

1 Ecological Informatics 44: 1-6. (2018)

2

3 **Through the jungle of methods quantifying multiple-site resemblance**

4

5

6 Dénes Schmera^{1,2} & János Podani^{3,4}

7

8

9 ¹MTA Centre for Ecological Research, Balaton Limnological Institute, Klebelsberg K. u. 3, H-
10 Tihany, Hungary, E-mail: schmera.denes@okologia.mta.hu

11 ²MTA Centre for Ecological Research, GINOP Sustainable Ecosystem Group, Klebelsberg K. u.
12 3, H-Tihany, Hungary

13 ³Department of Plant Systematics, Ecology and Theoretical Biology, Institute of Biology, L.
14 Eötvös University, Budapest, Hungary

15 ⁴MTA-ELTE-MTM Ecology Research Group, Budapest, Hungary

16

17

18 **Abstract**

19 Methods that quantify multiple-site resemblance are basic toolkits of ecology for studying
20 community variation in space and time. Although both pairwise and multiple-site
21 coefficients have received increasing attention in the past decade, the high variety of
22 methodologies combined with the absence of a systematic review prevents full
23 understanding and comprehension. To illuminate the situation, we compare and classify
24 methods that use incidence data and propose a unified terminology. The methods can be
25 grouped according to families, approaches and forms. The examination of algebraic
26 expressions and analyses of artificial and actual data sets suggest that inference drawn
27 about communities strongly depends on the methodology applied. We found that the
28 impact of mimicking the original pairwise indices (i.e. the impact of families) was stronger
29 than the impact of components used in formulating the coefficients (i.e. the impact of
30 approach). Our findings suggest that the measures examined quantify drastically different

31 facets of multiple-site resemblance and therefore they have to be selected with care in
32 community studies.

33

34 **Keywords**

35 community resemblance, community variation, dissimilarity, multiple-site resemblance
36 coefficients, pairwise resemblance coefficients, similarity

37

38

39 **1. Introduction**

40

41 Understanding spatial variation in species composition is one of the most fundamental
42 challenges of community ecology. This is promoted by testing hypotheses about the
43 processes that generate and maintain biodiversity in ecosystems (Legendre & De Cáceres,
44 2013). Invasion ecologists, for instance, examine the impact of alien species on native
45 communities, while conservation biologists rely on the measurement of compositional
46 variation in prioritizing areas. The spatial variation of communities can be viewed as either
47 compositional differentiation or similarity (Jost et al., 2011). Beta diversity (Whittaker, 1960,
48 1972), for instance, expresses compositional differentiation, while community overlap (Arita,
49 2017, Schmera, 2017) relates to compositional similarity – which are two sides of the same
50 coin.

51

52 Community variation has been traditionally studied by examining several pairs of sites from
53 the same locality (but see Legendre & De Cáceres, 2013, for alternative solutions) and
54 quantified by the average value of pairwise resemblance (i.e., similarity or dissimilarity)
55 coefficients (Koleff et al., 2003). Such averages may be used to express both compositional
56 similarity and differentiation. Recently, however, it has been suggested that inference drawn
57 from mean values may be misleading, because pairwise resemblance coefficients cannot
58 account properly for co-occurrence patterns of species in many sites and therefore special
59 indices are required (Diserud & Ødegaard, 2007; Baselga, 2013).

60

61 Although multiple-site resemblance coefficients have received increasing attention in
62 contemporary ecology, our knowledge on their relative merits and potential disadvantages is
63 still limited. A recent review on beta diversity deliberately omitted their discussion (Legendre
64 & De Cáceres, 2013) while an even more recent study deepened our understanding of
65 multiple-site overlap measures by providing novel measures and a unified terminology
66 (Arita, 2017). Unfortunately, however, the increasing number of methods, the application of
67 different and often overcomplicated mathematical equations, the ambiguous terminology,
68 as well as the parallel development of similarity and dissimilarity forms impede proper
69 measurement of multiple-site resemblance. Therefore, for the benefit of practicing
70 ecologists, we review the methods quantifying multiple-site resemblance that are based on
71 incidence (presence-absence) data. First, we discuss some basic terms, then we overview
72 pairwise and multiple-site resemblance coefficients. Specifically, we identify and match
73 similarity and dissimilarity forms and simplify some equations. Finally, by using artificial and
74 actual data sets we compare the performance of multiple-site resemblance measures.

75

76

77 **2. Basic terms**

78

79 Originally, pairwise and multiple-site resemblance coefficients have been suggested to
80 measure the (dis)similarity of two or multiple sites based on the presence-absence of
81 species. Consequently, sites are the *objects* of such studies and species are the *descriptors*
82 which characterize the objects. Observed data are commonly arranged in matrix $\mathbf{X} \equiv \{x_{ij}\}$, in
83 which rows represent sites while columns correspond to species (e.g. Legendre & DeCáceres,
84 2013), a convention followed here as well. *Occurrence* (of species j in site i) means that
85 species j is present in site i , coded as $x_{ij} = 1$. In case of species absence, $x_{ij} = 0$. The *species*

86 *richness* of site i (t_i) is the number of occurrences in the given row (row total, $t_i = \sum_{j=1}^T x_{ij}$,

87 where T is the number of species). The *occurrence frequency* of species j (n_j) is the number of
88 sites in which the species is present (called also as range size and calculated as the column

89 total, $n_j = \sum_{i=1}^N x_{ij}$, where N is the number of sites). Whereas *co-occurrence* is traditionally

90 understood as the presence of a pair of species in a given site (Mackenzie et al. 2004, Bell
91 2005, Pollock et al. 2014 and references therein), Arita and co-workers (Trejo-Barocio & Arita
92 2013, Arita 2017) termed *co-diversity*, with a reference to Bell (2005), as the occurrence of a
93 species in two sites. It follows that the *number of co-occurrences* in a site is the number of
94 species pairs present there, while the *number of co-diversities* is the number of unique site-
95 pair occupancies of a given species. In a more formal way, the number of co-occurrences in
96 site i can be expressed as

$$97 \binom{t_i}{2}, \quad \text{Eq. 1}$$

98 while the number of co-diversities of species j as:

$$99 \binom{n_j}{2}. \quad \text{Eq. 2.}$$

100 Furthermore, following Schmera (2017) we consider *community overlap* as a phenomenon
101 that represents the intersection in the composition of sites, *overlapping species* as species
102 with at least two occurrences in a set of sites, *overlap size* as a quantitative property of
103 overlapping species that is quantified as the occurrence frequency of the given species
104 minus one:

$$105 n_j - 1, \quad \text{Eq. 3.}$$

106 and *total overlap size* as a quantitative property of community overlap

$$107 \sum_{j=1}^T n_j - T. \quad \text{Eq. 4.}$$

108

109

110 **3. Pairwise resemblance coefficients: a short overview**

111

112 The literature of numerical ecology abounds in resemblance coefficients (sensu Orlóci 1972)
113 for comparing pairs of sites based on their species composition. We are concerned here with

