

1 **Joint optimization of cluster number and abundance transformation for obtaining**  
2 **effective vegetation classifications**

3

4 **Attila Lengyel**<sup>1,2,\*</sup> (lengyel.attila@okologia.mta.hu)

5 **Flavia Landucci**<sup>3</sup> (flavia.landucci@gmail.com)

6 **Ladislav Mucina**<sup>4,5</sup> (laco.mucina@uwa.edu.au)

7 **James Tsakalos**<sup>4</sup> (james.tsakalos@research.uwa.edu.au)

8 **Zoltán Botta-Dukat**<sup>1,6</sup> (botta-dukát.zoltan@okologia.mta.hu)

9

10 <sup>1</sup> MTA Centre for Ecological Research, Institute of Ecology and Botany, Alkotmány u. 2-4,  
11 H-2163 Vácrátót, Hungary

12 <sup>2</sup> Department of Vegetation Ecology, University of Wrocław, ul. Przybyszewskiego 63, 51-  
13 148 Wrocław, Poland

14 <sup>3</sup> Department of Botany and Zoology, Masaryk University, Kotlářská 2, CZ-611 37 Brno,  
15 Czech Republic

16 <sup>4</sup> School of Biological Sciences, The University of Western Australia, 35 Stirling Hwy,  
17 Crawley WA 6009, Perth, Australia

18 <sup>5</sup> Department of Geography & Environmental Studies, Stellenbosch University, Private Bag  
19 X1, Matieland 7602, Stellenbosch, South Africa

20 <sup>6</sup> MTA Centre for Ecological Research, GINOP Sustainable Ecosystems Group, Klebelsberg  
21 Kuno u. 3, H-8237 Tihany, Hungary

22

23 \*Corresponding author

24

25 **Abstract**

26 **Question:** Is it possible to determine which combination of cluster number and taxon  
27 abundance transformation would produce the most effective classification of vegetation data?  
28 What is the effect of changing cluster number and taxon abundance weighting (applied  
29 simultaneously) on the stability and biological interpretation of vegetation classifications?

30 **Locality:** Europe, Western Australia, simulated data

31 **Methods:** Real data sets representing Hungarian submontane grasslands, European wetlands,  
32 and Western Australian kwongan vegetation, as well as simulated data sets were used. The  
33 data sets were classified using the partitioning around medoids method. We generated  
34 classification solutions by gradually changing the transformation exponent applied to the  
35 species projected covers and the number of clusters. The effectiveness of each classification  
36 was assessed by a stability index. This index is based on bootstrap resampling of the original  
37 data set with subsequent elimination of duplicates. The vegetation types delimited by the most  
38 stable classification were compared with other classifications obtained at local maxima of the  
39 stability values. The effect of changing the transformation power exponent on the number of  
40 clusters, indexed according to their stability, was evaluated.

41 **Results:** The optimal number of clusters varied with the power exponent in all cases, both  
42 with real and simulated data sets. With the real data sets, optimal cluster numbers obtained  
43 with different data transformations recovered interpretable biological patterns. Using the  
44 simulated data, the optima of stability values identified the simulated number of clusters  
45 correctly in most cases.

46 **Conclusions:** With changing the settings of data transformation and the number of clusters,  
47 classifications of different stability can be produced. Highly stable classifications can be  
48 obtained from different settings for cluster number and data transformation. Despite similarly  
49 high stability, such classifications may reveal contrasting biological patterns, thus suggesting  
50 different interpretations. We suggest testing a wide range of available combinations to find  
51 the parameters resulting in the most effective classifications.

52

53 **Keywords**

54 Clustering; Cluster validation; Community similarity; Cover scale; Data type; Multivariate  
55 data analysis; Numerical classification; Stability of classification

56

57 **Abbreviations**

58 MSL = mean standardized lambda; PAM = partitioning around medoids; PCoA = principal  
59 coordinate analysis

60

61 **Nomenclature**

62 The names of high-rank European syntaxa follow Mucina et al. (2016).

63

64 **Introduction**

65 Numerical methods are applied in vegetation classification studies to reduce the  
66 dimensionality of the data in seeking patterns, to increase objectivity in the analyses, and thus  
67 to enhance the reproducibility of results. Still, classification protocols often rely on subjective  
68 decisions that can significantly influence the results (De Cáceres et al. 2015). Subjective  
69 choices can hardly be avoided, yet they should be well-informed and logical to make the  
70 analytical procedures reliable and repeatable. In numerical classifications, according to  
71 Lengyel & Podani (2015), the choice of the number of clusters and the weight attributed to  
72 abundant species relative to scarce species (hence the data transformation), are among the  
73 most influential decisions that have to be considered carefully. If the aim of the classification  
74 is to delimit a pre-set number of vegetation types within the data set, then the choice of the  
75 resulting clusters should be guided by practical considerations. In certain cases there is  
76 reasonable external information available for selecting a transformation function as well. For  
77 instance, if the abundance estimations are deemed inaccurate, only presence/absence data  
78 should be used. Equally, if the purpose of the study is to analyse vegetation types  
79 characterised by dominant species, it is more logical to apply a transformation giving high

80 emphasis to differences in species abundance. However, if the aim of the classification is to  
81 explore variation by separating and differentiating vegetation types, classifications using a  
82 suite of contrasting parameters should be produced. These should be evaluated *a posteriori* in  
83 order to identify the optimal parameter values yielding in the ‘best’ (according to the set  
84 criteria) classification.

85 The optimal number of clusters can be sought for by calculating *cluster effectiveness* (or  
86 *validity*) *index* for classifications with increasing number of clusters. Thus, the optimal  
87 number of clusters is the one where the effectiveness index reaches maximum or minimum,  
88 depending on scaling. This procedure is widely known and regularly applied in classification  
89 studies (e.g. Botta-Dukát et al. 2005; Tichý et al. 2010, 2011). However, we are aware of only  
90 a few examples when authors evaluated different data transformations for finding the optimal  
91 weighting of abundances that would reveal biological patterns most effectively or would lead  
92 to the most stable results. Jensen (1978) evaluated the effect of several data transformations  
93 on classifications and ordinations of a lake vegetation data, and concluded that ‘extreme  
94 transformations’ (i.e. those giving high weight either to high abundance values or, in reverse,  
95 to presence/absence data) can yield significantly different results. This finding was  
96 corroborated by Campbell (1978) and van der Maarel (1979). Wilson (2012) compared the  
97 stability of ordination analyses performed on various vegetation samples using different  
98 transformations of abundance and concluded that the ‘optimal’ transformations depend on  
99 context, such as geographical extent, environmental heterogeneity, disturbance status of the  
100 study area, and quality of abundance estimations. Although, any ‘optimal’ parameterization  
101 supposed to produce a robust classification is specific for the actual data set, the low interest  
102 of researchers in finding them, or at least in assessing the performance of methods they apply,  
103 is surprising, given that vastly different results can be achieved by application of different  
104 abundance scales in multivariate analyses – a fact well known for long time (Austin & Greig-  
105 Smith 1968; Noy-Meir et al. 1975; van der Maarel 1979).

106 In this paper, we introduce a procedure for choosing the combination of two factors, namely  
107 (1) the number of clusters and (2) varying scale of transformation power, assisting in  
108 identification of the most effective classification outcome. Like other approaches aimed at  
109 determination of the optimal number of clusters (e.g. Aho et al. 2008), a general guideline for  
110 finding the optimal transformation would be to find the function that leads to the most stable  
111 of several possible classifications produced by differently parameterized transformation

112 functions. We show that changing one of these two factors has an impact on the optimal  
113 values of the other, which influences the biological interpretation of the classification result,  
114 and therefore we promote their joint optimization. We test this approach using real and  
115 simulated data sets.

116

## 117 **Materials and methods**

### 118 *Grasslands data set*

119 The Grasslands data set consists of phytosociological plots collected in the colline and  
120 montane belts of northern Hungary. This data set represents different types of mesic,  
121 unproductive to moderately productive, grazed, mown, and recently abandoned grasslands on  
122 neutral to acidic soils. Several types can be recognized by their dominant species, e.g.  
123 *Agrostis capillaris*, *Arrhenatherum elatius*, *Danthonia decumbens*, *Festuca rubra* and *Nardus*  
124 *stricta*. However, these types are not floristically distinctly separated, and stands with  
125 different dominant species can be similar in the overall species composition.

### 126 *Wetlands data set*

127 The Wetlands data set was extracted from the WetVegEurope database (Landucci et al. 2015).  
128 It contains plots from Austria, Czech Republic, Germany, Hungary, Poland, Slovakia, and the  
129 Netherlands. In these plots the diagnostic species of the class *Phragmito-Magnocaricetea*  
130 (according to Mucina et al. 2016) should have dominance of at least 25% of the total cover.  
131 Only plots having at least five species and plot sizes between 15 and 50 m<sup>2</sup> were included.  
132 The data set was subject to geographical stratification and to heterogeneity-constrained  
133 random resampling (Lengyel et al. 2011) as modified by Wiser & De Cáceres (2013) in order  
134 to avoid pseudo-replications and maximally diversify the dataset. In this data set, several  
135 types can be distinguished on basis of dominant species, however many of these communities  
136 share similar species pool. Therefore, classifications are expected to vary with changing  
137 power of the data transformation.

