attributes should be considered as  
585 fixed or random changes classification of the hypothesis from category 2 (with recommended  
586 standard parametric or permutation test) into category 3 (with two-step permutation test). For  
587 example, trait studies, which are testing whether environment is filtering the species into  
588 community via their functional traits, routinely assume that such traits are functional, and in  
589 weighted-mean approach are therefore considered as fixed (category 2). However, this  
590 assumption may not always be justified; traits included in these analyses are often those readily  
591 available in databases and/or relatively easy to measure, but these do not necessarily need to be  
592 really the functional ones (Fox 2012, Mlambo 2014). In case of compositionally relatively  
593 homogeneous datasets, even the traits with no ecological meaning may show high and significant  
594 relationship to environmental variables if tested by standard tests. I believe that this calls for  
595 revision of such commonly applied practice.

596         Differences in effective degrees of freedom among datasets complicate comparison of  
597 results between studies based on datasets of different compositional heterogeneities. Studies  
598 conducted on datasets of relatively low beta diversity may obtain stronger and more likely  
599 significant relationship between weighted-mean of species attributes and sample attributes than  
600 studies on datasets of relatively higher beta diversity, even in case that the real link of species  
601 attributes to species composition is missing (Figs. 3 & 4, Scenario 2). This situation is analogous  
602 to biased estimates of model parameters and inflated Type I error rate in analysis of spatially  
603 autocorrelated variables. An option how to deal with this problem is to routinely report, in each  
604 case-study using weighted-mean approach, some standardized value of compositional  
605 heterogeneity. Although this would not remove bias in results of these studies, it would at least

606 allow for comparison of the potential for bias among different studies. Good metric for this  
607 purpose should be independent on the sample size, and should pertain dissimilarity in species  
608 composition which is relevant for differences in community-weighted means; here, I suggested  
609 beta diversity measure based on Whittaker's index of association (Appendix S3) following the  
610 approach summarized by Legendre & De Cáceres (2013).

611         Specific question is how to deal with missing values of species attributes for some of the  
612 species. Should species with missing species attributes remain in the matrix of species attributes  
613 and species composition? And in a case of the modified and two-step permutation tests, should  
614 the missing values be kept and permuted among species? The analogy to spatial autocorrelation  
615 issue offers clear answers for these questions. Species with missing attribute values are not used  
616 for weighted-mean calculation, so they do not contribute to the compositional autocorrelation of  
617 weighted-mean values. The point of the modified (and subsequently also two-step) permutation  
618 test is to generate random variables with the same compositionally autocorrelated structure as the  
619 weighted mean calculated from the original species attributes. For this, matrix of species  
620 composition, which inherits the compositional autocorrelation into weighted-mean values,  
621 should remain the same also for calculation of weighted-means from randomly generated species  
622 attribute values. This would not be the case if the species with missing attribute values remains  
623 in both matrices, because permuting missing values would cause the weighted mean of permuted  
624 species attributes being calculated every time with different species composition matrix (the  
625 species which in given permutation run would be assigned missing values will not be included in  
626 this weighted-mean calculation). The solution is hence to remove species with missing species  
627 attributes from both species attributes and species composition matrix, and in the case of

628 modified permutation test to permute only existing species attribute values. In case that more  
629 species attributes are analyzed (e.g. three different functional traits, or six different species  
630 indicator values) and species has missing species attribute value for some attributes and not for  
631 the others, the species should be removed from species composition matrix only for the purpose  
632 of calculating and testing weighted mean of that species attributes for which the species value is  
633 missing, and not for the others.

634 Power test using simulated dataset showed that the power of both two-step as well as  
635 modified permutation test decreases with decreasing number of species in the dataset (and less  
636 strongly also with decreasing number of samples). This makes these tests less suitable for  
637 smaller and relatively homogeneous datasets with few species (e.g. less than 40), since the  
638 probability of Type II error (i.e. not rejecting the null hypothesis which is false) strongly  
639 increases. Similarly, both two-step and modified permutation tests are overly conservative for  
640 scenario 3. For modified permutation test this is not a problem, since for the hypotheses for  
641 which the scenario 3 is null hypothesis (category 2) the modified permutation test is not  
642 recommended method (see Table 1). The two-step permutation test is also overly conservative,  
643 but only only in case of the most homogeneous community dataset, and with increasing  
644 compositional heterogeneity this issue diminishes (Fig. 6 and Table S2).

645 In this study, I explicitly ignored intraspecific variation in species attributes, focusing  
646 only on use of dataset-wide mean species attribute values. Indeed, intraspecific variation may be  
647 important; e.g. in the context of functional traits, the intra-specific variation gains an increasing  
648 attention (Albert et al. 2012), and relevant question is whether the inclusion of intra-specific  
649 variation (e.g. by including trait values which are sample-specific, not dataset-wide) influences

650 the potential bias reported in this study or not. This question requires further examination, which  
651 goes beyond this study, but in my opinion including another source of variation (species-level  
652 variation in species attributes) does not remove the problem of the bias itself, but makes the  
653 estimation of the bias and its correction more complex.