114 similarity (s) and dissimilarity (d) forms which are bounded between 0 and 1, and are
115 therefore complements ($d + s = 1$). Presence-absence versions are commonly expressed in
116 terms of a 2 x 2 contingency table in which a refers to the number of species present in both
117 sites being compared (shared species, or the number of overlaps in species composition), b
118 to the number of species present only in the first and c to the number of species in the
119 second. That is, with respect to a given pair of sites there are b and c species unique to the
120 first and to the second site, respectively, so that the total number of species in the two sites
121 equals to $a + b + c$. We shall focus on three well-known resemblance coefficients, namely the
122 Jaccard, the Simpson and the Sørensen indices (Table 1, see Koleff et al. 2003 for further
123 indices).

124

125

126 **4. A proposal for a unified terminology to classify methods quantifying multiple-site** 127 **resemblance**

128

129 Here we suggest a unified terminology that allows the classification of methods quantifying
130 multiple-site resemblance. The multiple-site indices (see next paragraph for details)
131 mimicking some properties of the original pairwise Jaccard, Simpson and Sørensen indices
132 (Table 2) are termed as different groups (Legendre 2014), types (Arita 2017) or families
133 (Baselga, 2012, Baselga & Leprieur, 2015, Podani & Schmera, 2016). This confused
134 nomenclature, however, does not support the development of the field. We therefore
135 suggest, following the terminology of the first classifier (Baselga 2012), that classes of
136 methods mimicking some properties of the original pairwise coefficients should be termed
137 as *families*. Accordingly, the methods in question can be classified into Jaccard, Simpson and
138 the Sørensen families.

139

140 Families, however, do not provide the only way for classifying multiple-site resemblance
141 measures. The next feature on which further grouping is made depends on the type of
142 *components* (mathematical terms) incorporated into the coefficient. Some of the measures

143 rely only upon pairwise components, some use only general components of the studied
144 presence-absence matrix such as the total overlap size, others use co-diversity and, finally,
145 further ones combine general and pairwise components (see next paragraph for details). We
146 suggest that classes of methods formed according to the components used should be
147 termed as *approaches*, and we can distinguish among mean pairwise, general, co-diversity
148 and mixed components approaches (see below). Finally, as said above, each coefficient can
149 be expressed as similarity or dissimilarity. We will refer to this property of coefficients as
150 *forms*.

151

152 Consequently, we suggest a classification of methods quantifying multiple-site resemblance
153 according to families, approaches and forms. The terminology becomes even more complex
154 if we consider that dissimilarity forms (also used as measures of beta diversity) may be
155 partitioned into additive components to separate the effect of various background factors
156 influencing dissimilarity. There are two different *frameworks* for such a partitioning,
157 intensively discussed and debated in the relevant literature (Baselga, 2010, Carvalho et al.,
158 2013, Cardoso et al., 2014, Ensing & Pither, 2015, Chen, 2016, Podani & Schmera, 2016).

159

160

161 **5. Multiple-site resemblance: a new classification**

162

163 Here we suggest a classification of methods assessing multiple-site resemblance by
164 considering families, approaches and forms. In this, we do not suggest any hierarchy among
165 these categories. The classification includes both pairwise and multiple-site coefficients.
166 Pairwise coefficients are used for quantifying multiple-site resemblance by calculating the
167 mean of pairwise coefficients, referred here as "mean pairwise approach".

168

169 Multiple-site coefficients express resemblance of more than two sites simultaneously (Table
170 2). Although the first multiple-site index dates back to the 1950's (Koch 1957, see also Eq.

171 Tab2/1), further elaboration of such coefficients has started only recently. Some of the new
172 indices follow the logic of pairwise indices and therefore we categorize them into the
173 Jaccard, Simpson and Sørensen families (Table 2). However, the coefficients in either family
174 use different components in quantifying similarity or dissimilarity. In studying the overlap of
175 multiple sites, for instance, Arita (2017) suggested general overlap indices, which use only
176 some general components of the incidence matrix, as well as co-diversity indices, which use
177 the occurrence of two species at a particular site. In other studies, Baselga and co-workers
178 (Baselga et al. 2007, Baselga 2010, 2012) used both general and pairwise components in
179 expressing multiple site resemblance or, in other words, they used mixed components. As
180 said above, we refer to this property of coefficients as *approach* and distinguish among
181 mean pairwise (Table 1), general, co-diversity and mixed components approaches (Table 2).
182 Note that we use the term *general* instead of *general overlap* (sensu Arita 2017), because
183 “general” can reflect both similarity and dissimilarity, whereas “general overlap” intuitively
184 relates to similarity only. Thus, we can distinguish three families (Jaccard, Simpson and
185 Sørensen), four approaches (mean pairwise, general, co-diversity and mixed components)
186 and two forms (similarity and dissimilarity) of methods quantifying multiple site resemblance
187 (Tables 1 & 2).

188

189 General similarity indices belonging to different families may be formalized in different ways
190 (Table 2). The observed total overlap size (Eq. 4) may be divided by the maximum number of
191 total overlap size with N sites and T species (Jaccard family, Eq. Tab2/1), or by the maximum
192 number of total overlap sizes possible if the sites show a nested design (Simpson family, Eq.
193 Tab2/3). Thirdly, average overlap size of species may be divided by the average species
194 richness of sites (Sørensen family, Eq. Tab2/4).

195

196 Baselga and co-workers (Baselga et al. 2007, Baselga 2010, 2012), following Diserud &
197 Ødegaard (2007), used $\sum_i t_i - T$ as the "number of shared species" in the multiple-site

198 situation. Since $\sum_{i=1}^N t_i = \sum_{j=1}^T n_j = G$, the grand total of \mathbf{X} (see also Arita et al., 2008, 2012; Arita

2017), we can call $\sum_i t_i - T$ as total overlap size (Eq. 4). Multiple-site "unique species",
however, were quantified as the sum of unique species for pairs of sites. It follows that it is a
mixed components approach having both pairwise and general constituents. A possible
theoretical problem with this is that total overlap size (from the general approach) and the
number of site pairs in which the same species occur (pairwise component, called also as co-
diversity [Arita 2017]) in the data matrix are not the same (Arita 2017), and therefore the
ecological interpretation of these indices is less straightforward.