### 138 *Kwongan data set*

139 The Kwongan data set is composed of 375 plots of natural shrubland (heath-like) vegetation  
140 of the Geraldton Sandplains (surrounds of the Eneabba township), Western Australia. This  
141 unique, endemic-rich vegetation is supported by sandy soils extremely depleted in phosphorus  
142 (and also nitrogen) – a product of prolonged tectonic quiescence of the Western Australian  
143 landscapes spanning hundreds of millions of years, resulting in lack of soil rejuvenation and  
144 progressive nutrient leaching, combined with relatively stable and predictable climatic  
145 seasonality, and predictable natural fire disturbance (Lambers 2014). This data set exemplifies  
146 an unusual, yet real situation: both alpha and beta diversity are high, resulting in high regional  
147 species pool (gamma diversity). Species dominance (in terms of biomass and projected cover)  
148 in this vegetation is suppressed. We expect that the classification outcomes would be quite  
149 resistant to changes of the magnitude of the data transformation.

150 Characteristics of the three data sets are summarized in Table 1. A more in-depth analysis of  
151 the Grasslands data set is presented, while we focused on the relationship between the  
152 examined methodological decisions and classification stability in the Wetlands and the  
153 Kwongan data sets.

#### 154 *Simulated data*

155 Simulated data matrices consist of  $N$  plots (in the rows) and  $S$  species (in the columns). Plots  
156 belong to  $K$  clusters of equal size, thus the number of plots is  $N/K = n$  in each cluster, and  $n$  is  
157 a pre-defined integer. Ten species occur in each cluster and each species occurs in two  
158 clusters, thus  $S = 10 \times K/2$ . Each species has constant abundance across plots within a cluster,  
159 while the abundances may differ among clusters. The abundances of species within one of the  
160 two clusters where they occur, are drawn from a Poisson-lognormal distribution (Bulmer  
161 1974) where the mean and the standard deviation (SD) of the lognormal distribution are (2; 1)  
162 on log scale. For the other cluster, the order of abundances is reversed, thus if a species was  
163 the most abundant in one of the clusters where it occurs, then this species will be the least  
164 abundant in the other one (considering only species occurring in this cluster). These matrices,  
165 therefore, consist of plots of  $K$  clusters according to raw abundances of species, but  $K/2$   
166 clusters according to presence/absence data because pairs of clusters share the same species  
167 occurring with different abundances. We expect the optimal number of clusters to be  $K/2$  with  
168 low exponents, while with high exponents optimal solution should comprise  $K$  clusters.  
169 Notably, abundance-based clusters are nested within clusters based on presence/absence data.  
170 Within each cluster, plots are identical, thus the clustered structure is initially perfect. An

171 exemplary matrix is shown in Appendix S1. Then, noise was added to this initial matrix  
172 following the method of Gotelli (2000) used for ‘noise test’, but applied to abundances instead  
173 of presence/absence data. This procedure applies a swapping algorithm to introduce noise. In  
174 a single swap, the rows and columns of the original matrix are permuted, and a  $2 \times 2$   
175 submatrix with positive values in the diagonal is chosen randomly. Then the two diagonal  
176 cells are decreased by 1, while abundances in the two off-diagonal cells are increased by 1  
177 individual, thus the sum and the marginal totals of the submatrix do not change. Finally, the  
178 original order of rows and columns is restored. A single swap would affect a sparse matrix  
179 more than one with high fill. Also, large matrices are more ‘resistant’ to the same number of  
180 swaps than small ones. Therefore, noise is added to the matrices in discrete levels, one level  
181 consisting of as many swaps as the number of non-zero elements in the matrix. Our  
182 preliminary analyses suggested that in this way a comparable amount of stochasticity can be  
183 added to matrices of different size and fill.

184 Five simulation series were performed, each of them with five different set-ups. In these  
185 series, one or two parameters were changed systematically in order to generate simulated  
186 matrices that would differ in: i) noise level; ii) size of clusters with number of clusters fixed;  
187 iii) number of clusters with cluster sizes fixed; iv) number and size of clusters with total  
188 number of plots fixed; v) dominance of species. The dominance was changed by modifying  
189 the SD of the lognormal distribution used as input for the Poisson process of species  
190 abundances. When SD is high, there is one or a few highly dominant species within a plot and  
191 many very scarce species, while with lower SD species abundances should be balanced.

### 192 *Classification method*

193 For classifying the data sets, we used the partitioning around medoids method (PAM;  
194 Kaufman & Rousseeuw 1990) using Marczewski-Steinhaus index as the measure of  
195 dissimilarity (Appendix S2). For the Grasslands and Kwongan data set covers of species were  
196 directly estimated on percentage scale in the field, while for the Wetlands data set,  
197 abundances were mostly recorded on Braun-Blanquet or finer ordinal scales. These ordinal  
198 categories were replaced by their midpoint percentages. Cover percentages were power  
199 transformed using the function  $x' = x^a$ , where  $x$  is the original cover value on percentage  
200 scale,  $a$  is the power exponent, and  $x'$  is the transformed cover value. The power exponent  
201 was gradually changed from 0 to 1, with 21 steps by 0.05 in between in case of real data, and  
202 with steps of 0.1 in case of simulations where simpler patterns were expected. Low values of

203 the exponent reduce the effect of differences between species abundances, thus giving more  
204 weight to rare species, while values near 1 give more weight to abundant species. The lowest  
205 number of clusters examined was 2. The highest number of examined clusters was 10 for the  
206 Grasslands data, 40 for the Wetlands and for the Kwongan data, and it varied in simulations  
207 according to the pre-defined number of clusters and sample size. The maximal number of  
208 clusters was arbitrarily determined to balance between computation time and the number of  
209 practically distinguishable vegetation types.

### 210 *Evaluation of classifications*

211 Several approaches for evaluating classifications exist, and each of them involves numerous  
212 indices (e.g. Milligan & Cooper 1985; Vendramin et al. 2010). These approaches include  
213 correlating the original distances between objects and their representations in the  
214 classification (e.g. Rohlf 1974), measuring compactness, connectedness, and separation of  
215 clusters (e.g. Popma et al. 1983), assessing the robustness of the results to changes in  
216 methodological decisions and choice of variables (e.g. Chiang & Mirkin 2010), repetitiveness  
217 (e.g. McIntyre & Blashfield 1980), stability (e.g. Hennig 2007), interpretability (e.g. Tichý et  
218 al. 2010), and predictive power (e.g. Lyons et al. 2016) of the classification, and degree of  
219 divergence from a random classification (e.g. Hunter & McCoy 2004).

220 A family of classification effectiveness (or validity) measures called geometric indices (Aho  
221 et al. 2008) rely on dissimilarities between plots which involve a decision on the weighting of  
222 species abundances. For example, if an effectiveness index uses resemblances calculated by  
223 the Jaccard index (Podani 2000) using presence/absence data, then the classifications  
224 produced on the basis of binary occurrences of species are likely to seem to be ‘better’ than  
225 classifications based on cover percentages. However, not only geometric indices need  
226 decisions on data transformation. The non-geometric OptimClass indices (Tichý et al. 2010),  
227 which use the number of characteristic species of clusters as the measure of effectiveness, can  
228 be calculated from both presence/absence and cover percentage data. As the form of cover  
229 transformation is known to strongly affect the fidelity values of species (Willner et al. 2009),  
230 it is expected that classifications based on presence/absence data would have more character  
231 species, if only binary occurrences are considered for fidelity calculations, while  
232 classifications using cover data would seem less effective.



233 For an unbiased comparison of effectiveness among classifications based on different data  
234 transformations and cluster numbers, it is necessary to compare all classifications to a  
235 standardized reference. The stability index, introduced by Tichý et al. (2011), meets this  
236 criterion. It compares the classification of plots in the original data set with classifications of  
237 its subsets selected by bootstrap resampling with subsequent elimination of duplicates (Tichý  
238 et al. 2011). The similarity between the cluster assignments of resampled plots in the original  
239 classification and in the classification of the subset is calculated using the *mean standardized*  
240 *lambda* (hereafter called MSL), the standardized version of Goodman & Kruskal's lambda  
241 index (Goodman & Kruskal 1954; Appendix S2). In our analysis, we used 50 without-  
242 replacement bootstrap samples for each classification produced by different cluster numbers  
243 and data transformations. MSL was plotted on a so-called *heat map*, in which the colour of  
244 the respective segment of the space defined by two explanatory variables (i.e. the power  
245 exponent and cluster number) refers to the magnitude of the dependent variable (i.e. MSL).

246 The marginal distribution of the heat map can also be examined for determining those  
247 parameter values which are likely to provide the most effective classification outcomes, or the  
248 lowest or highest variation in classification stability. If one of the parameters, e.g. the  
249 exponent, is fixed to an actual value, the mean of the MSL values obtained with changing the  
250 other parameter, that is the number of clusters, gives how stable the classifications obtained  
251 with the actual exponent are on average. By using the SD instead of the mean, the variation of  
252 stability can be expressed, too. Therefore, the SD is a measure of how important the decision  
253 is about one of the two parameters if the other one is fixed to an actual value. The use of  
254 marginal distributions is showed only for the Grasslands data set.

