5

6

7

8

9   **Joint optimization of cluster number and abundance transformation for obtaining**

10   **effective vegetation classifications**

11

12   **Attila Lengyel** [1,2,*] (lengyel.attila@okologia.mta.hu)

13   **Flavia Landucci** [3] (flavia.landucci@gmail.com)

14   **Ladislav Mucina** [4, 5] (laco.mucina@uwa.edu.au)

15   **James Tsakalos** [4] (james.tsakalos@research.uwa.edu.au)

16   **Zoltán Botta-Dukát** [1,6] (botta-dukat.zoltan@okologia.mta.hu)

17

18   [1] MTA Centre for Ecological Research, Institute of Ecology and Botany, Alkotmány u. 2-4,

19   H-2163 Vácrátót, Hungary

20   [2] Department of Vegetation Ecology, University of Wrocław, ul. Przybyszewskiego 63, 51-

21   148 Wrocław, Poland

22 [3] Department of Botany and Zoology, Masaryk University, Kotlářská 2, CZ-611 37 Brno,

23 Czech Republic

24 [4] School of Biological Sciences, The University of Western Australia, 35 Stirling Hwy,

25 Crawley WA 6009, Perth, Australia

26 [5] Department of Geography & Environmental Studies, Stellenbosch University, Private Bag

27 X1, Matieland 7602, Stellenbosch, South Africa

28 [6]MTA Centre for Ecological Research, GINOP Sustainable Ecosystems Group, Klebelsberg

29 Kuno u. 3, H-8237 Tihany, Hungary

30

31 *Corresponding author

32

## Abstract

34 **Question:** Is it possible to determine which combination of cluster number and taxon

35 abundance transformation would produce the most effective classification of vegetation data?

36 What is the effect of changing cluster number and taxon abundance weighting (applied

37 simultaneously) on the stability and biological interpretation of vegetation classifications?

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

39 **Methods**: Real data sets representing Hungarian submontane grasslands, European wetlands,

40 and Western Australian kwongan vegetation, as well as simulated data sets were used. The

41 data sets were classified using the partitioning around medoids method. We generated

42 classification solutions by gradually changing the transformation exponent applied to the

43 species projected covers and the number of clusters. The effectiveness of each classification

44 was assessed by a stability index. This index is based on bootstrap resampling of the original

45 data set with subsequent elimination of duplicates. The vegetation types delimited by the most

46 stable classification were compared with other classifications obtained at local maxima of the

47 stability values. The effect of changing the transformation power exponent on the number of

48 clusters, indexed according to their stability, was evaluated.

49  **Results:** The optimal number of clusters varied with the power exponent in all cases, both

50  with real and simulated data sets. With the real data sets, optimal cluster numbers obtained

51  with different data transformations recovered interpretable biological patterns. Using the

52  simulated data, the optima of stability values identified the simulated number of clusters

53  correctly in most cases.

54  **Conclusions:** With changing the settings of data transformation and the number of clusters,

55  classifications of different stability can be produced. Highly stable classifications can be

56  obtained from different settings for cluster number and data transformation. Despite similarly

57  high stability, such classifications may reveal contrasting biological patterns, thus suggesting

58  different interpretations. We suggest testing a wide range of available combinations to find

59  the parameters resulting in the most effective classifications.

60

## Keywords

62  Clustering; Cluster validation; Community similarity; Cover scale; Data type; Multivariate

63  data analysis; Numerical classification; Stability of classification

64

## Abreviations

66  MSL = mean standardized lambda; PAM = partitioning around medoids; PCoA = principal

67  coordinate analysis

68

## Nomenclature

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

71

## Introduction

73  Numerical methods are applied in vegetation classification studies to reduce the

74  dimensionality of the data in seeking patterns, to increase objectivity in the analyses, and thus

75  to enhance the reproducibility of results. Still, classification protocols often rely on subjective

76  decisions that can significantly influence the results (De Cáceres et al. 2015). Subjective

77  choices can hardly be avoided, yet they should be well-informed and logical to make the

78  analytical procedures reliable and repeatable. In numerical classifications, according to

79  Lengyel & Podani (2015), the choice of the number of clusters and the weight attributed to

80  abundant species relative to scarce species (hence the data transformation), are among the

81  most influential decisions that have to be considered carefully. If the aim of the classification

82  is to delimit a pre-set number of vegetation types within the data set, then the choice of the

83  resulting clusters should be guided by practical considerations. In certain cases there is

84  reasonable external information available for selecting a transformation function as well. For

85  instance, if the abundance estimations are deemed inaccurate, only presence/absence data

86  should be used. Equally, if the purpose of the study is to analyse vegetation types

87  characterised by dominant species, it is more logical to apply a transformation giving high

88  emphasis to differences in species abundance. However, if the aim of the classification is to

89  explore variation by separating and differentiating vegetation types, classifications using a

90  suite of contrasting parameters should be produced. These should be evaluated *a posteriori* in

91  order to identify the optimal parameter values yielding in the 'best' (according to the set

92  criteria) classification.