Moreover, in addition to general indices, Arita (2017) developed a new approach of multiple-
site similarity measures he called the co-diversity indices. These indices, in fact, count the
sum of the two-site occurrences (co-diversity) of species which is divided either by the sum
of the co-diversities when site compositions show a nested design (Simpson family, Eq.
Tab2/10), by the possible number of co-diversities when N sites are occupied with T species
(Jaccard family, Eq. Tab2/9) or, finally, by the sum of average species richness for each pair
of sites (Sørensen family, Eq. Tab2/11).

6. Simplification of some equations

A couple of mixed-component resemblance coefficients have originally been published with
extensive mathematical equations. To make their use easier, we suggest the simplification of
two functions. The mixed component Jaccard dissimilarity (Eq. Tab2/6) suggested by Baselga
(2012) can be simplified to

$$\frac{\sum_{k<l} (b_{kl} + b_{lk})}{(\sum_i t_i - T) + \sum_{k<l} (b_{kl} + b_{lk})}, \quad \text{Eq. 5.}$$

while the mixed component Sørensen dissimilarity (Eq. Tab2/8) suggested by Baselga (2010)
reduces to

225
$$\frac{\sum_{k<l} (b_{kl} + b_{lk})}{2(\sum_i t_i - T) + \sum_{k<l} (b_{kl} + b_{lk})}$$
 Eq. 6.

226

227

228 **7. Comparison of methods quantifying multiple-site similarity**

229

230 *7.1 Methods to compare*

231 Here we compare methods that allow quantification multiple-site resemblance. Although we
 232 use similarity forms, our conclusions are not restricted to similarity because it is
 233 complementary to dissimilarity. Although pairwise coefficients are designed for examining
 234 pairs of sites, the mean values of these coefficients are frequently used for assessing
 235 multiple site similarity. We will refer to this as mean pairwise approach. We examined also
 236 general, mixed components and co-diversity approaches, as well as the Jaccard, Simpson and
 237 Sørensen families. In sum, we define any particular method as the combination of an
 238 approach and a family, and thus compared 12 methods.

239

240 *7.2 Artificial data 1*

241 To compare the performance of methods, we examined all possible communities that can be
 242 produced by the co-occurrence of 4 species in 4 sites. In order to calculate the number of
 243 possibilities, we have to first determine how many ways a single species can be distributed in
 244 N sites. Since there are two outcomes for each site (the species is present or absent), the
 245 possible number of occurrence patterns equals 2^N . However, this includes the situation
 246 when the species is absent from all sites. Therefore, the number of occurrence patterns
 247 reduces to $2^N - 1$ (for $N = 4$ we have 15 different patterns). When we have T species, then the
 248 possible number of co-occurrence patterns increases dramatically $(2^N - 1)^T$ (for $N = T = 4$ we
 249 get 50,625). However, these co-occurrence patterns include empty sites (those without
 250 species) as well. After removing degenerate matrices, the number of meaningful co-
 251 occurrence patterns reduces to 41,503 in the example.

252

253 We calculated multiple-site similarities by the different methods (i.e. the combinations of
254 families and approaches) for each of the 41,503 occurrence patterns. When no similarity
255 form was given (the Jaccard and Sørensen families of mixed components), we used the
256 complement of dissimilarity. The resulting scores served as a data set to calculate the
257 Pearson correlation between different methods, in order to express agreement in trends
258 among the measures. We transformed the correlations to distances (distance = 1 –
259 correlation) and analyzed the distance matrix by UPGMA clustering to obtain a dendrogram.
260 The same distance matrix was analyzed by principal coordinates analysis (PCoA). Thus, in
261 these multivariate studies, each object represents a given measure. We used the *gtools*
262 (Warnes et al. 2014), the *betapart* (Baselga et al. 2013) packages in *R* (R Core Team, 2015)
263 and the SYN-TAX 2000 package (Podani 2001) for computations.

264

265 The dendrogram (Fig. 1) shows that methods belonging to the Simpson family constitute one
266 group, separated from the methods of the Jaccard and Sørensen families grouped in the
267 other. Within the latter, general coefficients are well-separated and grouping is more
268 strongly influenced by the choice of approach than by the family (Fig. 1). The PCoA
269 ordination of the methodologies (Fig. 2) supports these conclusions. The first axis separates
270 the Simpson family from the Jaccard and Sørensen families, while the second separates the
271 general approach from the others. Since these axes account for 44% and 29% of the total
272 variance, respectively, we can conclude that choice between families had stronger impact on
273 the results than another decision between the general overlap approach and the others.

274

275 *7.3 Artificial data 2*

276

277 Artificial data set 1 allowed examining all theoretical possibilities in a matrix with very few
278 sites and species. To obtain a more realistic picture on the relationships among measures,
279 we generated a second artificial data set that is closer to actual community data. We
280 produced 150 sets of 10 sites by 10 species incidence matrices, in which the probability of

281 the occurrence of a species in a particular site was 0.5. We removed degenerate matrices
282 (i.e. those with zero row or column totals) and used the first 100 matrices. We followed the
283 multivariate exploration procedure applied to artificial data set 1. The dendrogram (Fig. 3)
284 shows that methods belonging to the Simpson family form one group, and methods
285 belonging to the Jaccard and Sørensen families appear in another. Within the second group,
286 general coefficients are well-separated. The difference between Figs. 1 and 3 are that Fig. 3
287 shows larger distances among some groups of methods (the maximum distance is larger
288 than 0.5) and at the same time smaller distances among similar methods (the behavior of
289 MC.JAC and MCSOR is similar). The PCoA ordination of the measures (Fig. 4) resulted in
290 much the same conclusions. The separation of the G.SIM from the other methods is clear.
291 On the first axis the Simpson family is distinguished from the Jaccard and Sørensen families,
292 while the second axis separates the general approach from the others. Since these axes
293 account for 68% and 16% of the total variance, respectively, we can conclude that choice
294 between families had stronger impact on the results than the decision between the general
295 overlap approach and the others. We may thus derive the final conclusion from clustering
296 and ordination that the general indices, especially those belonging to the Simpson family,
297 present a rather unique way of calculating multiple-site similarities.

298

299 *7.4 Actual data set*

300

301 Rey (1981) examined the recolonization of islets by arthropods after defaunization by
302 insecticides. The fauna was recorded every week for more than a year; we took the data
303 from the 10th, 13th, 20th and 53rd weeks after treatment. These four data matrices,
304 published in Atmar & Patterson (1995) contain 6 sites and 25, 27, 33 and 33 species,
305 respectively. An analysis equivalent to the mean pairwise Jaccard method indicated a
306 monotonic increase of similarity over the study period (Podani & Schmera 2011).