255 The most stable classification of a real data set (i.e. the classification with settings resulting in  
256 the absolute maximum of MSL and the darkest segment on the heat map) was evaluated by  
257 creating a synoptic table containing frequency, average percentage cover, and fidelity of  
258 species. The fidelity of species to clusters was calculated using the phi coefficient on 0 to 100  
259 scale (Chytrý et al. 2002). Species with phi value over 20 were considered 'characteristic',  
260 and only species with Fisher exact test  $p < 0.001$  were considered. Classifications at the  
261 optimal cluster level obtained by different exponents, with special attention to the commonly  
262 used values ( $\alpha = 0, 0.5$  or  $1$ ) and local peaks in stability, were compared on basis of the group  
263 memberships of plots using cross-tabulations, as well as by contrasting their biological  
264 interpretation with the help of characteristic species.

265 Data analyses were performed in the R software environment (version 3.1.2, [www.r-](http://www.r-project.org)  
266 [project.org](http://www.r-project.org)) using the *vegan* (Oksanen et al., <http://cran.r-project.org/package=vegan>), *cluster*  
267 (Maechler et al., <http://cran.r-project.org/package=cluster>), *rapport* (Blagotić & Daróczy,  
268 <http://cran.r-project.org/package=rapport>), and *fields* (Nychka et al., [http://cran.r-](http://cran.r-project.org/package=fields)  
269 [project.org/package=fields](http://cran.r-project.org/package=fields)) packages. R scripts for data simulation, swapping and the  
270 optimization procedure are available in the Appendix S3. We used Juice (Tichý 2002) for data  
271 management and construction of synoptic tables.

272

## 273 **Results**

### 274 *Grasslands data set*

275 The heat map (Fig. 1) showed that the MSL values varied considerably across cluster number  
276 and power exponent. With presence/absence data ( $a = 0$ ), stability was the highest at the five-  
277 cluster solution. From  $a = 0.05$  to  $a = 0.25$ , the three-cluster level was the most stable,  
278 including  $a = 0.15$  where the second highest stability value was obtained (MSL = 0.804).  
279 Between  $a = 0.3$  and  $a = 0.4$ , the stability peaked at two clusters, then from  $a = 0.45$  the four-  
280 cluster solution was optimal until  $a = 0.90$ , while for the higher exponent values again three  
281 clusters were shown to be the best. The absolute maximum value was found with  $a = 0.55$  and  
282 the four-cluster solution, where the stability of the classification was MSL = 0.824. Exponents  
283 between  $a = 0.25$  and  $0.50$  resulted in the highest stability values on average, and the SD of  
284 stability was also the lowest in this interval (Fig. 2). Nevertheless, a second local optimum  
285 was found at  $a = 0.8$ , although the SD was much bigger here. Across the cluster levels, the  
286 three- and four-cluster solutions were the most stable on average, while stability values did  
287 not vary much, except for 2 clusters where SD was the highest.

288 We used the most stable classification (i.e. four clusters and exponent 0.55; hereafter called  
289 ‘Partition A’) as the baseline for the interpretation of all clusters and classifications (Appendix  
290 S4). This classification was identical with what was obtained by  $a = 0.50$ , that is, square-root  
291 transformation. Clusters A1, A2, A3, and A4 are the elements of the Partition A. Cluster A1  
292 represents grasslands of the alliance *Violion caninae*, but some species of the mesic meadows  
293 of the order *Arrhenatheretalia* are also frequent. Cluster A2 contains plots of the  
294 *Arrhenatherion*. This type was recently described as the *Diantho-Arrhenatheretum*

295 association by Lengyel et al. (2016); it represents nutrient-poor, acidic grasslands overgrown  
296 by taller grasses (e.g. *Helictotrichon pubescens*, *Arrhenatherum elatius*) after abandonment or  
297 changing management to mowing. Cluster A3 comprises unproductive meadows and pastures  
298 dominated by *Agrostis capillaris*, *Festuca rubra*, and *Galium verum*. These stands are similar  
299 in species composition to the *Anthoxantho-Agrostietum*, known also from Slovakia and the  
300 Czech Republic. Cluster A3 is also intermediate between *Arrhenatheretalia* and *Violion*  
301 *caninae*. Cluster A4 contains grasslands dominated by *Nardus stricta*, in which species of  
302 waterlogged soils are also present. This type is traditionally also called ‘*Hygro-Nardetum*’  
303 (e.g. Borhidi et al. 2012).

304 In the presence/absence case ( $a = 0$ ), five clusters were differentiated. Hereafter, this  
305 classification is called ‘Partition B’. Cluster B1 included many plots of Cluster A1 and A3,  
306 thus representing mesic meadows with some species of the *Violion caninae*, and matching the  
307 species composition of *Anthoxantho-Agrostietum*. Cluster B2 and B3 contained mostly plots  
308 previously classified to A2, thus differentiating between two subtypes of *Diantho-*  
309 *Arrhenatheretum*: one with more hygrophilous, and one with more forest-steppe species,  
310 respectively. Cluster B4 represents the ‘*Hygro-Nardetum*’ type, thus is similar to Cluster A4.  
311 Cluster B5 contains only two plots similar to the *Anthoxantho-Agrostietum*.

312 With  $a = 0.15$  and three clusters a local peak was detected, to be referred to as Partition C.  
313 Cluster C1 contains many plots representing the types mediating between the  
314 *Arrhenatheretalia* and *Violion caninae*, formerly classified to Clusters A1 and A3. Cluster C2  
315 represents the *Diantho-Arrhenatheretum*, and it is very similar to Cluster A2. Cluster C3  
316 represents the ‘*Hygro-Nardetum*’ and matches with Cluster A4.

317 With  $a = 1$  (= no data transformation), three clusters provided the most stable resolution. This  
318 classification was called Partition D. Cluster D1 represents grasslands on nutrient-poor soils,  
319 including the ‘*Hygro-Nardetum*’ and other types related to the *Violion caninae* and containing  
320 *Nardus stricta*. It contains plots of Cluster A1 and A4. Cluster D2 represents mesic hay  
321 meadows with *Arrhenatherum elatius*, and it shares many plots with Cluster A2. Cluster D3  
322 represents unproductive meadows and pastures with the dominance of *Agrostis capillaris*,  
323 *Briza media* and *Festuca rubra*. Most of its plots were assigned to Cluster A3 and C2.  
324 Therefore, the Partitions C and D similarly separated the *Diantho-Arrhenatheretum* from  
325 other types, but differed in how they delimited two other clusters in the rest of the data set.

326 The cross-tabulation of Partition A against Partitions B, C and D, as well as Partition C  
327 against Partition D are shown in Appendix S5.

### 328 *Wetlands data set*

329 The optimal number of clusters ranged between 3 and 7 when the exponent ranged between 0  
330 and 0.20 (Fig. 3). With higher exponents, the optimal cluster levels increased, too; from  $a =$   
331 0.35 the most stable classifications were found at levels of more than 30 clusters. In the binary  
332 case ( $a = 0$ ), the optimal cluster level was 6, with the square-root transformation ( $a = 0.5$ ) it  
333 was 30, with no transformation ( $a = 1$ ) it was 39. The most stable classification was the one  
334 with  $a = 0.80$  and 40 clusters where MSL was 0.933. At this level clusters were distinguished  
335 according to dominant species that were both constant and character species in many cases.  
336 Using other high exponents (e.g.  $a = 0.50$  or  $a = 1$ ) resulted in very similar classifications,  
337 thus only the comparison of solutions with  $a = 0$  (hereafter called 'Partition W') and  $a = 0.80$   
338 ('Partition Z') are presented using synoptic tables (Appendix S6 and S7, respectively). Since  
339 many phytosociological associations and alliances of wetland vegetation are defined by  
340 dominant species, classifications with high exponents (Partition Z) showed a good  
341 correspondence with low-rank syntaxa. With low exponents, the most stable classifications  
342 revealed markedly different patterns that were difficult to interpret, yet these local optima  
343 possessed much lower stability. With  $a = 0$  (Partition W) differences in species pools offered  
344 some, although not fully satisfactory explanation for the distinction of clusters. Cluster W1  
345 contained many plots of tall-sedge vegetation with short submerged periods and eutrophic  
346 soils (supporting mostly *Magnocaricion gracilis* vegetation). Cluster W2 included mostly  
347 plots of tall-sedge vegetation on sites with poorer nutrient supply (mostly *Magnocaricion*  
348 *gracilis* and *Magnocaricion elatae*). Cluster W3 is characterised, to a large part, by reed  
349 vegetation belonging to the *Phragmition* and *Phalaridion*. Clusters W4 and W5 contained  
350 many plots sampled in wetlands characterised by fluctuating shallow waters (mostly  
351 *Eleocharito-Sagittario*, *Phragmition*, *Glycerio-Sparganion*), however no clear ecological  
352 difference could be recognized between them. Cluster W6 included plots from nutrient-poor  
353 mire vegetation often classified as the *Scheuchzerio-Caricetea*. Obviously, Partition W  
354 showed very low congruence with the syntaxonomical system and Partition Z (Appendix S8).