93  The optimal number of clusters can be sought for by calculating *cluster effectiveness* (or

94  *validity*) *index* for classifications with increasing number of clusters. Thus, the optimal

95  number of clusters is the one where the effectiveness index reaches maximum or minimum,

96  depending on scaling. This procedure is widely known and regularly applied in classification

97  studies (e.g. Botta-Dukát et al. 2005; Tichý et al. 2010, 2011). However, we are aware of only

98  a few examples when authors evaluated different data transformations for finding the optimal

99  weighting of abundances that would reveal biological patters most effectively or would lead

100 to the most stable results. Jensen (1978) evaluated the effect of several data transformations

101 on classifications and ordinations of a lake vegetation data, and concluded that 'extreme

102 transformations' (i.e. those giving high weight either to high abundance values or, in reverse,

103 to presence/absence data) can yield significantly different results. This finding was

104 corroborated by Campbell (1978) and van der Maarel (1979). Wilson (2012) compared the

105 stability of ordination analyses performed on various vegetation samples using different

transformations of abundance and concluded that the 'optimal' transformations depend on context, such as geographical extent, environmental heterogeneity, disturbance status of the study area, and quality of abundance estimations. Although, any 'optimal' parameterization supposed to produce a robust classification is specific for the actual data set, the low interest of researchers in finding them, or at least in assessing the performance of methods they apply, is surprising, given that vastly different results can be achieved by application of different abundance scales in multivariate analyses – a fact well known for long time (Austin & Greig-Smith 1968; Noy-Meir et al. 1975; van der Maarel 1979).

In this paper, we introduce a procedure for choosing the combination of two factors, namely (1) the number of clusters and (2) varying scale of transformation power, assisting in identification of the most effective classification outcome. Like other approaches aimed at determination of the optimal number of clusters (e.g. Aho et al. 2008), a general guideline for finding the optimal transformation would be to find the function that leads to the most stable of several possible classifications produced by differently parameterized transformation functions. We show that changing one of these two factors has an impact on the optimal values of the other, which influences the biological interpretation of the classification result, and therefore we promote their joint optimization. We test this approach using real and simulated data sets.

**Materials and methods**

*Grasslands data set*

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

*Wetlands data set*

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

*Kwongan data set*

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

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

*Simulated data*

Simulated data matrices consist of $N$ plots (in the rows) and $S$ species (in the columns). Plots belong to $K$ clusters of equal size, thus the number of plots is $N/K = n$ in each cluster, and $n$ is

165   a pre-defined integer. Ten species occur in each cluster and each species occurs in two

166   clusters, thus $S = 10 \times K/2$. Each species has constant abundance across plots within a cluster,

167   while the abundances may differ among clusters. The abundances of species within one of the

168   two clusters where they occur, are drawn from a Poisson-lognormal distribution (Bulmer

169   1974) where the mean and the standard deviation (SD) of the lognormal distribution are (2; 1)

170   on log scale. For the other cluster, the order of abundances is reversed, thus if a species was

171   the most abundant in one of the clusters where it occurs, then this species will be the least

172   abundant in the other one (considering only species occurring in this cluster). These matrices,

173   therefore, consist of plots of $K$ clusters according to raw abundances of species, but $K/2$

174   clusters according to presence/absence data because pairs of clusters share the same species

175   occurring with different abundances. We expect the optimal number of clusters to be $K/2$ with

176   low exponents, while with high exponents optimal solution should comprise $K$ clusters.

177   Notably, abundance-based clusters are nested within clusters based on presence/absence data.

178   Within each cluster, plots are identical, thus the clustered structure is initially perfect. An

179   exemplary matrix is shown in Appendix S1. Then, noise was added to this initial matrix

180   following the method of Gotelli (2000) used for 'noise test', but applied to abundances instead

181   of presence/absence data. This procedure applies a swapping algorithm to introduce noise. In

182   a single swap, the rows and columns of the original matrix are permuted, and a $2 \times 2$

183   submatrix with positive values in the diagonal is chosen randomly. Then the two diagonal

184   cells are decreased by 1, while abundances in the two off-diagonal cells are increased by 1

185   individual, thus the sum and the marginal totals of the submatrix do not change. Finally, the

186   original order of rows and columns is restored. A single swap would affect a sparse matrix

187   more than one with high fill. Also, large matrices are more 'resistant' to the same number of

188   swaps than small ones. Therefore, noise is added to the matrices in discrete levels, one level

189   consisting of as many swaps as the number of non-zero elements in the matrix. Our

190   preliminary analyses suggested that in this way a comparable amount of stochasticity can be

191   added to matrices of different size and fill.

192   Five simulation series were performed, each of them with five different set-ups. In these

193   series, one or two parameters were changed systematically in order to generate simulated

194   matrices that would differ in: i) noise level; ii) size of clusters with number of clusters fixed;

195   iii) number of clusters with cluster sizes fixed; iv) number and size of clusters with total

196   number of plots fixed; v) dominance of species. The dominance was changed by modifying

197   the SD of the lognormal distribution used as input for the Poisson process of species

abundances. When SD is high, there is one or a few highly dominant species within a plot and many very scarce species, while with lower SD species abundances should be balanced.