307

308 We found that all methods compared here indicate a monotonically increasing similarity
309 over time (Fig. 5). Nonetheless, the methods show considerable differences regarding the
310 multiple site similarity in the four assemblages. For instance, in week 53, the Jaccard family

311 co-diversity index (CD.JAC) yields a similarity value of 0.131, while the Simpson family
312 general index (G.SIM) produces 0.698. This suggests that selection of the methodology (i.e.
313 the choice of the family together with the approach) has significant impact on our inference
314 about community pattern (here similarity). It is important to note that the traditionally used
315 mean of pairwise indices (here abbreviated as mean pairwise method) and the other "true"
316 multiple-site indices produced very different results, suggesting that the pairwise and
317 multiple-site measures are complementary.

318

319

320 **8. Conclusions**

321

322 We emphasized that understanding and interpreting the multiple-site community patterns
323 pose relevant methodological issues of contemporary ecology and biogeography. Our review
324 demonstrated that a wide variety of methods have been available for quantifying multiple-
325 site resemblance patterns. To help the ecologist navigating among them, we suggested a
326 classification of methodology. Accordingly, a method is a combination of an index family and
327 an approach. Analyses of simulated and actual data sets revealed that inference drawn on
328 community pattern strongly depends on the applied method: multiple-site incidence
329 coefficients quantify different facets of multiple-site community patterns. In particular, we
330 found that the impact of choosing from original pairwise index families was stronger on
331 quantifying multiple-site resemblance patterns than the impact of selecting different
332 approaches. Thus, any methodology used for studying multiple-site community patterns
333 should be carefully evaluated before use.

334

335

336 **Acknowledgements**

337

338 We thank Hector Arita and Carlo Ricotta for their constructive comments on an earlier draft
339 of the manuscript. This research was supported by GINOP-2.3.2-15-2016-00019 project.

340

341

342 **References**

343

344 Arita, H.T. (2017) Multisite and multispecies measures of overlap, co-occurrence, and co-
345 diversity. *Ecography*, 40: 709-718.

346 Arita, H.T, Christen, J.A., Rodriguez, P. & Soberon, J. (2008) Species diversity and distribution
347 in presence-absence matrices: mathematical relationships and biological implications.
348 *The American Naturalist*, 172, 519-532.

349 Arita, H.T, Christen, J.A., Rodriguez, P. & Soberon, J. (2012) The presence-absence matrix
350 reloaded: the use and interpretation of range-diversity plots. *Global Ecology and*
351 *Biogeography*, 21, 282-292.

352 Atmar, W & Patterson, B.D. (1995) The nestedness temperature calculator: a visual basic
353 program, including 294 presence-absence matrices. AICS Research Inc., Univ. Park, NM
354 and The Field Museum, Chicago, IL

355 Baselga, A. (2010) Partitioning the turnover and nestedness components of beta diversity.
356 *Global Ecology and Biogeography*, 19, 134-143.

357 Baselga, A. (2012) The relationship between species replacement, dissimilarity derived from
358 nestedness and nestedness. *Global Ecology and Biogeography*, 21, 1223-1232.

359 Baselga, A. (2013) Multiple site dissimilarity quantifies compositional heterogeneity among
360 several sites, while average pairwise dissimilarity might be misleading. *Ecography*, 36,
361 124-128.

362 Baselga, A., Jiménez-Valverde, A. & Niccolini, G. (2007) A multiple-site similarity measures
363 independent of richness. *Biology Letters*, 3, 642-645.

364 Baselga, A. & Leprieur, F. (2015) Comparing methods to separate components of beta
365 diversity. *Methods in Ecology and Evolution*, 6, 1069-1079.

366 Baselga, A., Orme, D., Villeger, S., de Bortoli, J & Leprieur, F. (2013) betapart: Partitioning
367 beta diversity into turnover and nestedness components. R package version
368 1.3.<http://CRAN.R-project.org/package=betapart>

369 Bell, G. (2005) The co-distribution of species in relation to the neutral theory of community
370 ecology. *Ecology*, 86, 1757-1770.

371 Cardoso, P., Rigal, F., Carvalho, J.C., Fortelius, M., Borges, P.A.V., Podani, J. & Schmera, D.
372 (2014) Partitioning taxon, phylogenetic and functional beta diversity into replacement
373 and richness difference components. *Journal of Biogeography*, 41, 749-761.

374 Carvalho J.C., Cardoso, P., Borges, P.A.V., Schmera, D. & Podani, J (2013) Measuring fractions
375 of beta diversity and their relationships to nestedness: a theoretical and empirical
376 comparison of novel approaches. *Oikos*, 122, 825-834.

377 Chao, A., Chiu, C.C. & Hsieh, T.C. (2012) Proposing resolution to debates on diversity
378 partitioning. *Ecology*, 93, 2037-2051.

379 Chen, Y. (2016) Partitioning multiple-site tree-like beta diversity into turnover and
380 nestedness components without pairwise comparisons. *Ecological Indicators*, 61, 413-
381 417.

382 Diserud, O.H. & Ødegaard, F. (2007) A multiple-site similarity measure. *Biology Letters*, 3, 20-
383 22.

384 Gotelli, N.J. & Chao, A. (2013) Measuring and estimating species richness, species diversity,
385 and biotic similarity from sampling data. *Encyclopedia of Biodiversity*, vol. 5, pp. 195-
386 211.

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

388 Jost, L., Chao, A. & Chazdon R.R. (2011) Compositional similarity and β (beta) diversity. In
389 Magurran A.E. & McGill B.J. (eds) *Biological diversity. Frontiers in measurement and*
390 *assessment*. Oxford University Press, pp. 66-84.

391 Koch, L.F. (1957) Index of biotic dispersity. *Ecology*, 38, 145-148.

392 Koleff, P., Gaston, K.J. & Lennon, J.L. (2003) Measuring beta diversity for presence-absence
393 data. *Journal of Animal Ecology*, 72, 367-382.

394 Legendre, P. (2014) Interpreting the replacement and richness difference components of
395 beta diversity. *Global Ecology and Biogeography*, 23, 1324-1334.

396 Legendre, P. & De Cáceres M. (2013) Beta diversity as the variance of community data:
397 dissimilarity coefficients and partitioning. *Ecology Letters*, 16, 951-963.

398 MacKenzie D.I., Bailey L.L & Nichols J.D. (2004) Investigating species co-occurrence patterns
399 when species are detected imperfectly. *Journal of Animal Ecology*, 73, 546-555.

400 Orlóci, L. 1972. On objective functions of phytosociological resemblance. *American Midland*
401 *Naturalist*, 88, 28-55.

402 Podani, J. (2001) SYN-TAX 2000. Computer programs for data analysis in ecology and
403 systematics. Scientia, Budapest.

404 Podani, J. & Schmera, D. (2011) A new conceptual and methodological framework for
405 exploring and explaining pattern in presence-absence data. *Oikos*, 120, 1625-1638.