654 Finally, relevant consideration is whether the weighted-mean approach is actually the best  
655 analytical solution for question which is being explored. In some cases, the question is explicitly  
656 focused on relating community-level values of species attributes, like mean Ellenberg indicator  
657 values (serving as an estimates of ecological conditions for individual sites) or CWM of traits (as  
658 one of the functional-diversity metrics and as a community-level trait value), and use of  
659 weighted-mean approach is fully justified. Yet, in other cases, when the question is focused on  
660 relating individual species-attributes to sample attributes, weighted-mean approach may not be  
661 the best analytical choice. Use of alternative options, like fourth-corner or RLQ analysis, for  
662 which the problem of inflated Type I error rate and choice of suitable permutation test have been  
663 already solved, can be a better solution.

664

## 665 **Conclusions**

666 In this study, I attempted to draw attention to the problem in weighted-mean approach which I  
667 believe is largely overlooked and generally not acknowledged, although it represents a source of  
668 potentially serious misinterpretations. Since in certain fields the weighted-mean approach gains  
669 increasing momentum (e.g. in functional ecology with CWM of species functional traits as one  
670 of the functional-diversity indices), I suggest that time is ripe to critically asses in which  
671 situations and for which types of hypotheses the commonly used standard parametric or

672 permutation tests are not appropriate, since they yield results which may be overly optimistic. I  
673 offer simple guidelines how to decide whether in given context of a study the standard  
674 methodology gives correct or biased results, and suggest solutions available in case that it does  
675 not.

676

## 677 **Acknowledgements**

678 This study was supported by the Czech Science Foundation (P505/12/1022). My thanks go to  
679 Bill Shipley and Cajo ter Braak for critical comments on the previous versions of this manuscript,  
680 which motivated me to heavily rework it.

681

## 682 **Literature cited**

683 Albert, C. H., F. de Bello, S. Lavorel, and W. Thuiller. 2012. On the importance of intraspecific  
684 variability for the quantification of functional diversity. *Oikos* 121:116-126.

685 Axmanová, I., et al. 2012. Estimation of herbaceous biomass from species composition and  
686 cover. *Applied Vegetation Science* 15:580-589.

687 Birks, H. J. B., J.M. Line, S. Juggins, A. C. Stevenson, and C. J. F. ter Braak. 1990. Diatoms and  
688 pH reconstruction. *Philosophical Transactions of the Royal Society B Biological Sciences*  
689 327:263–278.

690 Bivand, R. 1980. A Monte Carlo study of correlation coefficient estimation with spatially  
691 autocorrelated observations. *Quaestiones Geographicae* 6:5–10.



- 692 Blackburn, T. M., and B. A. Hawkins. 2004. Bergmann's rule and the mammal fauna of northern  
693 North America. *Ecography* 27:715–724.
- 694 Brown, A. M., D. I. Warton, N. R. Andrew, M. Binns, G. Cassis, and G. Helois. 2014. The  
695 fourth-corner solution – using predictive models to understand how species traits interact with  
696 the environment. *Methods in Ecology and Evolution* 5:344–352.
- 697 Clavero, M., and L. Brotons. 2010. Functional homogenization of bird communities along  
698 habitat gradients: accounting for niche multidimensionality. *Global Ecology and Biogeography*  
699 19:684–696.
- 700 Deblauwe, V., P. Kennel, and P. Coueron. 2012. Testing pairwise association between spatially  
701 autocorrelated variables: a new approach using surrogate lattice data. *Plos One* 7:e48766.
- 702 Díaz, S., M. Cabido, and F. Casanoves. 1998. Plant functional traits and environmental filters at  
703 a regional scale. *Journal of Vegetation Science* 9:113–122.
- 704 Díaz, S., S. Lavorel, F. de Bello, F. Quétler, K. Grigulis, and T. M. Robson. 2007. Incorporating  
705 plant functional diversity effects in ecosystem service assessments. *Proceedings of the National*  
706 *Academy of Sciences USA* 104:20684–20689.
- 707 Dray, S., and P. Legendre. 2008. Testing the species traits-environment relationships: the fourth-  
708 corner problem revisited. *Ecology* 89 3400–3412.
- 709 Dolédec, S., D. Chessel, C. J. F. ter Braak, and S. Champely. 1996. Matching species traits to  
710 environmental variables: a new three-table ordination method. *Environmental and Ecological*  
711 *Statistics* 3:143–166.