355 Classifications with  $a = 0$  and  $a = 0.80$  do not differ only in the resolution. As it is shown in  
356 Appendix S8, clusters of the latter are not nested within the former, instead, it is very common  
357 that plots classified to the same cluster at  $a = 0.80$  are assigned to different clusters at  $a = 0$ .

358 *Kwongan data set*

359 MSL values varied much at low levels of cluster numbers (up to 6 clusters) and showed much  
360 less (and also less predictable) variability at cluster levels above 6 (Fig. 4). The highest MSL  
361 values occurred at the cluster levels 2 and 4. The highest classification stability was detected  
362 at the 4-cluster level (for exponents spanning 0.0 and 0.75) or the 2-cluster level (for  
363 exponents spanning 0.8 and 1.0). The most stable classification was obtained with  $\alpha = 0.95$ ,  
364 cluster number = 2, with stability MSL = 0.843.

365 At  $\alpha = 0$ , four clusters were distinguished (Partition K; Appendix S9). Cluster K1 represented  
366 a community with typical species *Hakea candolleana* and *Allocasuarina humilis* found on  
367 free-draining soils. Cluster K2 was identified as *Xylomelum angustifolium-Banksia menziesii*  
368 community thriving on sandy soils on dune swells. Cluster K3 included plots from  
369 *Ecdeiocolea monostachya-Scholtzia laxiflora* community occurring on sandy soils with  
370 slightly elevated clay content in inter-dune depressions, while Cluster K4 represented *Banksia*  
371 *shuttleworthiana-Cristonia biloba* confined to regolith composed of depositional lateritic  
372 scree and sand. Therefore, these clusters represented an edaphic gradient spanning Cluster K2  
373 (deep sandy soils from the sand dune swells) and Cluster K3 (depressions showing elevated  
374 clay content), with Clusters K1 and K4 occupying intermediate position along the gradient. At  
375  $\alpha = 0.95$ , the 2-cluster solution was the most stable one (Partition L; Appendix S10). The  
376 cross-tabulation tables (Appendix S11) showed that all plots of the Cluster K3 were assigned  
377 to the Cluster L1 - the only cluster whose plots were assigned to the same cluster in Partitions  
378 K and L. The Cluster K1 was concentrated in Cluster L1, while most plots of the Clusters K2  
379 and K4 belonged to L2. Partitions K and L similarly recovered the gradient between  
380 vegetation types supported by soils having elevated clay content (represented by Clusters K1  
381 & K3, as well as L1) and sandy soils (as Clusters K2 & K4, and L2) on the basis of  
382 characteristic species of the clusters. The relative position of the clusters in a PCoA ordination  
383 also supports the notion that the main compositional patterns are similarly revealed by  
384 different abundance weighting (Appendix S12).

385 *Simulations*

386 At the noise level 1, where abundances were strongly down-weighted ( $\alpha = 0$  or  $\alpha = 0.1$ ), the  
387 stability was highest at the pre-defined number of four species-pool based clusters (Fig. 5).  
388 From  $\alpha = 0.2$  to  $\alpha = 0.7$ , two peaks were found, namely at the 4- and 8-cluster levels, the latter

389 being of higher stability, and with one intermediate peak at  $a = 0.3$  and seven clusters. Where  
390 abundance differences were not or only slightly reduced ( $a > 0.7$ ), only the 8-cluster peak was  
391 obvious. From the noise level 2 and higher, the stability peaked at the 8-cluster level. As more  
392 levels of noise were added, classifications with low exponent were becoming less and less  
393 stable.

394 Two optimal cluster levels were found where the number of plots in each cluster was 5 (Fig.  
395 6). From  $a = 0$  to  $a = 0.4$ , the 4-cluster peak (corresponding the species-pool-based number of  
396 clusters) was higher, but from  $a = 0.5$  to  $a = 1$  the 8-cluster solution (i.e. the abundance based  
397 optimum) was the most stable one. The pattern of stability was similar, although, less distinct,  
398 with clusters of 10 and 25 plots. However, with 50 plots per cluster, the locations of the  
399 optima were more irregular, with several peaks between four and eight clusters. With 100  
400 plots per cluster, the optima were detected at four clusters for most of the exponent values,  
401 except for  $a = 0.3$  and  $a = 0.4$ .

402 When the number of clusters increased from four with constant cluster sizes, the typical  
403 pattern of lower optima at low exponents and higher optima at high exponents were found in  
404 most cases, yet with some exceptions (Fig. 7). Where the species-pool based cluster number  
405 was two and the abundance-based cluster number was four, three clusters were the most stable  
406 with low exponent and four with high exponent. With higher number of true clusters, the most  
407 stable classification identified the pre-defined cluster numbers correctly: 8, 12, 16, and 24  
408 clusters with higher exponents, and 4, 6, 8, and 12 clusters with lower exponents,  
409 respectively. The point of inflection, when the observed optima shifted from the species-pool-  
410 based level to the abundance-based level, was variable. Yet a broad interval with at least two  
411 local peaks of stability was detectable in all heat maps at intermediate exponent values.  
412 Cluster numbers between the species-pool-based and the abundance-based optima also came  
413 out as optimal in some cases, especially with exponents near the inflection value.

414 A very similar pattern was found when the number of clusters and cluster sizes were changed  
415 with constant sample size (Appendix S13). The species-pool-based and the abundance-based  
416 cluster numbers were recovered correctly as local or global peaks. Between them,  
417 intermediate levels also gained high stability values, but they were identified as optimal only  
418 in a few cases.

419 With  $SD = 0.1$  the optimal cluster level was four clusters irrespective of the exponent value  
420 (Appendix S13). Using  $a > 0.5$  classifications of 7 and 8 groups showed local peaks. With  
421 increasing  $SD$ , the stability of classifications with eight clusters and high exponent also  
422 increased. With  $SD = 4$ , the 8-cluster solutions appeared the most stable, except for when  $a =$   
423 0, that is, in the binary case.

## 424 **Discussion**

### 425 *Evaluation of the real data*

426 The choice of data transformation and cluster number influences the delimitation of  
427 vegetation types, as concluded in several other studies (e.g. Jensen 1978; Lengyel & Podani  
428 2015). Certain types (e.g. *Diantho-Arrhenatheretum* in the Grasslands data set) are relatively  
429 robust to changes in the examined parameters, while others (e.g. transitional types between  
430 *Arrhenatheretalia* and *Violion caninae*) are more sensitive. When it comes to making an  
431 unambiguous distinction between vegetation types for practical (such as management)  
432 purposes or syntaxonomical revision, it is crucial to consider that different weighting of  
433 abundant species may have implications for the delimitation of vegetation units, and thus for  
434 the future applicability of the classification.

435 The Wetlands data set showed that the optimal cluster level can markedly differ if different  
436 data transformations are used. While presence/absence data yielded six stable clusters that  
437 represented types with more or less different species pools, accounting for differences in  
438 abundances raised the optimal levels over 30, where each cluster is separated according to the  
439 dominant species. The fact that the high number of stable clusters obtained using high  
440 exponent were not nested within the few stable clusters based on presence/absence data, is a  
441 clear indication that different data transformations can reveal different types of biological  
442 patterns. With low exponents, classifications were best explained by patterns generated by  
443 habitat-specific species-pools, while with high exponents, community types differing in fine-  
444 scale environmental variation, temporal variability and site history were revealed. It is of  
445 interest, that in our study, 40 clusters was the finest classification level examined due to a  
446 compromise between practical and scientific reasons, but in reality the optimal number of  
447 clusters in the Wetlands data set could have been even higher.

448 The Kwongan data provided a special insight into the interaction of data transformation and  
449 cluster number. Changing the exponent changed the optimal number of clusters as well, and  
450 the resulting stable classifications were moderately congruent. However, even these,  
451 seemingly less similar classifications revealed the most important ecological pattern on the  
452 basis of faithful species — the soil gradient, although fine patterns of transitional subtypes  
453 between the extremes were not detected equally well. The Kwongan data set, due to its high  
454 beta diversity and balanced within-plot abundance distribution, was less sensitive to changes  
455 in data transformation and cluster number in terms of biological interpretation, even though  
456 the assignment of plots showed some variation.

#### 457 *Lessons from the simulations*

458 In the simulations, we generated data structure with contrasting patterns with respect to  
459 occurrence information. If abundance information were emphasized, the true number of  
460 clusters (vegetation types) was twice as high as in cases where only presence/absence data  
461 were considered, hence we differentiated a ‘species-pool-based’ and an ‘abundance-based’  
462 number of clusters. In reality, however, also an opposite can be observed, where a few species  
463 can be dominant in habitats with different species pools. In such a case the number of  
464 abundance-based clusters could be lower than those based on species-pools, as it was seen  
465 with the Kwongan data set.

466 We expected that *weak* data transformations (the exponent being close to 1) which preserve  
467 the differences in original abundance patterns, would yield a higher cluster number, while  
468 *strong* transformations (the exponent approaching 0) which significantly reduce abundance  
469 differences would find the half of this number of clusters optimal. Our results confirmed this  
470 expectation.