406 Podani, J & Schmera, D. (2016) Once again on the components of pairwise beta diversity.
407 *Ecological Informatics*, 32, 63-68.

408 Pollock L.J., Tingley R., Morris W.K., Golding N., O'Hara R.B., Parris K.M., Vesik P.A. &
409 McCarthy M.A. (2014) Understanding co-occurrence by modelling species
410 simultaneously with a Joint Species Distribution Model (JSDM). *Methods in Ecology*
411 *and Evolution*, 5, 397-406

412 R Core Team (2015). R: A language and environment for statistical computing. R Foundation
413 for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

414 Rey, J.R. (1981) Ecological biogeography of arthropods on *Spartina* islands in northwestern
415 Florida. *Ecological Monographs*, 51, 237-265.

416 Ricotta, C. & Pavoine, S. (2015) A multiple-site dissimilarity measure for species
417 presence/absence data and its relationship with nestedness and turnover. *Ecological*
418 *Indicators*, 54, 203-206.

419 Schmera D (2017) On the operative use of community overlap in analyzing incidence data.
420 *Community Ecology*, 18, 117-119.

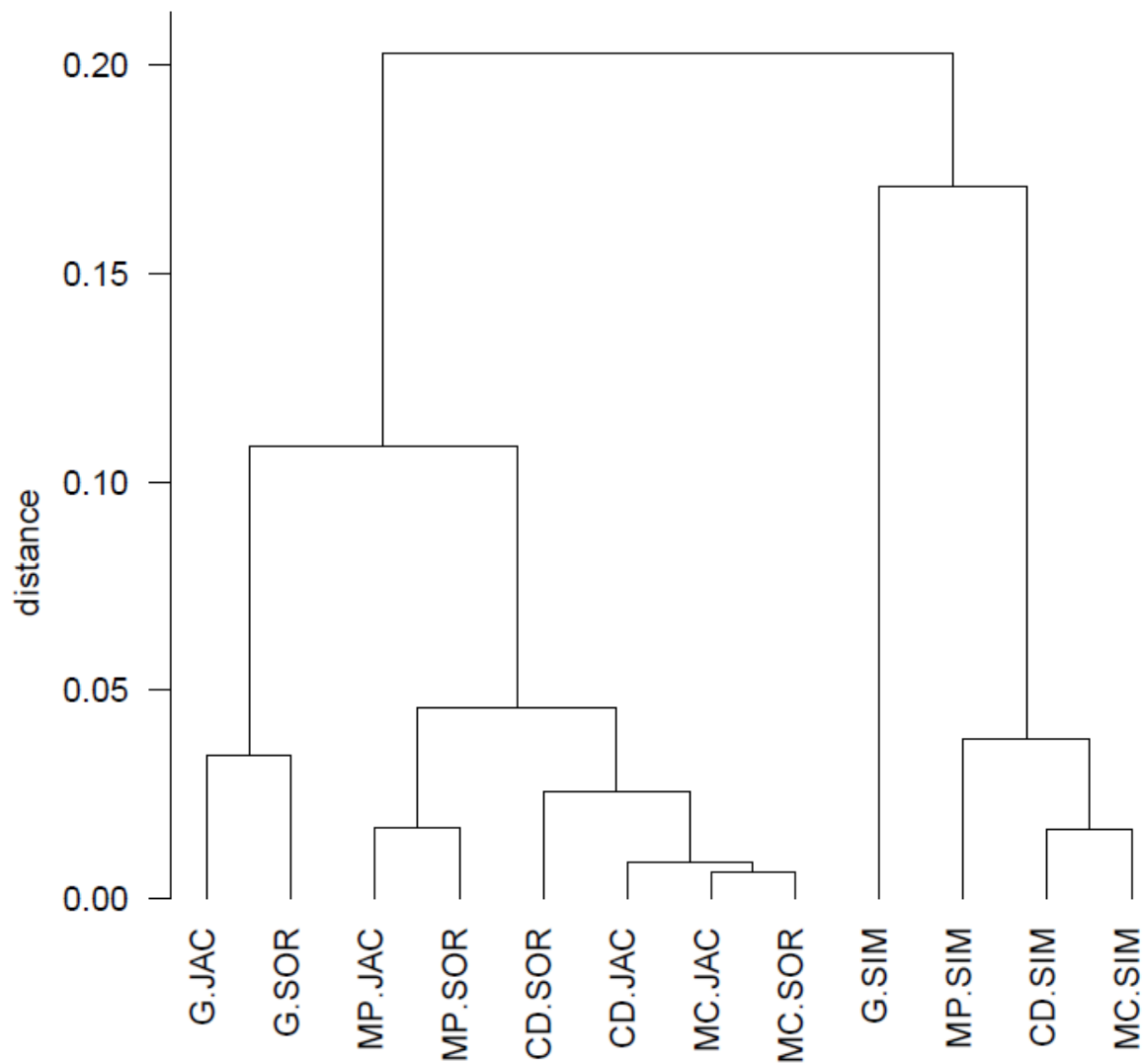
- 421 Simpson, G.G. (1943) Mammals and the nature of continents. American Journal of Science,
422 241, 1-31.
- 423 Sørensen, T.A. (1948) A method of establishing groups of equal amplitude in plant sociology
424 based on similarity of species content, and its application to analyses of the vegetation
425 on Danish commons. Kongelige Danske Viden- skabernes Selskabs Biologiske Skrifter,
426 5, 1–34.
- 427 Trejo-Barocio, P. & Arita, H.T. (2013) The co-occurrence of species and the co-diversity of
428 sites in neutral models of biodiversity. Plos One, 8, e79918.
- 429 Warnes GR, Bolker B, Lumley T (2014) gtools: Various R programming tools. R package
430 version 3.4.1. <http://CRAN.R-project.org/package=gtools>
- 431 Whittaker, R.H. (1960) Vegetation of the Siskiyou Mountains, Oregon and California.
432 Ecological Monographs, 30, 279-338.
- 433 Whittaker, R.H. (1972) Evolution and measurement of species diversity. Taxon, 21, 213-251.
434