- 712 Dutilleul, P. 1993. Modifying the t test for assessing the correlation between two spatial  
713 processes. *Biometrics* 49:305–314.
- 714 Ellenberg, H., H. E. Weber, R. Düll, V. Wirth, W. Werner, and D. Paulissen. 1992. *Zeigerwerte*  
715 *von Pflanzen in Mitteleuropa*. Second Edition. *Scripta Geobotanica* 18:1–248.
- 716 Fajmonová, Z., D. Zelený, V. Syrovátka, G. Vončina, and M. Hájek. 2013. Distribution of  
717 habitat specialists in semi-natural grasslands. *Journal of Vegetation Science* 24:616–627.
- 718 Fortin, M.-J., and M. R. T. Dale. 2005. *Spatial Analysis. A Guide for Ecologists*. Cambridge  
719 University Press, New York, USA
- 720 Fox, J. W. 2012. When should we expect microbial phenotypic traits to predict microbial  
721 abundances? *Frontiers in Microbiology* 3:268.
- 722 Fridley, J. D., D. B. Vandermast, D. M. Kuppinger, M. Manthey, R. K. Peet. 2007. Co-  
723 occurrence based assessment of habitat generalists and specialists: a new approach for the  
724 measurement of niche width. *Journal of Ecology* 95:707–722.
- 725 Garnier, E., et al. 2004. Plant functional markers capture ecosystem properties during secondary  
726 succession. *Ecology* 85:2630–2637.
- 727 Gelman, A. 2005. Analysis of variance – why is it more important than ever. *The Annals of*  
728 *Statistics* 33:1–53.
- 729 Hawkins, B. A., and J. A. F. Diniz-Filho. 2006. Beyond Rapoport’s rule: evaluating range size  
730 patterns of New World birds in a two-dimensional framework. *Global Ecology and*  
731 *Biogeography* 15:461–469.

- 732 Hawkins, B. A., M. Rueda, T. F. Rangel, R. Field, and J. A. F. Diniz-Filho. 2014. Community  
733 phylogenetics at the biogeographic scale: cold tolerance, niche conservatism and the structure of  
734 North American forests. *Journal of Biogeography* 41:23–28.
- 735 Jamil, T., W. A. Ozinga, M. Kleyer, and C. J. F. ter Braak. 2013. Selecting traits that explain  
736 species-environment relationships: a generalized linear mixed model approach. *Journal of*  
737 *Vegetation Science* 24:988–1000.
- 738 Juggins, S. 2013. Quantitative reconstructions in palaeolimnology: new paradigm or sick science?  
739 *Quaternary Science Reviews* 64:20–32.
- 740 Kelly, M. G., and B. A. Whitton. 1995. Biological monitoring of eutrophication in rivers.  
741 *Hydrobiologia* 384:55–67.
- 742 Kleyer, M., S. Dray, F. de Bello, J. Lepš, R. J. Pakeman, B. Strauss, W. Thuiller, and S. Lavorel.  
743 2012. Assessing species and community functional responses to environmental gradients: which  
744 multivariate methods? *Journal of Vegetation Science* 23:805–821.
- 745 Lambers, H., F. S. Chapin III, and T. L. Pons. 2008. *Plant Physiological Ecology*, 2nd Edition.  
746 Springer, New York, USA.
- 747 Laliberté, E, B. Shipley, D. A. Norton, and D. Scott. 2012. Which plant traits determine  
748 abundance under long term shifts in soil resource availability and grazing intensity? *Journal of*  
749 *Ecology* 100:662–677.
- 750 Landolt, E. 1977. *Ökologische Zeigerwerte zur Schweizer Flora*. Veröffentlichungen des  
751 Geobotanischen Institutes der Eidgenössischen Technischen Hochschule, Stiftung  
752 Rübél, Zürich, 64:1–208.

- 753 Legendre, P. 1993. Spatial autocorrelation: trouble or new paradigm? *Ecology* 74:1659–1673.
- 754 Legendre, P., R. Galzin, and M. L. Harmelin-Vivien. 1997. Relating behavior to habitat:  
755 solutions to the fourth-corner problem. *Ecology* 78:547–562.
- 756 Legendre, P, and L. Legendre. 2012. *Numerical Ecology*, Third English Edition. Elsevier  
757 Science, Amsterdam, The Netherlands.
- 758 Legendre, P., and M. De Cáceres. 2012. Beta diversity as the variance of community data:  
759 dissimilarity coefficients and partitioning. *Ecology Letters* 16:951–963.
- 760 Lengyel, A., M. Chytrý, and L. Tichý. 2001. Heterogeneity-constrained random resampling of  
761 phytosociological databases. *Journal of Vegetation Science* 22:175–183.
- 762 Minchin, P. R. 1987. Simulation of multidimensional community patterns: towards a  
763 comprehensive model. *Vegetatio* 71:145–156.
- 764 Mlambo, M. C. 2014. Not all traits are ‘functional’: insights from taxonomy and biodiversity-  
765 ecosystem functioning research. *Biodiversity and Conservation* 23:781–790.
- 766 Peres-Neto, P. R., M. A. Leibold, and S. Dray. 2012. Assessing the effects of spatial contingency  
767 and environmental filtering on metacommunity phylogenetics. *Ecology* 93:S14–S30.
- 768 Reich, P. B., M. B. Walters, and D. S. Ellsworth. 1992. Leaf life-span in relation to leaf, plant  
769 and stand characteristics among diverse ecosystems. *Ecological Monographs* 62:365–392.
- 770 Schaffers, A. P., and K. V. Sýkora. 2000. Reliability of Ellenberg indicator values for moisture,  
771 nitrogen and soil reaction: comparison with field measurements. *Journal of Vegetation Science*  
772 11:225–244.