471 We introduced stochasticity to artificial data using a similar method as that by Gotelli (2000)  
472 called ‘noise test’. This type of noise made classifications with stronger transformations less  
473 stable than those involving weak transformations. This result can be understood by recalling  
474 how we generated species abundances and noise. The species abundances had been drawn  
475 from a Poisson-lognormal distribution, which resulted in many scarce and few abundant  
476 species. Considering that the artificial matrices are designed in a way that their matrix fill is  
477 low, swapping individuals can moderately reduce the abundance of species in a plot, or it can  
478 slightly increase less abundant species, or make absent species present with low abundance.



479 However, it is unlikely to make an abundant species absent in a plot, or to make an absent  
480 species very abundant. As a result, the applied noise affected binary information more than  
481 the proportions of abundances which determine classifications involving weak data  
482 transformations. We believe that this type of noise simulates a common form of stochasticity  
483 in nature that is caused by random death of individuals followed by random colonization.

484 The simulations have revealed several tendencies in classification stability as related to cluster  
485 number, data transformation, and sample properties. With increasing size of clusters, the  
486 number of abundance-based clusters was underestimated, while the number of clusters based  
487 on species pools was detected correctly. Despite this observation with both fixed and  
488 changing total sample size, we cannot offer a clear explanation for this finding.

489 Based on the tests with modified pre-defined number of clusters with fixed cluster sizes, the  
490 stability as optimality criterion seems to track the changes correctly in most cases. However,  
491 when the number of clusters based on presence/absence data was two, the most stable  
492 classifications were obtained at the three-cluster level with strong transformation. (With weak  
493 transformations, the abundance-based number of clusters was correctly found at the level of  
494 four clusters.) Moreover, in a few cases, optima were indicated between the species-pool-  
495 based and the abundance-based levels. When the total sample size was fixed, but number and  
496 size of clusters changed, stability performed similarly well. Some inconsistency was found at  
497 four abundance-based clusters, where the most stable level was found at two clusters for all  
498 but one value of the exponent. Surprisingly, the exception was the binary case ( $a = 0$ ) where  
499 all classifications were generally less stable and the optimum was at the pre-defined number  
500 of clusters based on abundance, i.e. four clusters. This contradicts our expectation and we  
501 have no clear explanation for this. Despite the above mentioned spurious exceptions, the  
502 stability seemed rather robust and accurate across a wide range of cluster numbers with PAM.  
503 In real situations, mapping a goodness of classification measure as a function of data  
504 transformation and cluster number would help avoiding less effective parameter  
505 combinations.

506 Testing the effect of community dominance on stability by changing the logarithm of SD of  
507 species abundances revealed that at the lowest dominance (i.e. low SD), the number of  
508 clusters based on species pool was optimal regardless of data transformation. As dominance  
509 increased, abundance-based cluster number became more stable and was identified as optimal.  
510 This is in line with the common experience that in monodominant vegetation types (e.g.

511 aquatic and marsh vegetation) classifications based on abundance data are more effective and  
512 can markedly differ from presence/absence-based classifications, while when the species  
513 abundances are more balanced, accounting for abundance differences does not give  
514 significantly different or more effective classification than what is obtained by species  
515 composition.

#### 516 *Concluding remarks*

517 Classification stability depends both on cluster number and data transformation. The trend of  
518 stability along increasing power exponent varies across cluster numbers, and vice versa, the  
519 number of clusters resulting in the most stable classifications depends on data transformation.  
520 Slight changes in any of these two factors may change the stability of a classification, hence  
521 different biological conclusions can be reached. At the same time, similarly effective  
522 classifications can be produced using different combinations of parameters. Finding such  
523 local optima contributes to the thorough understanding of biological patterns in the sample.

524 Stability, as proposed by Tichý et al. (2011), is a standardized measure of classification  
525 effectiveness because every single classification is compared to classifications of its without-  
526 replacement bootstrap subsamples obtained with exactly the same methods. We have chosen  
527 this index in our study because of this advantage. However, there are many other measures of  
528 effectiveness, but we have chosen not to evaluate them experimentally in this paper. For  
529 answering specific research questions, other indices may be more appropriate than stability. In  
530 such cases the workflow of testing the effect of data transformation and cluster number on  
531 classification effectiveness, and the visualization of results should be the same as we  
532 presented, only the measure of effectiveness should be replaced by an alternative. Moreover,  
533 it is also possible to perform the optimization analysis using several different effectiveness  
534 measures, and then combine the results in order to identify the classification which is the most  
535 effective on average across the applied indices.

536 Apart from the cluster number and the power exponent, we see no obstacles to test the effect  
537 of other types of methodological decisions using our approach. For example, an effectiveness  
538 measure might be calculated for classifications obtained by different values for the  $\beta$   
539 parameter of the flexible clustering method by Lance & Williams (1967), and the  $\beta$  value  
540 providing the most stable classification might be determined. Moreover, our optimization  
541 approach can easily be adapted to ordinations, too. If the cluster effectiveness index applied

542 here is substituted by a measure of stability of ordinations (as done by Wilson 2012), the  
543 effect of data transformation on the stability of ordinations can be evaluated systematically.  
544 The extension of the optimization procedure presented here beyond data transformation and  
545 cluster number is a future direction of our research.

546

#### 547 **Acknowledgements**

548 The Authors are grateful to Miquel De Cáceres, David W. Roberts, Lubomír Tichý, and Ákos  
549 Bede-Fazekas for their helpful comments. A.L. and Z.B.D. were supported by the GINOP  
550 2.3.3-15-2016-00019 project. The research stay of A.L. at the University of Wrocław was  
551 supported by the POLONEZ programme (grant 2016/23/P/NZ8/04260). L.M. thanks the Iluka  
552 Chair in Vegetation Science and Biogeography for logistic support. The work of L.M. and  
553 J.T. was also supported by ARC Linkage grant LP150100339.

554

#### 555 **Authors contributions**

556 A.L. outlined the main idea, performed data analysis and wrote the initial manuscript, Z.B.D.  
557 contributed with discussion in all stages of the work, F.L. helped in preparation of the  
558 Wetlands data set and the evaluation of the analysis, L.M. and J.T. contributed by providing  
559 the Kwongan data set and evaluating the results, L.M. and J.T. performed linguistic revisions  
560 of early versions of the text. All authors critically commented on the manuscript and the  
561 supplementary materials.

562

#### 563 **References**

564 Aho, K., Roberts, D.W. & Weaver, T. 2008. Using geometric and non-geometric internal  
565 evaluators to compare eight vegetation classification methods. *Journal of Vegetation Science*  
566 19: 549–562.

- 567 Austin, M.P. & Greig-Smith, P. 1968. The application of quantitative methods to vegetation  
568 survey: II. Some methodological problems of data from rain forest. *Journal of Ecology* 56:  
569 827–844.
- 570 Borhidi, A., Kevey, B. & Lendvai, G. 2012. *Plant communities of Hungary*. Akadémiai  
571 Kiadó, Budapest, HU.
- 572 Botta-Dukát, Z., Chytrý, M., Hájková, P. & Havlová, M. 2005. Vegetation of lowland wet  
573 meadows along a climatic continentality gradient in Central Europe. *Preslia* 77: 89–111.
- 574 Bulmer, M.G. 1974. On fitting the Poisson lognormal distribution to species-abundance data.  
575 *Biometrics* 30: 101–110.
- 576 Campbell, B.M. 1978. Similarity coefficients for classifying plots. *Vegetatio* 37: 101–108.
- 577 Chytrý, M., Tichý, L., Holt, J. & Botta-Dukát, Z. 2002. Determination of diagnostic species  
578 with statistical fidelity measures. *Journal of Vegetation Science* 13: 79–90.
- 579 Chiang, M. & Mirkin, B. 2010. Intelligent choice of the number of clusters in k-means  
580 clustering: An experimental study with different cluster spreads. *Journal of Classification* 27:  
581 3–40.
- 582 De Cáceres, M., Chytrý, M., Agrillo, E., Attorre, F., Botta-Dukát, Z., Capelo, J., Czúcz, B.,  
583 Dengler, J., Ewald, J., (...) & Wiser, S.K. 2015. A comparative framework for broad-scale  
584 plot-based vegetation classification. *Applied Vegetation Science* 18: 543–560.
- 585 Goodman, L. & Kruskal, W. 1954. Measures of association for cross classifications. *Journal*  
586 *of the American Statistical Association* 49: 732–764.
- 587 Gotelli, N.J. 2000. Null model analysis of species co-occurrence patterns. *Ecology* 81: 2606–  
588 2621.
- 589 Hennig, C. 2007. Cluster-wise assessment of cluster stability. *Computational Statistics &*  
590 *Data Analysis* 52: 258–271.
- 591 Hill, M.O. 1973. Diversity and evenness: a unifying notation and its consequences. *Ecology*  
592 54: 427–432.