*Classification method*

For classifying the data sets, we used the partitioning around medoids method (PAM; Kaufman & Rousseeuw 1990) using Marczewski-Steinhaus index as the measure of dissimilarity (Appendix S2). For the Grasslands and Kwongan data set covers of species were directly estimated on percentage scale in the field, while for the Wetlands data set, abundances were mostly recorded on Braun-Blanquet or finer ordinal scales. These ordinal categories were replaced by their midpoint percentages. Cover percentages were power transformed using the function $x´ = x^a$, where $x$ is the original cover value on percentage scale, $a$ is the power exponent, and $x´$ is the transformed cover value. The power exponent was gradually changed from 0 to 1, with 21 steps by 0.05 in between in case of real data, and with steps of 0.1 in case of simulations where simpler patterns were expected. Low values of the exponent reduce the effect of differences between species abundances, thus giving more weight to rare species, while values near 1 give more weight to abundant species. The lowest number of clusters examined was 2. The highest number of examined clusters was 10 for the Grasslands data, 40 for the Wetlands and for the Kwongan data, and it varied in simulations according to the pre-defined number of clusters and sample size. The maximal number of clusters was arbitrarily determined to balance between computation time and the number of practically distinguishable vegetation types.

*Evaluation of classifications*

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

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

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

The marginal distribution of the heat map can also be examined for determining those parameter values which are likely to provide the most effective classification ourcomes, or the lowest or highest variation in classification stability. If one of the parameters, e.g. the exponent, is fixed to an actual value, the mean of the MSL values obtained with changing the other parameter, that is the number of clusters, gives how stable the classifications obtained with the actual exponent are on average. By using the SD instead of the mean, the variation of

stability can be expressed, too. Therefore, the SD is a measure of how important the decision is about one of the two parameters if the other one is fixed to an actual value. The use of marginal distributions is showed only for the Grasslands data set.

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

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

**Results**

*Grasslands data set*

The heat map (Fig. 1) showed that the MSL values varied considerably across cluster number and power exponent. With presence/absence data (a = 0), stability was the highest at the five-cluster solution. From a = 0.05 to a = 0.25, the three-cluster level was the most stable, including a = 0.15 where the second highest stability value was obtained (MSL = 0.804). Between a = 0.3 and a = 0.4, the stability peaked at two clusters, then from a = 0.45 the four-cluster solution was optimal until a = 0.90, while for the higher exponent values again three

289    clusters were shown to be the best. The absolute maximum value was found with a = 0.55 and

290    the four-cluster solution, where the stability of the classification was MSL = 0.824. Exponents

291    between a = 0.25 and 0.50 resulted in the highest stability values on average, and the SD of

292    stability was also the lowest in this interval (Fig. 2). Nevertheless, a second local optimum

293    was found at a = 0.8, although the SD was much bigger here. Across the cluster levels, the

294    three- and four-cluster solutions were the most stable on average, while stability values did

295    not vary much, except for 2 clusters where SD was the highest.


296    We used the most stable classification (i.e. four clusters and exponent 0.55; hereafter called

297    'Partition A') as the baseline for the interpretation of all clusters and classifications (Appendix

298    S4). This classification was identical with what was obtained by a = 0.50, that is, square-root

299    transformation. Clusters A1, A2, A3, and A4 are the elements of the Partition A. Cluster A1

300    represents grasslands of the alliance *Violion caninae*, but some species of the mesic meadows

301    of the order *Arrhenatheretalia* are also frequent. Cluster A2 contains plots of the

302    *Arrhenatherion*. This type was recently described as the *Diantho-Arrhenatheretum*

303    association by Lengyel et al. (2016); it represents nutrient-poor, acidic grasslands overgrown

304    by taller grasses (e.g. *Helictotrichon pubescens*, *Arrhenatherum elatius*) after abandonment or

305    changing management to mowing. Cluster A3 comprises unproductive meadows and pastures

306    dominated by *Agrostis capillatis*, *Festuca rubra*, and *Galium verum*. These stands are similar

307    in species composition to the *Anthoxantho-Agrostietum,* known also from Slovakia and the

308    Czech Republic. Cluster A3 is also intermediate between *Arrhenatheretalia* and *Violion*

309    *caninae*. Cluster A4 contains grasslands dominated by *Nardus stricta*, in which species of

310    waterlogged soils are also present. This type is traditionally also called '*Hygro-Nardetum*'

311    (e.g. Borhidi et al. 2012).


312    In the presence/absence case (a = 0), five clusters were differentiated. Hereafter, this

313    classification is called 'Partition B'. Cluster B1 included many plots of Cluster A1 and A3,

314    thus representing mesic meadows with some species of the *Violion caninae*, and matching the

315    species composition of *Anthoxantho-Agrostietum*. Cluster B2 and B3 contained mostly plots

316    previously classified to A2, thus differentiating between two subtypes of *Diantho-*

317    *Arrhenatheretum*: one with more hygrophilous, and one with more forest-steppe species,

318    respectively. Cluster B4 represents the '*Hygro-Nardetum*' type, thus is similar to Cluster A4.

319    Cluster B5 contains only two plots similar to the *Anthoxantho-Agrostietum*.

320    With a = 0.15 and three clusters a local peak was detected, to be referred to as Partition C.

321    Cluster C1 contains many plots representing the types mediating between the

322    *Arrhenatheretalia* and *Violion caninae*, formerly classified to Clusters A1 and A3. Cluster C2

323    represents the *Diantho-Arrhenatheretum,* and it is very similar to Cluster A2. Cluster C3

324    represents the *'Hygro-Nardetum'* and matches with Cluster A4.