436 Table 1. The most important properties of three well known pairwise resemblance coefficients

Family	Form	Equation	Eq. n.	Interpretation	Reference
Jaccard	similarity	$\frac{a}{a + b + c}$	Tab1/1	the ratio of the number of shared species to the total number of species	Jaccard (1912)
	dissimilarity	$\frac{b + c}{a + b + c}$	Tab1/2	the ratio of the number of unique species to the total number of species	
Simpson	similarity	$\frac{a}{a + \min(b, c)}$	Tab1/3	the ratio of the number of shared species to the number of species at the poorest site	Simpson (1943)
	dissimilarity	$\frac{\min(b, c)}{a + \min(b, c)}$	Tab1/4	the ratio of the number of species unique to the poorest site and the number of species in the poorest site	
Sørensen	similarity	$\frac{2a}{2a + b + c} = \frac{a}{\frac{1}{2}((a + b) + (a + c))}$	Tab1/5	the ratio of the number of shared species to the mean number of species in a single site	Sørensen (1948)
	dissimilarity	$\frac{b + c}{2a + b + c} = \frac{\frac{1}{2}(b + c)}{\frac{1}{2}((a + b) + (a + c))}$	Tab1/6	the ratio of the mean unique species to the mean number of species in a single site	

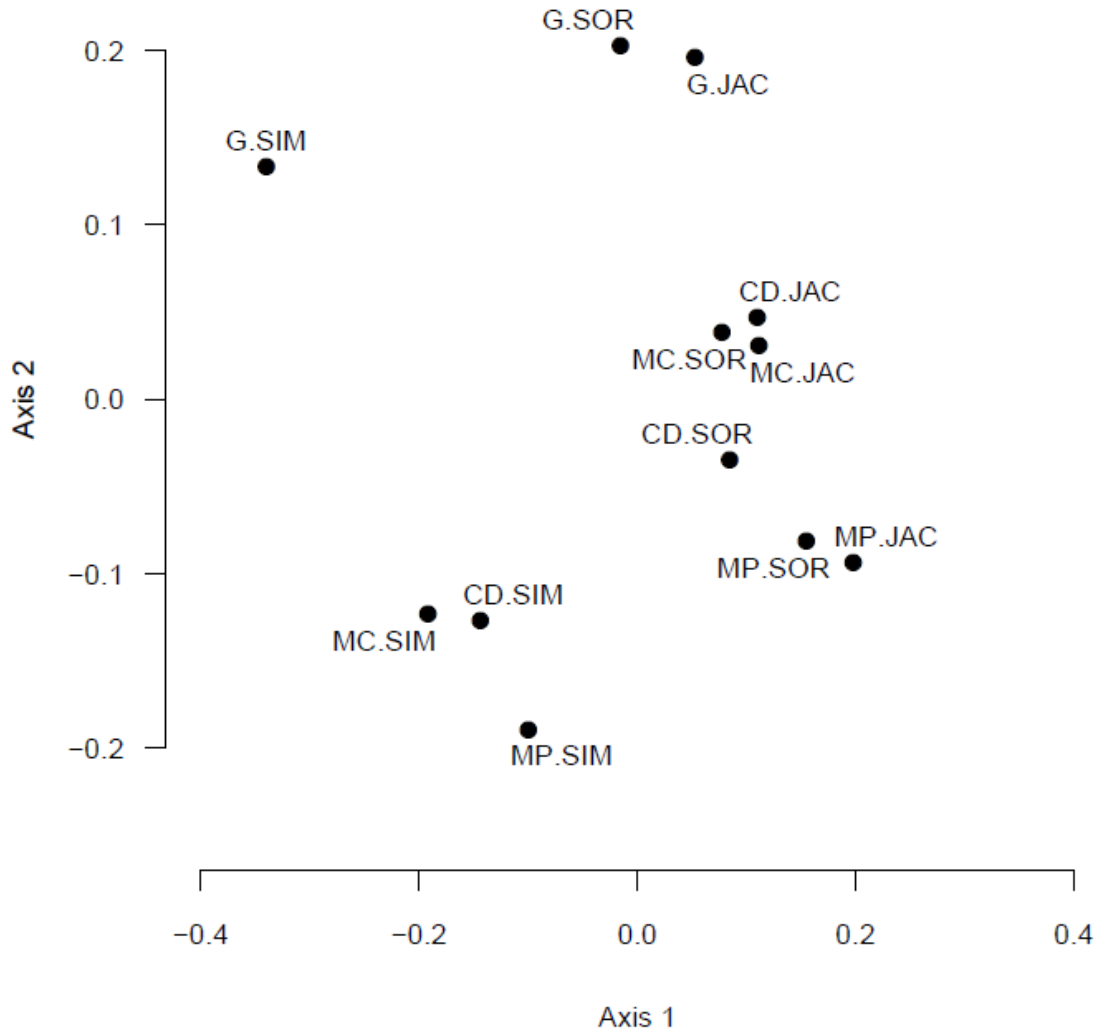
440 Table 2: Overview of multiple site resemblance coefficients (N : number of sites, T : total number of
 441 species, t_i : number of species at site i , n_j : number of sites where species j occurs, o : rank of a species
 442 richness value in the order from the smallest to the largest values, g_o : the frequency of sites with
 443 species richness of rank o , b_{kl} : number of species unique to site k in pairwise comparison with site l ,
 444 b_{lk} : number of species unique to site l in pairwise comparison with site k).

Approach	Family	Form	Equation	Eq. n.	Reference
General	Jaccard	similarity	$\frac{\sum_{i=1}^N t_i - T}{T(N-1)}$	Tab2/1	Koch (1957), Chao et al. (2012) Gotelli & Chao (2013), Arita (2017)
		dissimilarity	$\sum_{j=1}^N \frac{N - n_j}{T \times N} \times \frac{N}{N - 1}$	Tab2/2	Ricotta & Pavoine (2015, in their Appendix S2)
	Simpson	similarity	$\frac{\sum_{j=1}^N n_j - T}{\sum_{j=1}^N p_j - \max(t_i)}$	Tab2/3	Arita (2017)
Sørensen		similarity	$\frac{N}{N-1} \left(1 - \frac{T}{\sum_{i=1}^N t_i}\right)$	Tab2/4	Diserud & Ødegaard (2007), Chao et al. (2012), Gotelli & Chao (2013) and Arita (2017)
		dissimilarity	$\frac{\sum_i (T - t_i)}{(N-1) \sum_i t_i}$	Tab2/5	Ricotta & Pavoine (2015, in their Appendix S2)
Mixed components	Jaccard	dissimilarity	$\frac{[\sum_{k<l} \min(b_{kl}, b_{lk})] + [\sum_{k<l} \max(b_{kl}, b_{lk})]}{[\sum_i t_i - T] + [\sum_{k<l} \min(b_{kl}, b_{lk})] + [\sum_{k<l} \max(b_{kl}, b_{lk})]}$	Tab2/6	Baselga (2012)
		similarity	$\frac{\sum_i t_i - T}{(\sum_i t_i - T) + \sum_{k<l} \min(b_{kl}, b_{lk})}$	Tab2/7	Baselga et al. (2007)
		dissimilarity	$\frac{[\sum_{k<l} \min(b_{kl}, b_{lk})] + [\sum_{k<l} \max(b_{kl}, b_{lk})]}{2[\sum_i t_i - T] + [\sum_{k<l} \min(b_{kl}, b_{lk})] + [\sum_{k<l} \max(b_{kl}, b_{lk})]}$	Tab2/8	Baselga (2010)
Co-diversity	Jaccard	similarity	$\frac{\sum_{j=1}^N n^2 - \sum_{j=1}^N n}{TN(N-1)}$	Tab2/9	Arita (2017)
		similarity	$\frac{\sum_{j=1}^N n^2 - \sum_{j=1}^N n}{2 \sum_{o=1}^N (N - o) g_o}$	Tab2/10	Arita (2017)
		similarity	$\frac{\sum_{j=1}^N n^2 - \sum_{j=1}^N n}{(N-1) \sum_{i=1}^N t_i}$	Tab2/11	Arita (2017)



452 Fig. 1: UPGMA clustering of methods quantifying multiple-site similarities using $1 -$
 453 correlation as distance for 41,503 different data sets with 4 species and 4 sites.