- 593 Hunter, J.C. & McCoy, R.A. 2004. Applying randomization tests to cluster analyses. *Journal*  
594 *of Vegetation Science* 15: 135–138.
- 595 Jensen, S. 1978. Influences of transformation of cover values on classification and ordination  
596 of lake vegetation. *Vegetatio* 37: 19–31.
- 597 Kaufman, L. & Rousseeuw, P.J. 1990. *Finding groups in data: An introduction to cluster*  
598 *analysis*. John Wiley & Sons, New York, US.
- 599 Király, G. (ed.) 2009. *New Hungarian Herbal. The vascular plants of Hungary*. Identification  
600 key. Aggteleki Nemzeti Park Igazgatóság, Jósvalő, HU. (in Hungarian)
- 601 Lance, G.N. & Williams, W.T. 1967. A general theory of classificatory sorting strategies. I.  
602 Hierarchical systems. *Computer Journal* 9: 373–380.
- 603 Landucci, F., Řezníčková, M., Šumberová, K., Chytrý, M., Aunina L., Biță-Nicolae, C.,  
604 Bobrov, A., Borsukevych, L., Brisse, H., (...) & Willner W. 2015. WetVegEurope: a database  
605 of aquatic and wetland vegetation of Europe. *Phytocoenologia* 45: 187–194.
- 606 Lambers, H. (ed.) 2014. *Plant life on the sandplains in Southwest Australia: A global*  
607 *biodiversity hotspot*. UWA Publishing, Crawley, AU.
- 608 Lengyel, A., Chytrý, M. & Tichý, L. 2011. Heterogeneity-constrained random resampling of  
609 phytosociological databases. *Journal of Vegetation Science* 22: 175–183.
- 610 Lengyel, A. & Podani, J. 2015. Assessing the relative importance of methodological decisions  
611 in classifications of vegetation data. *Journal of Vegetation Science* 26: 804–815.
- 612 Lengyel, A., Illyés, E., Bauer, N., Csiky, J., Király, G., Purger, D. & Botta-Dukát, Z. 2016.  
613 Classification and syntaxonomical revision of mesic and semi-dry grasslands in Hungary.  
614 *Preslia* 88: 201–228.
- 615 Lötter, M.C., Mucina, L. & Witkowski, E. 2013. The classification conundrum: species  
616 fidelity as leading criterion in search of a rigorous method to classify a complex forest data  
617 set. *Community Ecology* 14: 121–132.

- 618 Lyons, M.B., Keith, D.A., Warton, D.I., Somerville, M. & Kingsford, R.T. 2016. Model-  
619 based assessment of ecological community classifications. *Journal of Vegetation Science* 27:  
620 704–715.
- 621 McIntyre, R.M. & Blashfield, R.K. 1980. A nearest-centroid technique for evaluating the  
622 minimum-variance clustering procedure. *Multivariate Behavioral Research* 15: 225–238.
- 623 Milligan, G.W. & Cooper, M.C. 1985. An examination of procedures for determining the  
624 number of clusters in a data set. *Psychometrika* 50: 159–179.
- 625 Mucina, L., Bültmann, H., Dierßen, K., Theurillat, J.-P., Raus, T., Čarni, A., Šumberová, K.,  
626 Willner, W., Dengler, J., (...) & Tichý, L. 2016. Vegetation of Europe: hierarchical floristic  
627 classification system of vascular plant, bryophyte, lichen, and algal communities. *Applied*  
628 *Vegetation Science* 19: 3–264.
- 629 Noy-Meir, I., Walker, D. & Williams, W.T. 1975. Data transformations in ecological  
630 ordination: II. On the meaning of data standardization. *Journal of Ecology* 63: 779–800.
- 631 Podani, J. 2000. *Introduction to the exploration of multivariate biological data*. Backhuys,  
632 Leiden, NL.
- 633 Podani, J. & Feoli, E. 1991. A general strategy for the simultaneous classification of variables  
634 and objects in ecological data tables. *Journal Vegetation Science* 2: 435–444.
- 635 Popma, J., Mucina, L., van Tongeren, O. & van der Maarel, E. 1983. On the determinants of  
636 optimal levels in phytosociological classification. *Vegetatio* 52: 65–75.
- 637 Roberts, D.W. 2015. Vegetation classification by two new iterative reallocation optimization  
638 algorithms. *Plant Ecology* 216: 741–758.
- 639 Rohlf, F.J. 1974. Methods of comparing classifications. *Annual Review of Ecology &*  
640 *Systematics* 5: 101–113.
- 641 Rozbrojová, Z., Hájek, M. & Hájek, O. 2010. Vegetation diversity of mesic meadows and  
642 pastures in the West Carpathians. *Preslia* 82: 307–332.

- 643 Tichý, L. 2002. JUICE, software for vegetation classification. *Journal of Vegetation Science*  
644 13: 451–453.
- 645 Tichý, L., Chytrý, M. & Šmarda, P. 2011. Evaluating the stability of the classification of  
646 community data. *Ecography* 34: 807–813.
- 647 Tichý, L., Chytrý, M., Hájek, M., Talbot, S.S. & Botta-Dukát, Z. 2010. OptimClass: Using  
648 species-to-cluster fidelity to determine the optimal partition in classification of ecological  
649 communities. *Journal of Vegetation Science* 21: 287–299.
- 650 van der Maarel, E. 1979. Transformation of cover-abundance values in phytosociology and its  
651 effects on community similarity. *Vegetatio* 39: 97–114.
- 652 Vendramin, L., Campello, R.J.G.B. & Hruschka, E.R. 2010. Relative clustering validity  
653 criteria: A comparative overview. *Statistical Analysis & Data Mining* 3: 209–235.
- 654 Willner, W., Tichý, L. & Chytrý, M. 2009. Effects of different fidelity measures and contexts  
655 on the determination of diagnostic species. *Journal of Vegetation Science* 20: 130–137.
- 656 Wilson, J.B. 2012. Species presence/absence sometimes represents a plant community as well  
657 as species abundances do, or better. *Journal of Vegetation Science* 23: 1013–1023.
- 658 Wiser, S.K. & De Cáceres, M. 2013. Updating vegetation classifications: an example with  
659 New Zealand's woody vegetation. *Journal of Vegetation Science* 24: 80–93.

660

## 661 **List of Appendices**

662 Appendix S1: Simulation data example

663 Appendix S2: Mathematical formulae

664 Appendix S3: R scripts

665 Appendix S4: Grasslands synoptic table (Partition A)

666 Appendix S5: Cross-tabulations of partitions of the Grasslands data set

- 667 Appendix S6: Wetlands synoptic table (Partition W)
- 668 Appendix S7: Wetlands synoptic table (Partition Z)
- 669 Appendix S8: Wetlands cross-tabulations
- 670 Appendix S9: Kwongan synoptic table (Partition K)
- 671 Appendix S10: Kwongan synoptic table (Partition L)
- 672 Appendix S11: Kwongan cross-tabulation
- 673 Appendix S12: Kwongan ordination
- 674 Appendix S13: Additional heat maps of the simulated data sets
- 675
- 676



677 Tables

678

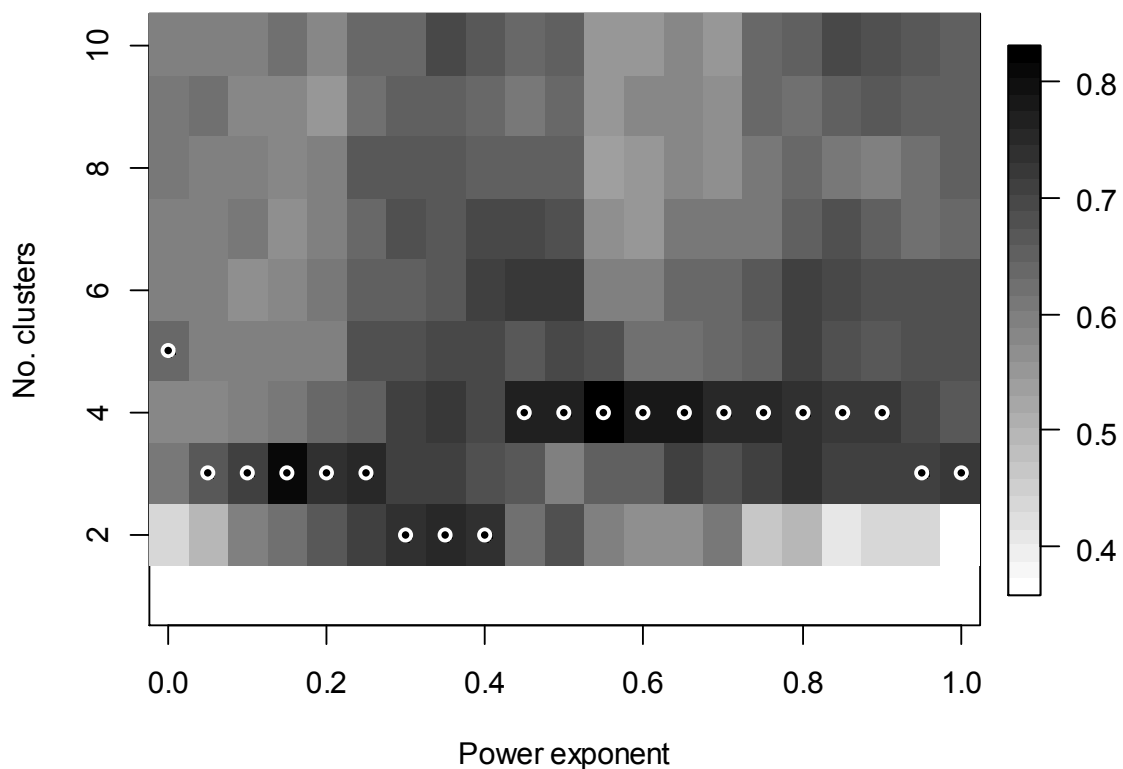
679 **Table 1.** Characteristics of the real vegetation data sets