325    With a = 1 (= no data transformation), three clusters provided the most stable resolution. This

326    classification was called Partition D. Cluster D1 represents grasslands on nutrient-poor soils,

327    including the *'Hygro-Nardetum'* and other types related to the *Violion caninae* and containing

328    *Nardus stricta*. It contains plots of Cluster A1 and A4. Cluster D2 represents mesic hay

329    meadows with *Arrhenatherum elatius*, and it shares many plots with Cluster A2. Cluster D3

330    represents unproductive meadows and pastures with the dominance of *Agrostis capillaris,*

331    *Briza media* and *Festuca rubra*. Most of its plots were assigned to Cluster A3 and C2.

332    Therefore, the Partitions C and D similarly separated the *Diantho-Arrhenatheretum* from

333    other types, but differed in how they delimited two other clusters in the rest of the data set.

334    The cross-tabulation of Partition A against Partitions B, C and D, as well as Partition C

335    against Partition D are shown in Appendix S5.

336    *Wetlands data set*

337    The optimal number of clusters ranged between 3 and 7 when the exponent ranged between 0

338    and 0.20 (Fig. 3). With higher exponents, the optimal cluster levels increased, too; from a =

339    0.35 the most stable classifications were found at levels of more than 30 clusters. In the binary

340    case (a = 0), the optimal cluster level was 6, with the square-root transformation (a = 0.5) it

341    was 30, with no transformation (a = 1) it was 39. The most stable classification was the one

342    with a = 0.80 and 40 clusters where MSL was 0.933. At this level clusters were distinguished

343    according to dominant species that were both constant and character species in many cases.

344    Using other high exponents (e.g. a = 0.50 or a = 1) resulted in very similar classifications,

345    thus only the comparison of solutions with a = 0 (hereafter called 'Partition W') and a = 0.80

346    ('Partition Z') are presented using synoptic tables (Appendix S6 and S7, respectively). Since

347    many phytosociological associations and alliances of wetland vegetation are defined by

348    dominant species, classifications with high exponents (Partition Z) showed a good

349    correspondence with low-rank syntaxa. With low exponents, the most stable classifications

350    revealed markedly different patterns that were difficult to interpret, yet these local optima

possessed much lower stability. With a = 0 (Partition W) differences in species pools offered some, although not fully satisfactory explanation for the distinction of clusters. Cluster W1 contained many plots of tall-sedge vegetation with short submerged periods and eutrophic soils (supporting mostly *Magnocaricion gracilis* vegetation). Cluster W2 included mostly plots of tall-sedge vegetation on sites with poorer nutrient supply (mostly *Magnocaricion gracilis* and *Magnocaricion elatae*). Cluster W3 is characterised, to a large part, by reed vegetation belonging to the *Phragmition* and *Phalaridion*. Clusters W4 and W5 contained many plots sampled in wetlands characterised by fluctuating shallow waters (mostly *Eleocharito-Sagittario*, *Phramition*, *Glycerio-Sparganion*), however no clear ecological difference could be recognized between them. Cluster W6 included plots from nutrient-poor mire vegetation often classified as the *Scheuchzerio-Caricetea*. Obviously, Partition W showed very low congruence with the syntaxonomical system and Parition Z (Appendix S8).

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

*Kwongan data set*

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

At a = 0, four clusters were distinguished (Partition K; Appendix S9). Cluster K1 represented a community with typical species *Hakea candolleana* and *Allocasuarina humilis* found on free-draining soils. Cluster K2 was identified as *Xylomelum angustifolium-Banskia menziesii* community thriving on sandy soils on dune swells. Cluster K3 included plots from *Ecdeiocolea monostachya-Scholtzia laxiflora* community occurring on sandy soils with slightly elevated clay content in inter-dune depressions, while Cluster K4 represented *Banksia shuttleworthiana-Cristonia biloba* confined to regolith composed of depositional lateritic scree and sand. Therefore, these clusters represented an edaphic gradient spanning Cluster K2 (deep sandy soils from the sand dune swells) and Cluster K3 (depressions showing elevated

clay content), with Clusters K1 and K4 occupying intermediate position along the gradient. At a = 0.95, the 2-cluster solution was the most stable one (Partition L; Appendix S10). The cross-tabulation tables (Appendix S11) showed that all plots of the Cluster K3 were assigned to the Cluster L1 - the only cluster whose plots were assigned to the same cluster in Partitions K and L. The Cluster K1 was concentrated in Cluster L1, while most plots of the Clusters K2 and K4 belonged to L2. Partitions K and L similarly recovered the gradient between vegetation types supported by soils having elevated clay content (represented by Clusters K1 & K3, as well as L1) and sandy soils (as Clusters K2 & K4, and L2) on the basis of characteristic species of the clusters. The relative position of the clusters in a PCoA ordination also supports the notion that the main compositional patterns are similarly revealed by different abundance weighting (Appendix S12).

*Simulations*

At the noise level 1, where abundances were strongly down-weighted (a = 0 or a = 0.1), the stability was highest at the pre-defined number of four species-pool based clusters (Fig. 5). From a = 0.2 to a = 0.7, two peaks were found, namely at the 4- and 8-cluster levels, the latter being of higher stability, and with one intermediate peak at a = 0.3 and seven clusters. Where abundance differences were not or only slightly reduced (a > 0.7), only the 8-cluster peak was obvious. From the noise level 2 and higher, the stability peaked at the 8-cluster level. As more levels of noise were added, classifications with low exponent were becoming less and less stable.