- 773 Shipley, B. 2010. From Plant Traits to Vegetation Structure. Chance and Selection in the  
774 Assembly of Ecological Communities. Cambridge University Press, Cambridge, UK.
- 775 Sládeček, V. 1973. System of water quality from the biological point of view. Archiv für  
776 Hydrobiologie 7:1–218.
- 777 Šmilauer, P, and J. Lepš. 2014. Multivariate analysis of ecological data using CANOCO 5.  
778 Second Edition. Cambridge University Press, Cambridge, UK.
- 779 ter Braak, C. J. F., and L. G. Barendregt. 1986. Weighted averaging of species indicator values:  
780 its efficiency in environmental calibration. Mathematical Biosciences 78:57–72.
- 781 ter Braak, C. J. F., and C. W. N. Looman. 1986. Weighted averaging, logistic regression and the  
782 Gaussian response model. Vegetatio 65:3–11.
- 783 ter Braak, C. J. F., A. Cormont, and S. Dray. 2012. Improved testing of species traits-  
784 environment relationships in the fourth-corner problem. Ecology 93:1525–1526.
- 785 Vile, D., B. Shipley, and E. Garnier. 2006. Ecosystem productivity can be predicted from  
786 potential relative growth rate and species abundance. Ecology Letters 9:1061–1067.
- 787 Wamelink, G. W. W., V. Joosten, H. F. van Dobben, and F. Berendse. 2002. Validity of  
788 Ellenberg indicator values judged from physico-chemical field measurements. Journal of  
789 Vegetation Science 13:269–278.
- 790 Wamelink, G. W. W., P. W. Goedhart, H. F. van Dobben, and F. Berendse. 2005. Plant species  
791 as predictors of soil pH: Replacing expert judgment with measurements. Journal of Vegetation  
792 Science 16:461–470.

793 Zelený, D., and A. P. Schaffers. 2012. Too good to be true: pitfalls of using mean Ellenberg  
794 indicator values in vegetation analyses. *Journal of Vegetation Science* 23:419–431.

795

796 **Supplementary materials**

797 **Appendix S1.** Description of an algorithm generating artificial community data along two  
798 environmental gradients.

799 **Appendix S2.** Weighted-mean approach applied on simulated data: additional results.

800 **Appendix S3.** Dissimilarity index between two CWM values and beta diversity assessment.

801 **Appendix S4.** Evaluation of permutation tests using simulated data from Dray & Legendre  
802 (2008).

803 **Appendix S5.** R-code for all analyses.

804 **Appendix S6.** Source code for the R library *weimea*, version v. 0.58 (actual version can be found  
805 on <https://github.com/zdealveindy/weimea/>).

806

807 *Table 1*

808 Overview of the characteristics for the three categories of hypotheses tested by *weighted-mean*  
 809 approach. For each situation, corresponding assumption about link between sample attributes (**R**)  
 810 or species attributes (**Q**) to species composition (**L**) is given, as well as null vs alternative  
 811 hypothesis, scenario relevant in the context of given category (see Fig. 2), and recommended test.

Category	Assumption	Null hypothesis	Alternative hypothesis	Relevant scenario	Recommended test
1	sample attributes fixed ( <b>R</b> <-> <b>L</b> )	<b>Q</b> <-> <b>L</b>	<b>Q</b> <-> <b>L</b>	Scenario 2	modified permutation test
2	species attributes fixed ( <b>Q</b> <-> <b>L</b> )	<b>R</b> <-> <b>L</b>	<b>R</b> <-> <b>L</b>	Scenario 3	standard parametric or permutation test
3	no assumptions	<b>R</b> <-> <b>Q</b> , i.e. <b>R</b> <-> <b>L</b> and/or <b>Q</b> <-> <b>L</b>	<b>R</b> <-> <b>Q</b> , i.e. <b>R</b> <-> <b>L</b> and <b>Q</b> <-> <b>L</b>	Scenarios 2, 3 and 4	two-step permutation test

812

813

814 **Figure captions**

815 **Figure 1.** Computational schema of the weighted-mean approach to analyse relationship between  
816 species attributes and sample attributes via matrix of species composition. **R** - matrix of sample  
817 attributes (e.g. environmental variables), **L** - matrix of species composition (**L<sub>s</sub>** – **L** standardized  
818 by sample totals to simplify the equation), **Q** - matrix of species attributes (e.g. traits, species  
819 indicator values), **M** - matrix of weighted means of species attributes (e.g. CWM). The colour  
820 gradient within the matrix **M** (weighted mean of species attributes) from dark to light grey  
821 illustrates that this matrix includes information from both matrix of species composition (dark  
822 grey) and matrix of species attributes (light grey).