	Grasslands	Wetlands	Kwongan
Vegetation type	mesic grasslands	reeds and sedge beds	sclerophyllous scrub
Geographical location	Northern Hungary	Central and Western Europe	Geraldton Sandplains, Western Australia
Nr. of plots	55	2725	379
Plot size (m <sup>2</sup> )	25	15 to 50	100
Number of species			
total	269	844	645
mean per plot	37.78	12.52	49.33
minimum per plot	18	5	20
maximum per plot	54	43	85
Mean diversity of order 1*	12.22	4.8	37
Mean evenness per plot**	0.32	0.38	0.75
Mean SD of species covers	8.77	20.60	1.79
Mean 25–75% quantiles of species covers	0.51–2.52	2.05–6.60	1.00–1.15

\*according to Hill (1973)

680 \*\*mean of diversity of order 1 divided by diversity of order 0, the latter being species  
681 richness

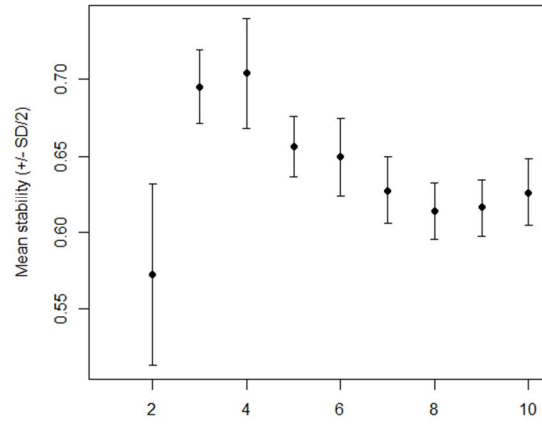
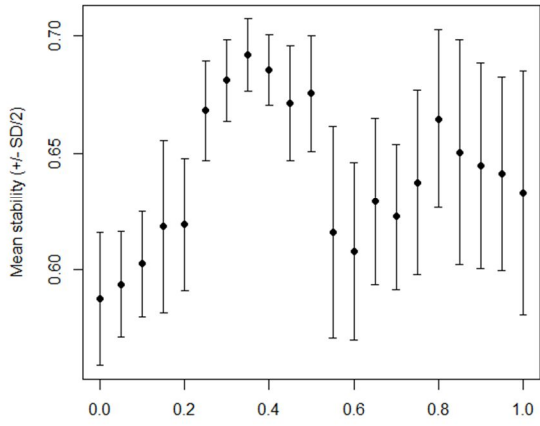
682



684

685 Fig. 1. Analysis of the Grasslands data set showing the heat map of classification stability  
 686 obtained using different parameters for number of clusters and power exponent. Darkness of  
 687 the segments correlate with the value of the mean standardized Goodman & Kruskal's lambda  
 688 (MSL), where the darkest segments marking the combinations of parameters leading to the  
 689 most stable classifications. White circles with black dots indicate the optimal number of  
 690 clusters for a given exponent.

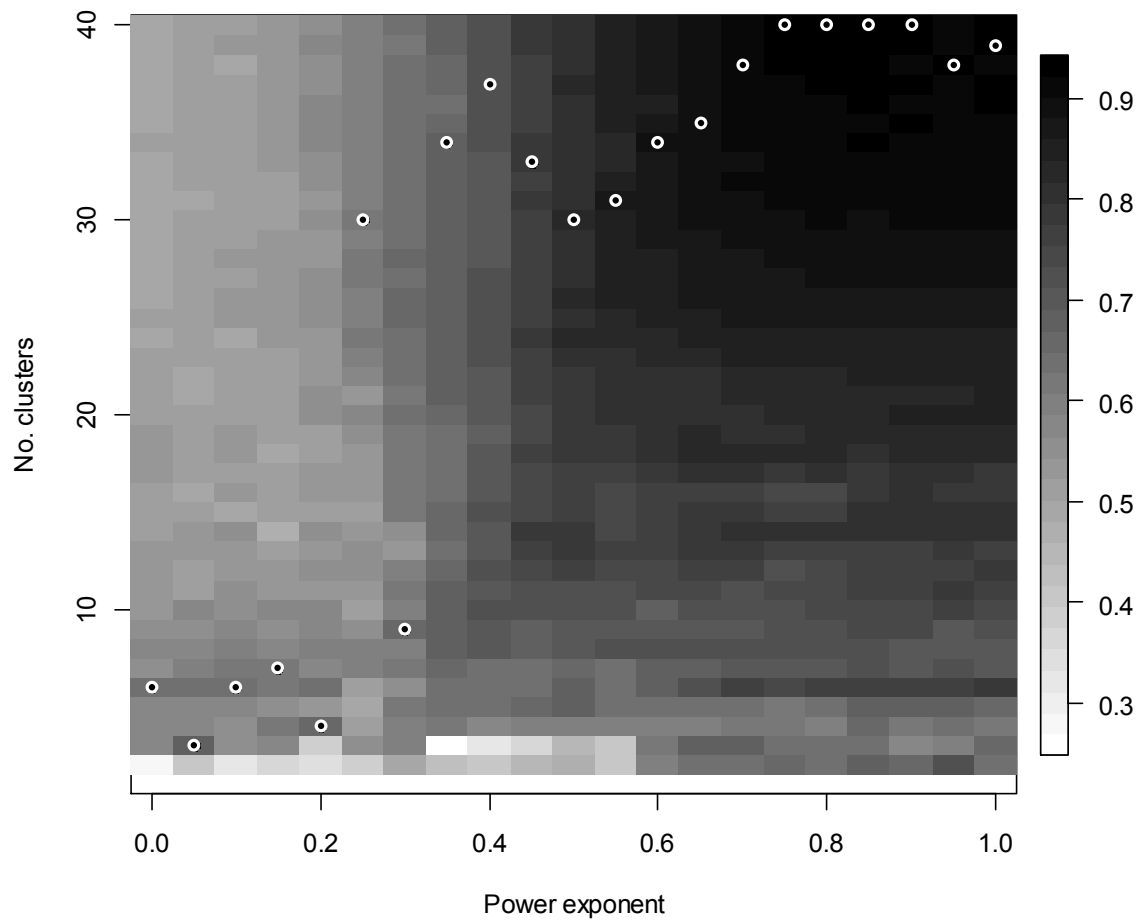
691



692

693 Fig. 2. Mean and standard deviation as error bars of the marginal of the heat map of the  
 694 Grasslands data set.

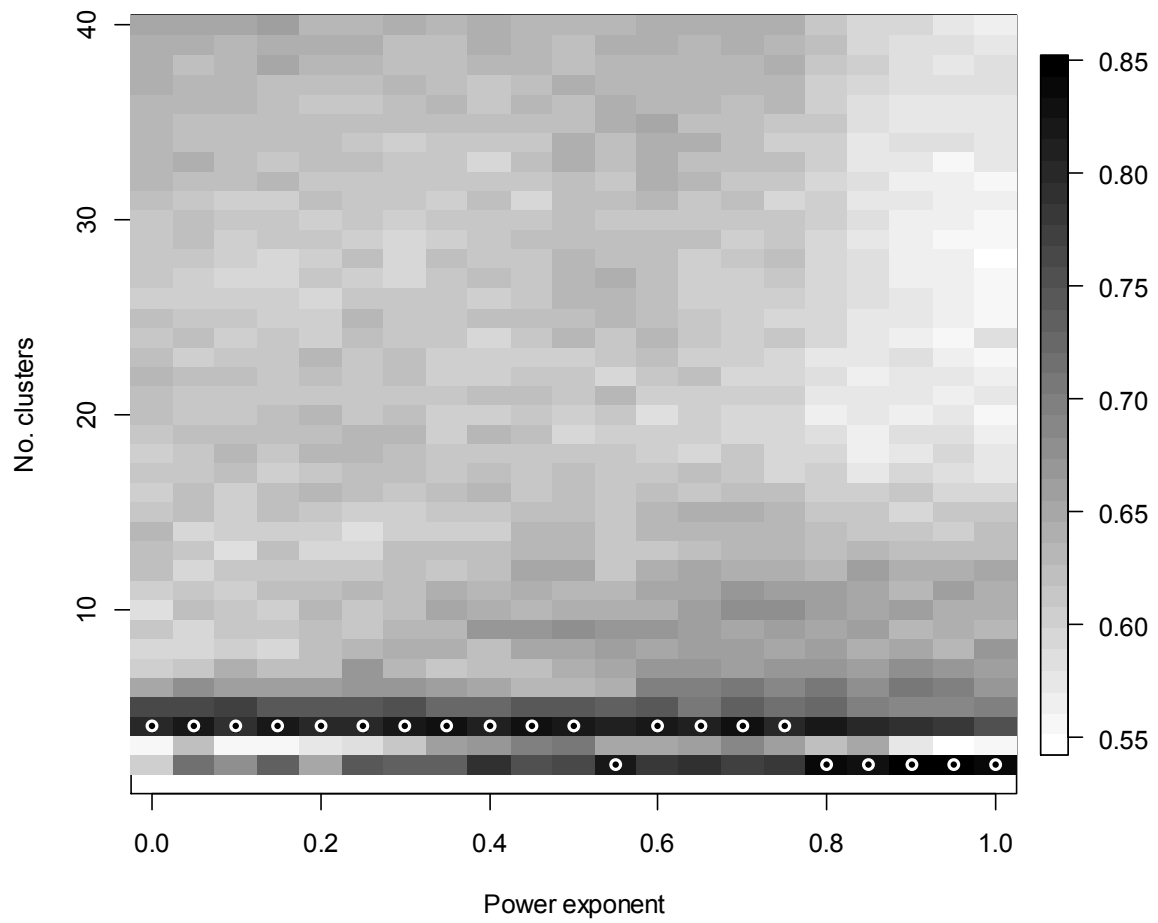
695



696

697 Fig. 3. Analysis of the Wetlands data set showing the heat map of classification stability  
 698 obtained using different parameters for number of clusters and power exponent. For the  
 699 meaning of shading and other symbols see Fig. 1.