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

When the number of clusters increased from four with constant cluster sizes, the typical pattern of lower optima at low exponents and higher optima at high exponents were found in most cases, yet with some exceptions (Fig. 7). Where the species-pool based cluster number

was two and the abundance-based cluster number was four, three clusters were the most stable with low exponent and four with high exponent. With higher number of true clusters, the most stable classification identified the pre-defined cluster numbers correctly: 8, 12, 16, and 24 clusters with higher exponents, and 4, 6, 8, and 12 clusters with lowers exponents, respectively. The point of inflection, when the observed optima shifted from the species-pool-based level to the abundance-based level, was variable. Yet a broad interval with at least two local peaks of stability was detectable in all heat maps at intermediate exponent values. Cluster numbers between the species-pool-based and the abundance-based optima also came out as optimal in some cases, especially with exponents near the inflection value.

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

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

**Discussion**

*Evaluation of the real data*

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

443 The Wetlands data set showed that the optimal cluster level can markedly differ if different

444 data transformations are used. While presence/absence data yielded six stable clusters that

445 represented types with more or less different species pools, accounting for differences in

446 abundances raised the optimal levels over 30, where each cluster is separated according to the

447 dominant species. The fact that the high number of stable clusters obtained using high

448 exponent were not nested within the few stable clusters based on presence/absence data, is a

449 clear indication that different data transformations can reveal different types of biological

450 patterns. With low exponents, classifications were best explained by patterns generated by

451 habitat-specific species-pools, while with high exponents, community types differing in fine-

452 scale environmental variation, temporal variability and site history were revealed. It is of

453 interest, that in our study, 40 clusters was the finest classification level examined due to a

454 compromise between practical and scientific reasons, but in reality the optimal number of

455 clusters in the Wetlands data set could have been even higher.

456 The Kwongan data provided a special insight into the interaction of data transformation and

457 cluster number. Changing the exponent changed the optimal number of clusters as well, and

458 the resulting stable classifications were moderately congruent. However, even these,

459 seemingly less similar classifications revealed the most important ecological pattern on the

460 basis of faithful species — the soil gradient, although fine patterns of transitional subtypes

461 between the extremes were not detected equally well. The Kwongan data set, due to its high

462 beta diversity and balanced within-plot abundance distribution, was less sensitive to changes

463 in data transformation and cluster number in terms of biological interpretation, even though

464 the assignment of plots showed some variation.

465 *Lessons from the simulations*

466 In the simulations, we generated data structure with contrasting patterns with respect to

467 occurrence information. If abundance information were emphasized, the true number of

468 clusters (vegetation types) was twice as high as in cases where only presence/absence data

469 were considered, hence we differentiated a 'species-pool-based' and an 'abundance-based'

470 number of clusters. In reality, however, also an opposite can be observed, where a few species

471 can be dominant in habitats with different species pools. In such a case the number of

472 abundance-based clusters could be lower than those based on species-pools, as it was seen

473 with the Kwongan data set.

474    We expected that *weak* data transformations (the exponent being close to 1) which preserve

475    the differences in original abundance patterns, would yield a higher cluster number, while

476    *strong* transformations (the exponent approaching 0) which significantly reduce abundance

477    differences would find the half of this number of clusters optimal. Our results confirmed this

478    expectation.

479    We introduced stochasticity to artificial data using a similar method as that by Gotelli (2000)

480    called 'noise test'. This type of noise made classifications with stronger transformations less

481    stable than those involving weak transformations. This result can be understood by recalling

482    how we generated species abundances and noise. The species abundances had been drawn

483    from a Poisson-lognormal distribution, which resulted in many scarce and few abundant

484    species. Considering that the artificial matrices are designed in a way that their matrix fill is

485    low, swapping individuals can moderately reduce the abundance of species in a plot, or it can

486    slightly increase less abundant species, or make absent species present with low abundance.

487    However, it is unlikely to make an abundant species absent in a plot, or to make an absent

488    species very abundant. As a result, the applied noise affected binary information more than

489    the proportions of abundances which determine classifications involving weak data

490    transformations. We believe that this type of noise simulates a common form of stochasticity

491    in nature that is caused by random death of individuals followed by random colonization.

492    The simulations have revealed several tendencies in classification stability as related to cluster

493    number, data transformation, and sample properties. With increasing size of clusters, the

494    number of abundance-based clusters was underestimated, while the number of clusters based

495    on species pools was detected correctly. Despite this observation with both fixed and

496    changing total sample size, we cannot offer a clear explanation for this finding.

497    Based on the tests with modified pre-defined number of clusters with fixed cluster sizes, the

498    stability as optimality criterion seems to track the changes correctly in most cases. However,

499    when the number of clusters based on presence/absence data was two, the most stable

500    classifications were obtained at the three-cluster level with strong transformation. (With weak

501    transformations, the abundance-based number of clusters was correctly found at the level of

502    four clusters.) Moreover, in a few cases, optima were indicated between the species-pool-

503    based and the abundance-based levels. When the total sample size was fixed, but number and

504    size of clusters changed, stability performed similarly well. Some inconsistency was found at

505    four abundance-based clusters, where the most stable level was found at two clusters for all

506  but one value of the exponent. Surprisingly, the exception was the binary case (a = 0) where

507  all classifications were generally less stable and the optimum was at the pre-defined number

508  of clusters based on abundance, i.e. four clusters. This contradicts our expectation and we

509  have no clear explanation for this. Despite the above mentioned spurious exceptions, the

510  stability seemed rather robust and accurate across a wide range of cluster numbers with PAM.