823 **Figure 2.** Schema showing conceptual differences between scenarios 1-4 in weighted-mean  
824 approach. In scenario 1, both sample attributes (**R**) and species attributes (**Q**) are fixed, linked to  
825 matrix of species composition (**L**), while in the other three scenarios one (or both) of attributes  
826 are considered random, without the link to species composition. In simulated data example, the  
827 link of attributes to species composition is cancelled by permuting the values of species attributes  
828 (scenario 2), sample attributes (scenario 3) or both (scenario 4). In the schema, matrix of species  
829 attributes is transposed (**Q'**) to match the dimension of matrix of species composition (**L**).

830 **Figure 3.** Pearson's *r* correlation coefficients among CWM and sample attributes for each of the  
831 four scenarios and ten levels of compositional heterogeneity of species matrix (1000 correlations  
832 for each combination have been conducted). Grey horizontal bars are outliers.

833 **Figure 4.** Proportion of significant correlations ( $P < 0.05$ ) between CWM and sample attributes,  
834 tested by standard parametric *t*-test. For each of the four scenarios and ten levels of  
835 compositional heterogeneity of species matrix, 1000 tests have been conducted.



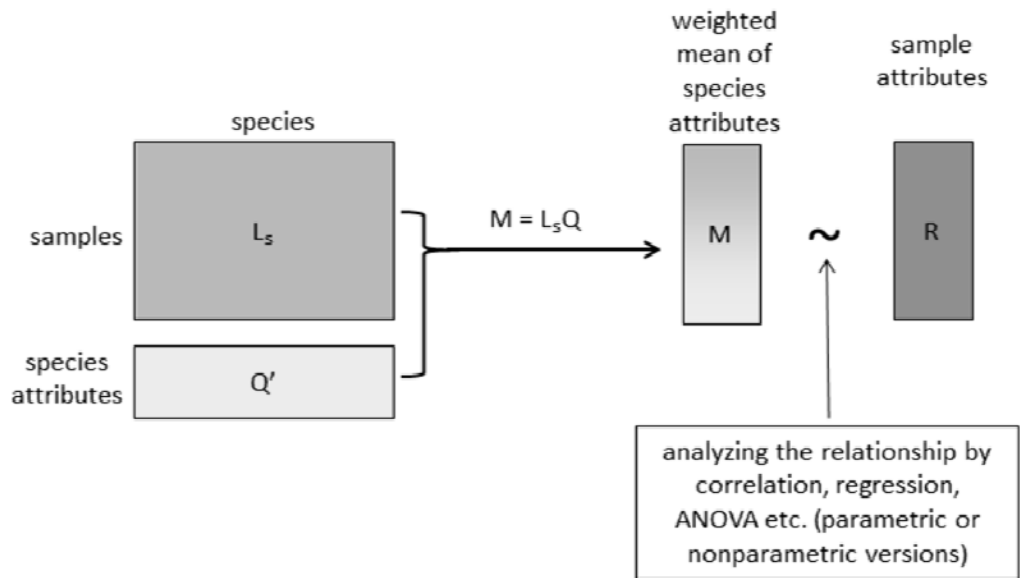
836 **Figure 5.** Simplified example of two community datasets (each with six samples) collected in  
837 plots of two different environmental conditions (A and B, e.g. open- vs closed-canopy forest).  
838 First sampling design (a) restricts choice into only two vegetation types (one open- and one  
839 closed-canopy forest), and results into three and three plots with identical species composition.  
840 Second sampling design (b) does not restrict the sampling by choice of community type (any  
841 forest with open- or closed-canopy can be sampled), resulting in situation when none of six  
842 samples share any species in common. For each dataset, three matrices are presented:  
843 sample  $\times$  species compositional matrix, matrix of sample attributes (in this case with two-level  
844 categorical variable) and matrix of species attributes (quantitative variable in range 1 to 5). x -  
845 presence of species in the sample.

846 **Figure 6.** Proportion of significant correlations ( $P < 0.05$ ) between CWM and sample attributes,  
847 tested by three different permutation tests: standard, modified and two-step. For each of the four  
848 scenarios and ten levels of compositional heterogeneity of species matrix, 1000 tests have been  
849 conducted.