700



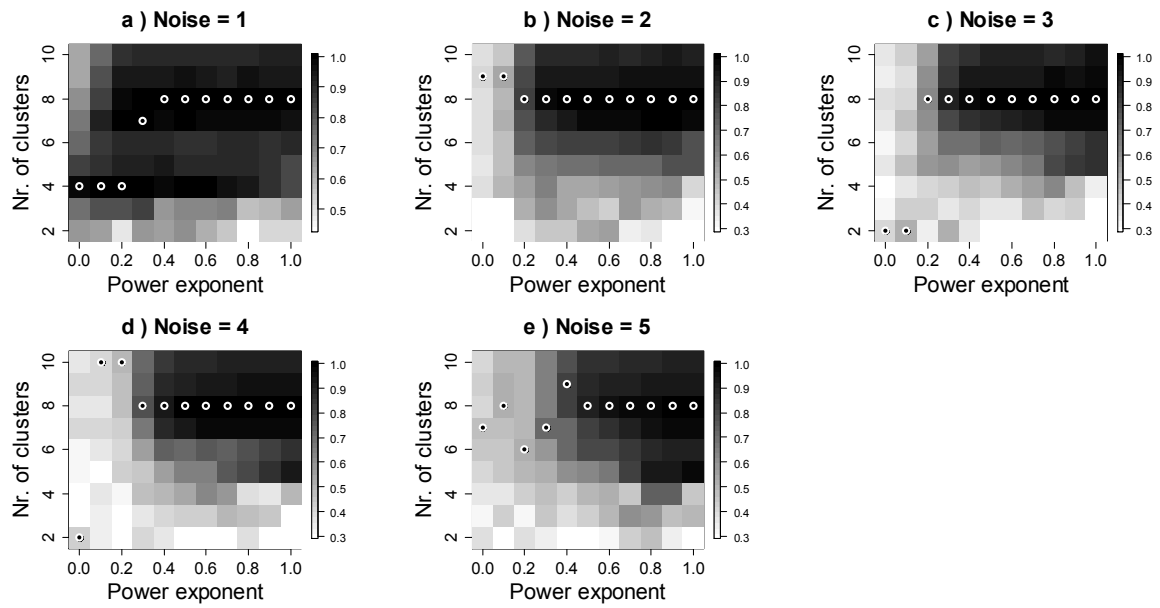
701

702 Fig. 4. Analysis of the Kwongan data set showing the heat map of the classification stability  
 703 obtained using different parameters for number of clusters and power exponent. For the  
 704 meaning of shading and other symbols see Fig. 1.

705

706

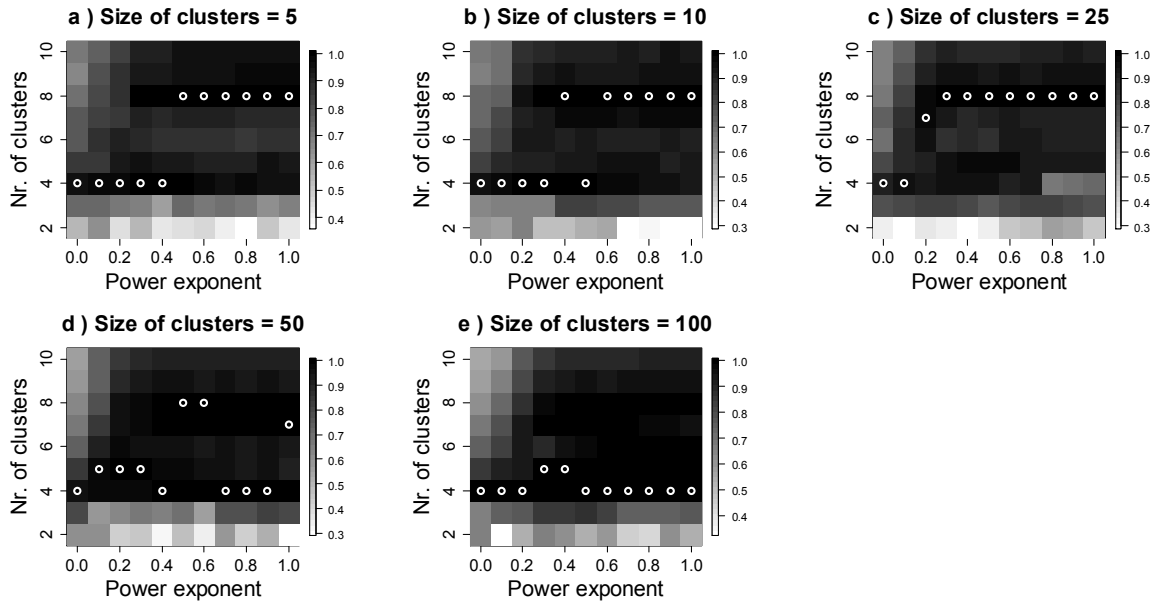
707



708

709 Fig. 5. Simulated data with different noise levels showing the heat maps of classification  
 710 stability obtained with different parameters for number of clusters and power exponent. For  
 711 the meaning of shading and other symbols see Fig. 1. The abundance-based numbers of  
 712 clusters is eight, and the species-pool-based number of clusters is four.

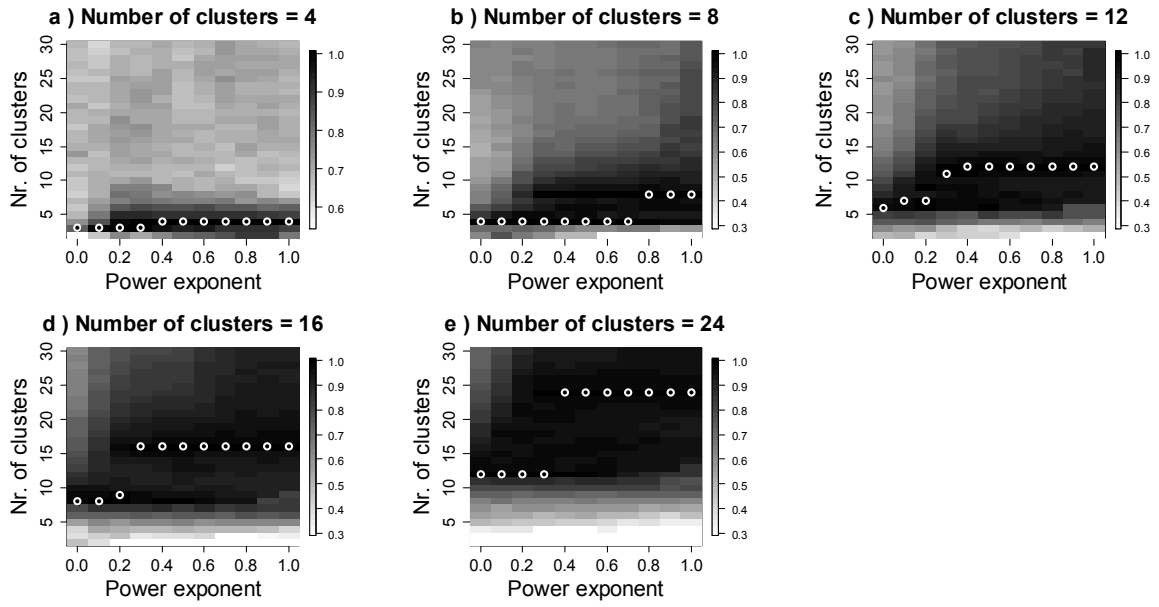
713



714

715 Fig. 6. Simulated data with different cluster sizes and fixed number of clusters showing the  
 716 heat maps of the classification stability obtained with different parameters for number of  
 717 clusters and power exponent. For the meaning of shading and other symbols see Fig. 1. The  
 718 abundance-based numbers of clusters is eight, and the species-pool-based number of clusters  
 719 is four.

720



721

722 Fig. 7. Simulated data with different numbers and fixed size of clusters showing the heat maps  
 723 of classification stability obtained with different parameters for number of clusters and power  
 724 exponents. For the meaning of shading and other symbols see Fig. 1.