511  In real situations, mapping a goodness of classification measure as a function of data

512  transformation and cluster number would help avoiding less effective parameter

513  combinations.

514  Testing the effect of community dominance on stability by changing the logarithm of SD of

515  species abundances revealed that at the lowest dominance (i.e. low SD), the number of

516  clusters based on species pool was optimal regardless of data transformation. As dominance

517  increased, abundance-based cluster number became more stable and was identified as optimal.

518  This is in line with the common experience that in monodominant vegetation types (e.g.

519  aquatic and marsh vegetation) classifications based on abundance data are more effective and

520  can markedly differ from presence/absence-based classifications, while when the species

521  abundances are more balanced, accounting for abundance differences does not give

522  significantly different or more effective classification than what is obtained by species

523  composition.

524  *Concluding remarks*

525  Classification stability depends both on cluster number and data transformation. The trend of

526  stability along increasing power exponent varies across cluster numbers, and vice versa, the

527  number of clusters resulting in the most stable classifications depends on data transformation.

528  Slight changes in any of these two factors may change the stability of a classification, hence

529  different biological conclusions can be reached. At the same time, similarly effective

530  classifications can be produced using different combinations of parameters. Finding such

531  local optima contributes to the thorough understanding of biological patterns in the sample.

532  Stability, as proposed by Tichý et al. (2011), is a standardized measure of classification

533  effectiveness because every single classification is compared to classifications of its without-

534  replacement bootstrap subsamples obtained with exactly the same methods. We have chosen

535  this index in our study because of this advantage. However, there are many other measures of

536  effectiveness, but we have chosen not to evaluate them experimentally in this paper. For

answering specific research questions, other indices may be more appropriate than stability. In such cases the workflow of testing the effect of data transformation and cluster number on classification effectiveness, and the visualization of results should be the same as we presented, only the measure of effectiveness should be replaced by an alternative. Moreover, it is also possible to perform the optimization analysis using several different effectiveness measures, and then combine the results in order to identify the classification which is the most effective on average across the applied indices.

Apart from the cluster number and the power exponent, we see no obstacles to test the effect of other types of methodological decisions using our approach. For example, an effectiveness measure might be calculated for classifications obtained by different values for the β parameter of the flexible clustering method by Lance & Williams (1967), and the β value providing the most stable classification might be determined. Moreover, our optimization approach can easily be adapted to ordinations, too. If the cluster effectiveness index applied here is substituted by a measure of stability of ordinations (as done by Wilson 2012), the effect of data transformation on the stability of ordinations can be evaluated systematically. The extension of the optimization procedure presented here beyond data transformation and cluster number is a future direction of our research.

**Authors contributions**

A.L. outlined the main idea, performed data analysis and wrote the initial manuscript, Z.B.D. contributed with discussion in all stages of the work, F.L. helped in preparation of the

566     Wetlands data set and the evaluation of the analysis, L.M. and J.T. contributed by providing

567     the Kwongan data set and evaluating the results, L.M. and J.T. performed linguistic revisions

568     of early versions of the text. All authors critically commented on the manuscript and the

569     supplementary materials.

570

571     **References**

572     Aho, K., Roberts, D.W. & Weaver, T. 2008. Using geometric and non-geometric internal

573     evaluators to compare eight vegetation classification methods. *Journal of Vegetation Science*

574     19: 549–562.

575     Austin, M.P. & Greig-Smith, P. 1968. The application of quantitative methods to vegetation

576     survey: II. Some methodological problems of data from rain forest. *Journal of Ecology* 56:

577     827–844.

578     Borhidi, A., Kevey, B. & Lendvai, G. 2012. *Plant communities of Hungary.* Akadémiai

579     Kiadó, Budapest, HU.

580     Botta-Dukát, Z., Chytrý, M., Hájková, P. & Havlová, M. 2005. Vegetation of lowland wet

581     meadows along a climatic continentality gradient in Central Europe. *Preslia* 77: 89–111.

582     Bulmer, M.G. 1974. On fitting the Poisson lognormal distribution to species-abundance data.

583     *Biometrics* 30: 101–110.

584     Campbell, B.M. 1978. Similarity coefficients for classifying plots. *Vegetatio* 37: 101–108.

585     Chytrý, M., Tichý, L., Holt, J. & Botta-Dukát, Z. 2002. Determination of diagnostic species

586     with statistical fidelity measures. *Journal of Vegetation Science* 13: 79–90.

587     Chiang, M. & Mirkin, B. 2010. Intelligent choice of the number of clusters in k-means

588     clustering: An experimental study with different cluster spreads. *Journal of Classification* 27:

589     3–40.

590    De Cáceres, M., Chytrý, M., Agrillo, E., Attorre, F., Botta-Dukát, Z., Capelo, J., Czúcz, B.,

591    Dengler, J., Ewald, J., (…) & Wiser, S.K. 2015. A comparative framework for broad-scale

592    plot-based vegetation classification. *Applied Vegetation Science* 18: 543–560.

593    Goodman, L. & Kruskal, W. 1954. Measures of association for cross classifications. *Journal*

594    *of the American Statistical Association* 49: 732–764.

595    Gotelli, N.J. 2000. Null model analysis of species co-occurrence patterns. *Ecology* 81: 2606–

596    2621.

597    Hennig, C. 2007. Cluster-wise assessment of cluster stability. *Computational Statistics &*

598    *Data Analysis* 52: 258–271.