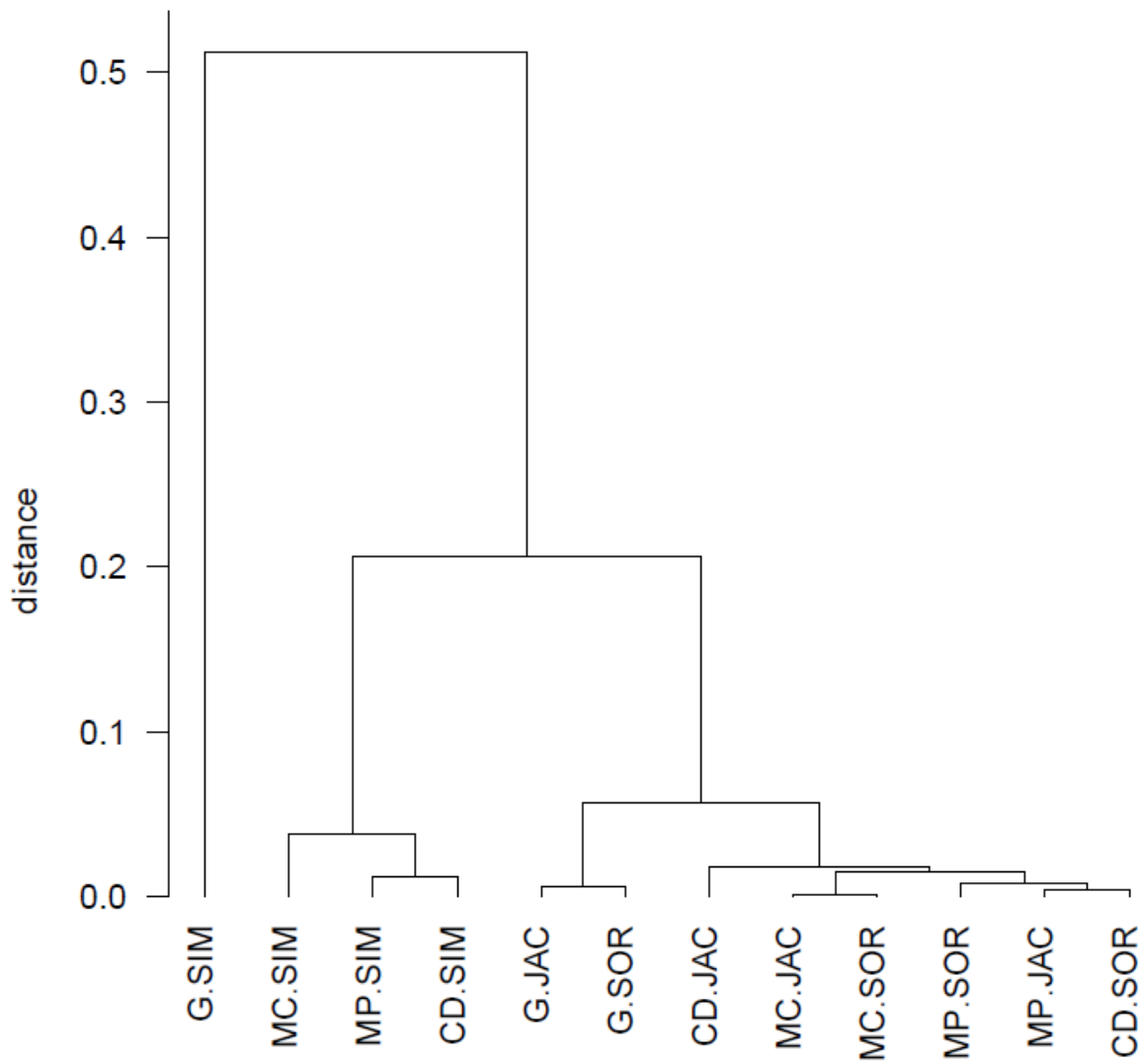
454 Abbreviations include the combination of an approach and a family, where one or two
 455 letters denote an approach (MP: mean pairwise, G: general, MC: mixed component and CO:
 456 co-diversity) and after a dot three letters denote a family (JAC: Jaccard, SIM: Simpson and
 457 SOR: Sørensen).



459

460 Fig 2: Principal coordinates analysis of methods quantifying multiple-site similarities using 1
 461 – correlation as distance for 41,503 different data sets with 4 species and 4 sites. For
 462 abbreviations, see caption to Fig. 1.