850

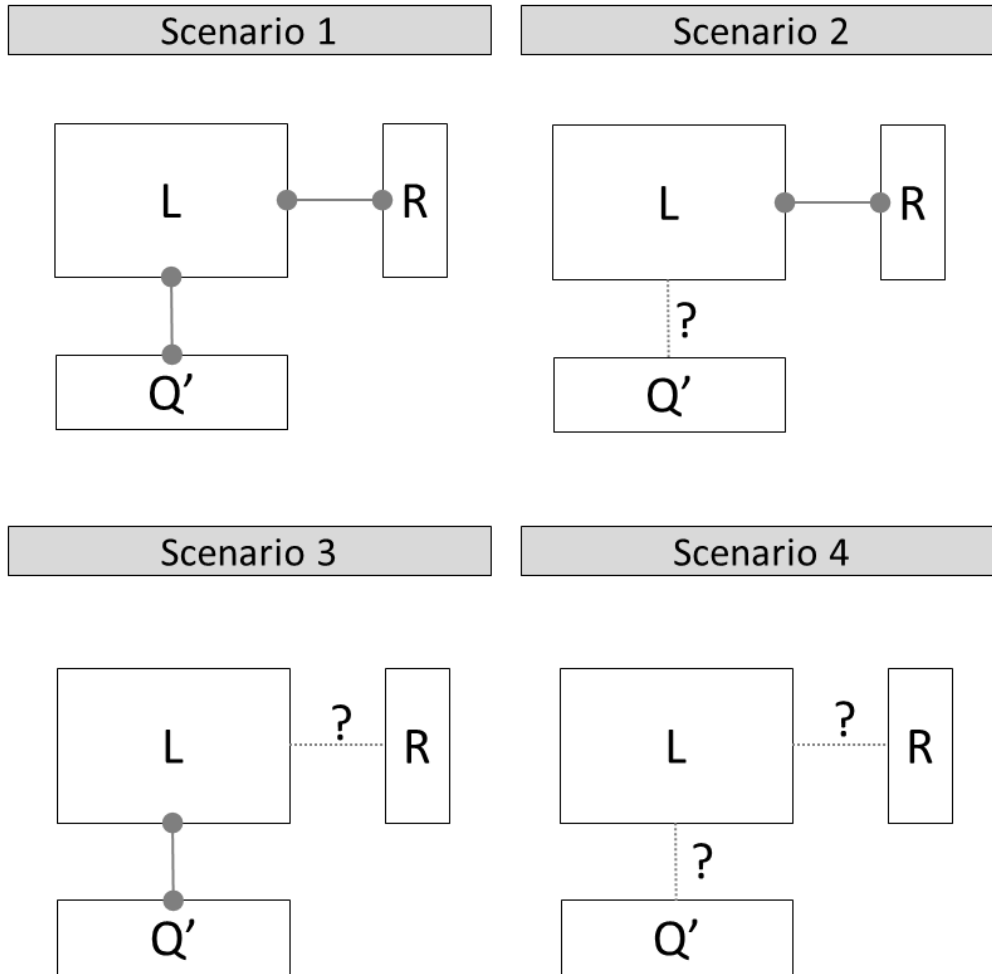
851 *Figure 1*

### *Weighted-mean approach*



852

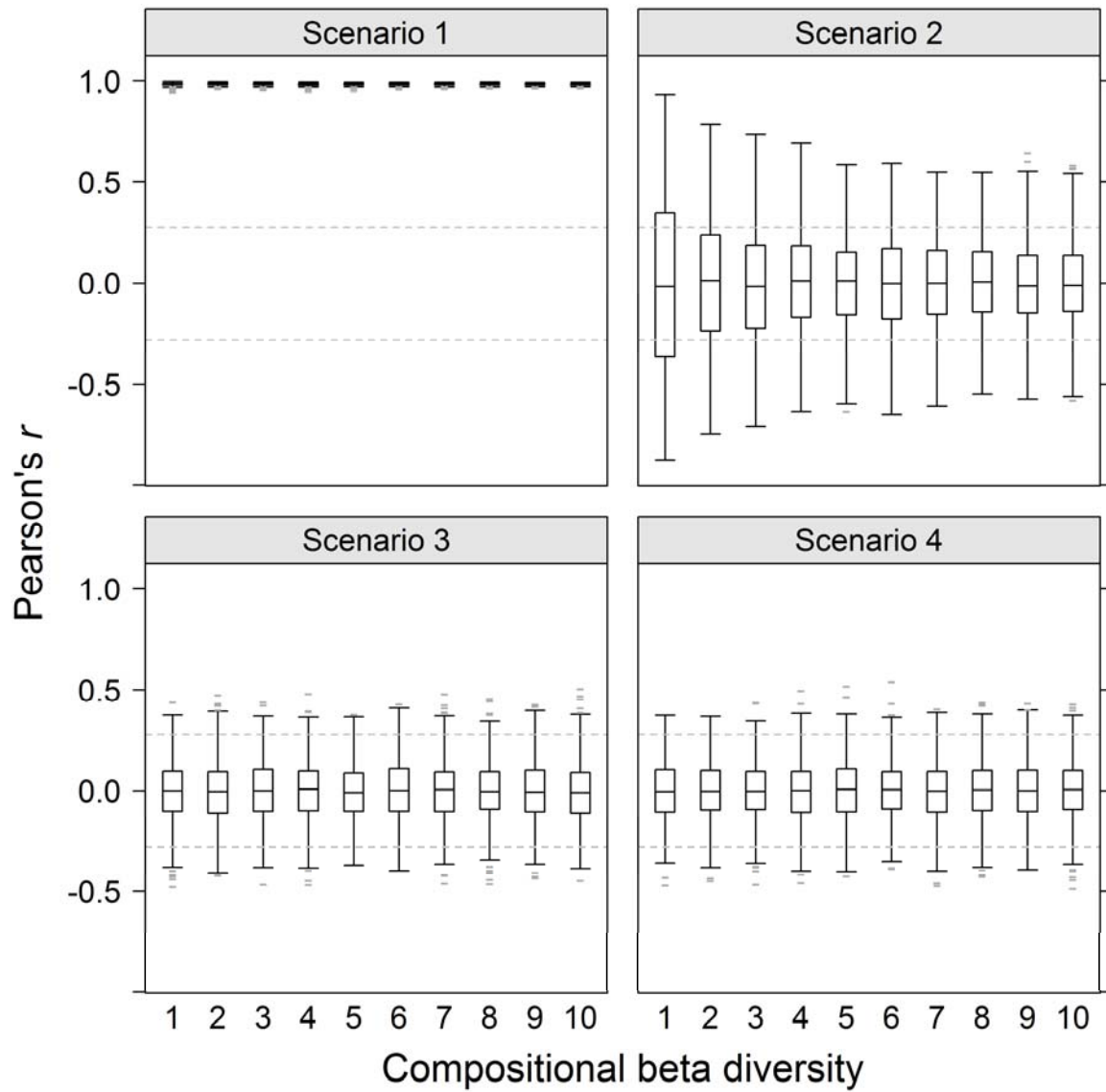
853 *Figure 2*



854

855

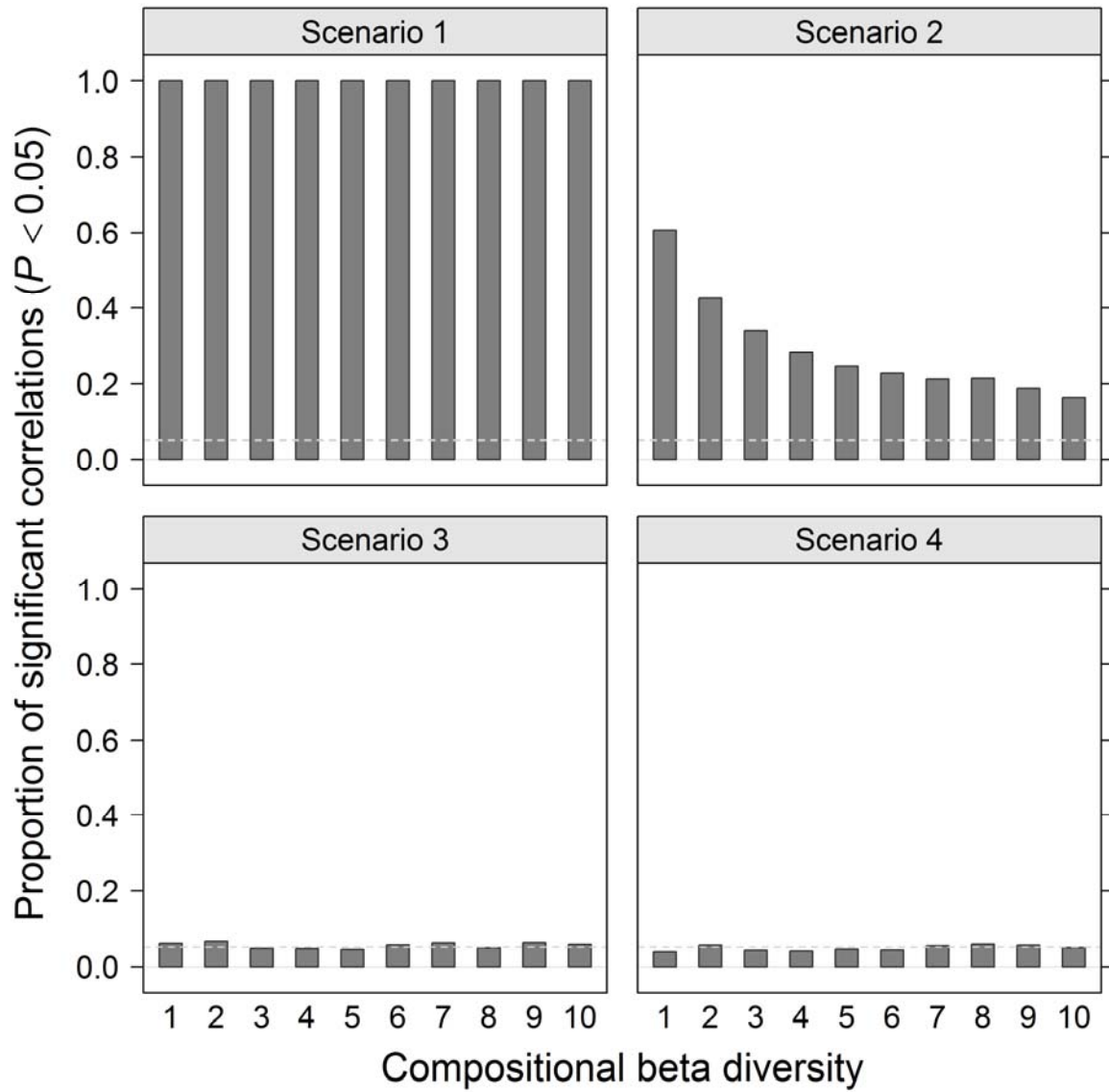
856 *Figure 3*



857

858

859 *Figure 4*



860

861

862 *Figure 5*

(a)

```
          species
          11111