599    Hill, M.O. 1973. Diversity and evenness: a unifying notation and its consequences. *Ecology*

600    54: 427–432.

601    Hunter, J.C. & McCoy, R.A. 2004. Applying randomization tests to cluster analyses. *Journal*

602    *of Vegetation Science* 15: 135–138.

603    Jensen, S. 1978. Influences of transformation of cover values on classification and ordination

604    of lake vegetation. *Vegetatio* 37: 19–31.

605    Kaufman, L. & Rousseeuw, P.J. 1990. *Finding groups in data: An introduction to cluster*

606    *analysis.* John Wiley & Sons, New York, US.

607    Király, G. (ed.) 2009. *New Hungarian Herbal. The vascular plants of Hungary.* Identification

608    key. Aggteleki Nemzeti Park Igazgatóság, Jósvafő, HU. (in Hungarian)

609    Lance, G.N. & Williams, W.T. 1967. A general theory of classificatory sorting strategies. I.

610    Hierarchical systems. *Computer Journal* 9: 373–380.

611    Landucci, F., Řezníčková, M., Šumberová, K., Chytrý, M., Aunina L., Biţă-Nicolae, C.,

612    Bobrov, A., Borsukevych, L., Brisse, H., (…) & Willner W. 2015. WetVegEurope: a database

613    of aquatic and wetland vegetation of Europe. *Phytocoenologia* 45: 187–194.

614    Lambers, H. (ed.) 2014. *Plant life on the sandplains in Southwest Australia: A global*

615    *biodiversity hotspot.* UWA Publishing, Crawley, AU.

616    Lengyel, A., Chytrý, M. & Tichý, L. 2011. Heterogeneity-constrained random resampling of
617    phytosociological databases. *Journal of Vegetation Science* 22: 175–183.

618    Lengyel, A. & Podani, J. 2015. Assessing the relative importance of methodological decisions
619    in classifications of vegetation data. *Journal of Vegetation Science* 26: 804–815.

620    Lengyel, A., Illyés, E., Bauer, N., Csiky, J., Király, G., Purger, D. & Botta-Dukát, Z. 2016.
621    Classification and syntaxonomical revision of mesic and semi-dry grasslands in Hungary.
622    *Preslia* 88: 201–228.

623    Lötter, M.C., Mucina, L. & Witkowski, E. 2013.  The classification conundrum: species
624    fidelity as leading criterion in search of a rigorous method to classify a complex forest data
625    set. *Community Ecology* 14: 121–132.

626    Lyons, M.B., Keith, D.A., Warton, D.I., Somerville, M. & Kingsford, R.T. 2016. Model-
627    based assessment of ecological community classifications. *Journal of Vegetation Science* 27:
628    704–715.

629    McIntyre, R.M. & Blashfield, R.K. 1980. A nearest-centroid technique for evaluating the
630    minimum-variance clustering procedure. *Multivariate Behavioral Research* 15: 225–238.

631    Milligan, G.W. & Cooper, M.C. 1985. An examination of procedures for determining the
632    number of clusters in a data set. *Psychometrika* 50: 159–179.

633    Mucina, L., Bültmann, H., Dierßen, K., Theurillat, J.-P., Raus, T., Čarni, A., Šumberová, K.,
634    Willner, W., Dengler, J., (…) & Tichý, L. 2016. Vegetation of Europe: hierarchical floristic
635    classification system of vascular plant, bryophyte, lichen, and algal communities. *Applied*
636    *Vegetation Science* 19: 3–264.

637    Noy-Meir, I., Walker, D. & Williams, W.T. 1975. Data transformations in ecological
638    ordination: II. On the meaning of data standardization. *Journal of Ecology* 63: 779–800.

639    Podani, J. 2000. *Introduction to the exploration of multivariate biological data.* Backhuys,
640    Leiden, NL.

641    Podani, J. & Feoli, E. 1991. A general strategy for the simultaneous classification of variables
642    and objects in ecological data tables. *Journal Vegetation Science* 2: 435–444.

643    Popma, J., Mucina, L., van Tongeren, O. & van der Maarel, E. 1983. On the determinants of

644    optimal levels in phytosociological classification. *Vegetatio* 52: 65–75.

645    Roberts, D.W. 2015. Vegetation classification by two new iterative reallocation optimization

646    algorithms. *Plant Ecology* 216: 741–758.

647    Rohlf, F.J. 1974. Methods of comparing classifications. *Annual Review of Ecology &*

648    *Systematics* 5: 101–113.

649    Rozbrojová, Z., Hájek, M. & Hájek, O. 2010. Vegetation diversity of mesic meadows and

650    pastures in the West Carpathians. *Preslia* 82: 307–332.

651    Tichý, L. 2002. JUICE, software for vegetation classification. *Journal of Vegetation Science*

652    13: 451–453.

653    Tichý, L., Chytrý, M. & Šmarda, P. 2011. Evaluating the stability of the classification of

654    community data. *Ecography* 34: 807–813.

655    Tichý, L., Chytrý, M., Hájek, M., Talbot, S.S. & Botta-Dukát, Z. 2010. OptimClass: Using

656    species-to-cluster fidelity to determine the optimal partition in classification of ecological

657    communities. *Journal of Vegetation Science* 21: 287–299.

658    van der Maarel, E. 1979. Transformation of cover-abundance values in phytosociology and its

659    effects on community similarity. *Vegetatio* 39: 97–114.