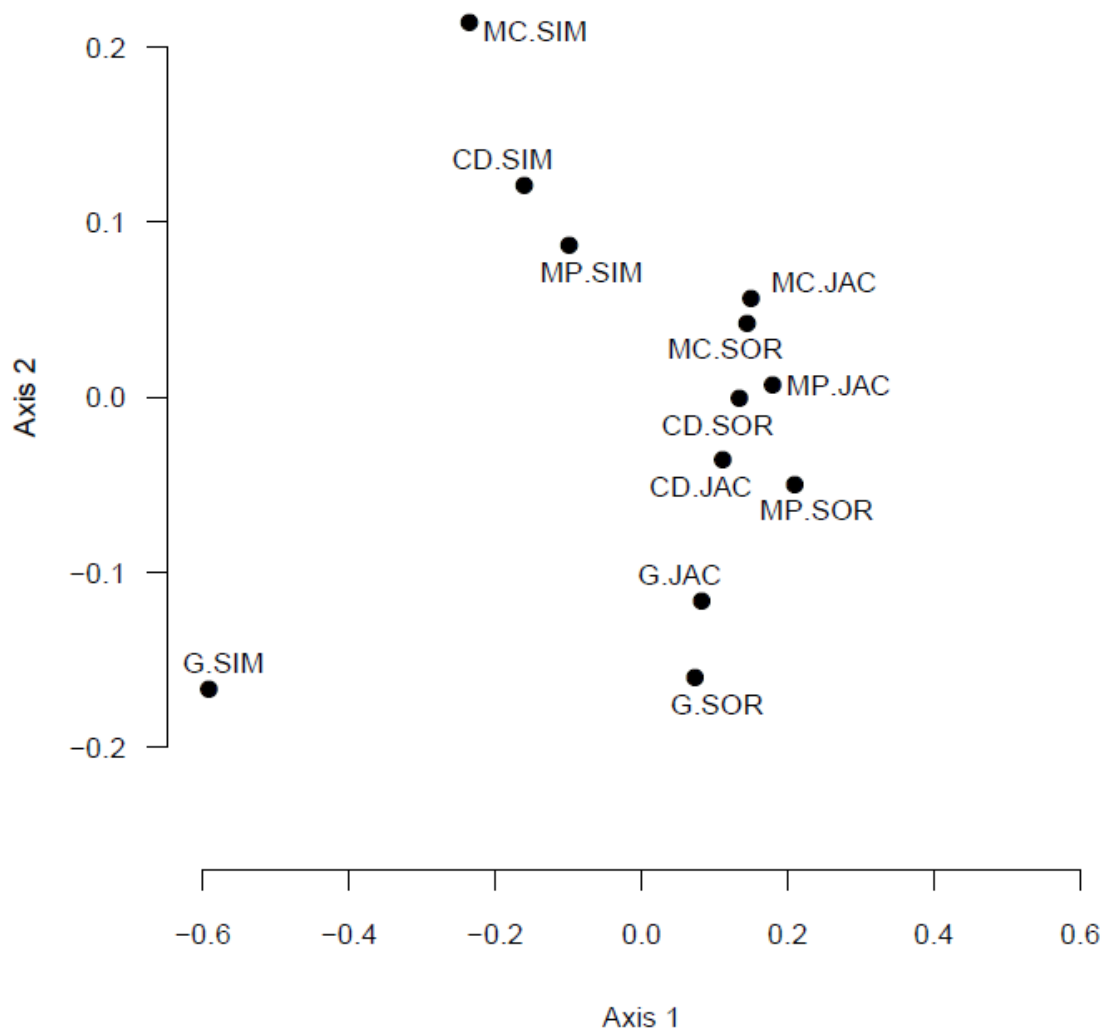
463



464

465 Fig. 3: UPGMA clustering of methods quantifying multiple-site similarities using 1 –
 466 correlation as distance for 100 different data sets with 10 species and 10 sites. For
 467 abbreviations, see caption to Fig. 1.

468

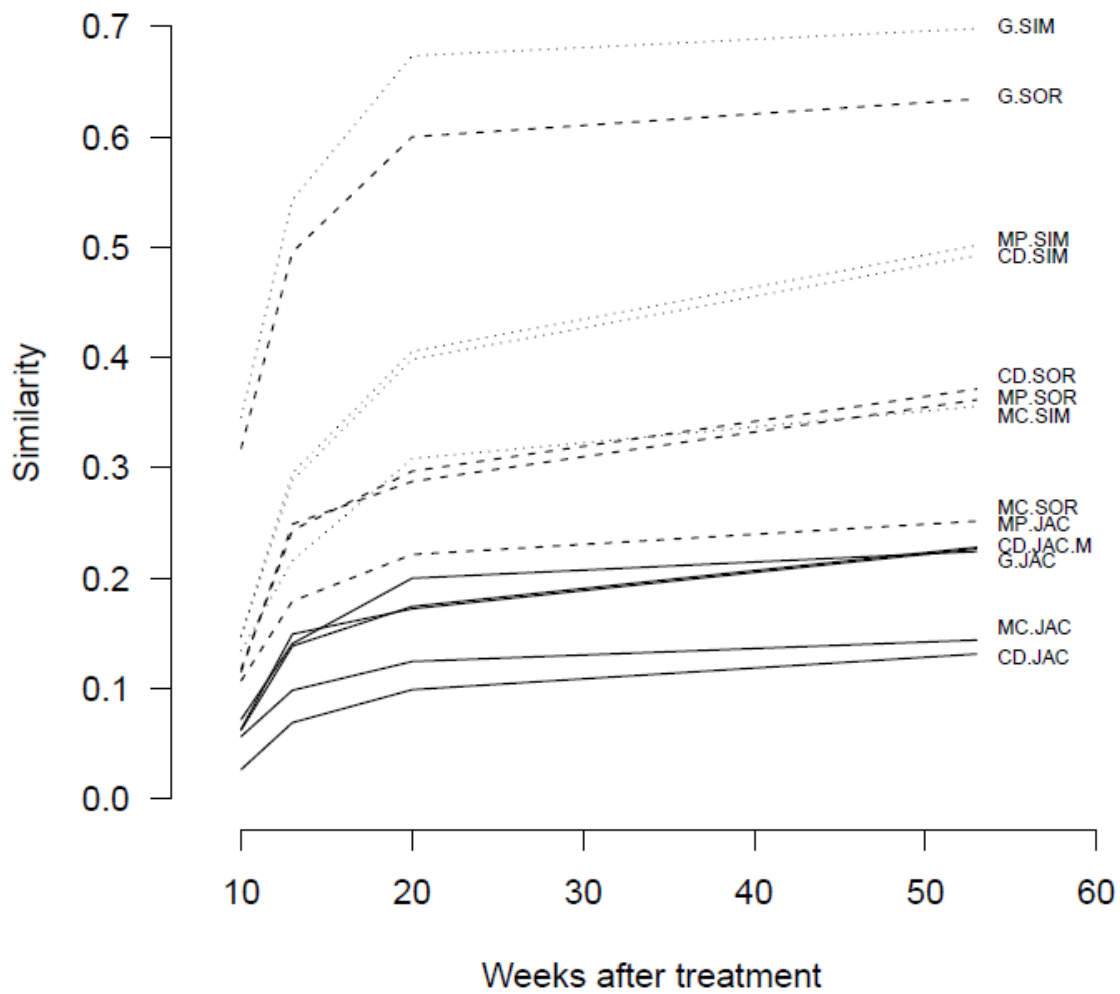


469

470 Fig 4: Principal coordinates analysis of methods quantifying multiple-site similarities using 1-
 471 correlation as distance for 100 different data sets with 10 species and 10 sites. For
 472 abbreviations, see caption to Fig. 1.

473

474



475

476 Fig. 5: Change of community similarity over time (in weeks) depicted by 13 multiple-site
 477 similarity indices. For clarity, data points are connected. For abbreviations, see caption to

478 Fig. 1.

479