12345678901234 env

sample_1 xxxxxxxx      A
sample_2 xxxxxxxx      A
sample_3 xxxxxxxx      A
sample_4          xxxxxxxx B
sample_5          xxxxxxxx B
sample_6          xxxxxxxx B

spec.attr. 11122233344455
```

(b)

```
          species
          1111111112222222222333333333444
123456789012345678901234567890123456789012 env

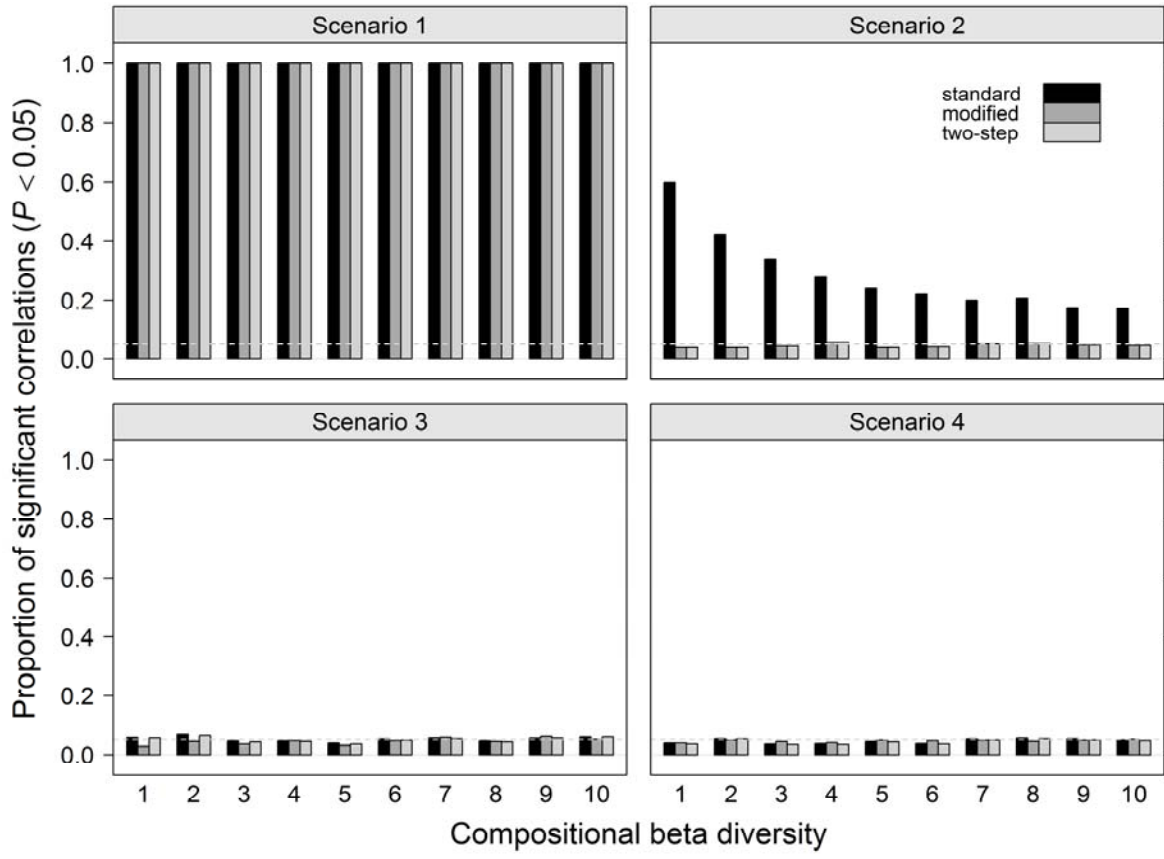
sample_1 xxxxxxxx      A
sample_2          xxxxxxxx      A
sample_3          xxxxxxxx      A
sample_4          xxxxxxxx      B
sample_5          xxxxxxxx      B
sample_6          xxxxxxxx      B

spec.attr. 111222311122231112223334445533444553344455
```

863

864

865 *Figure 6*



866