660    Vendramin, L., Campello, R.J.G.B. & Hruschka, E.R. 2010. Relative clustering validity

661    criteria: A comparative overview. *Statistical Analysis & Data Mining* 3: 209–235.

662    Willner, W., Tichý, L. & Chytrý, M. 2009. Effects of different fidelity measures and contexts

663    on the determination of diagnostic species. *Journal of Vegetation Science* 20: 130–137.

664    Wilson, J.B. 2012. Species presence/absence sometimes represents a plant community as well

665    as species abundances do, or better. *Journal of Vegetation Science* 23: 1013–1023.

666    Wiser, S.K. & De Cáceres, M. 2013. Updating vegetation classifications: an example with

667    New Zealand's woody vegetation. *Journal of Vegetation Science* 24: 80–93.

668

**List of Appendices**

Appendix S1: Simulation data example

Appendix S2: Mathematical formulae

Appendix S3: R scripts

Appendix S4: Grasslands synoptic table (Partition A)

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

Appendix S6: Wetlands synoptic table (Partition W)

Appendix S7: Wetlands synoptic table (Partition Z)

Appendix S8: Wetlands cross-tabulations

Appendix S9: Kwongan synoptic table (Partition K)

Appendix S10: Kwongan synoptic table (Partition L)

Appendix S11: Kwongan cross-tabulation

Appendix S12: Kwongan ordination

Appendix S13: Additional heat maps of the simulated data sets

683

684

685 Tables

686

687 **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 (m2) | 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)

688 **mean of diversity of order 1 divided by diversity of order 0, the latter being species
689 richness

690

Fig. 1. Analysis of the Grasslands data set showing the heat map of classification stability obtained using different parameters for number of clusters and power exponent. Darkness of the segments correlate with the value of the mean standardized Goodman & Kruskal's lambda (MSL), where the darkest segments marking the combinations of parameters leading to the most stable classifications. White circles with black dots indicate the optimal number of clusters for a given exponent.
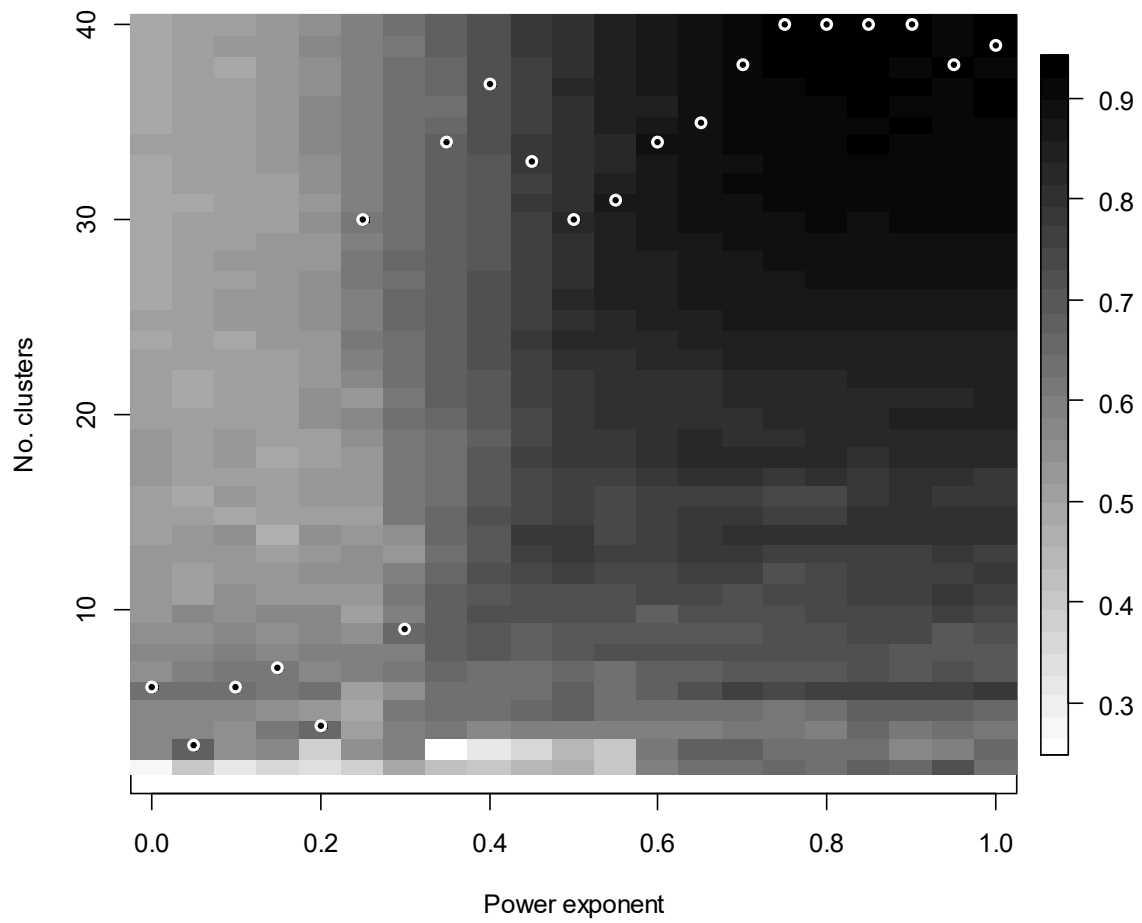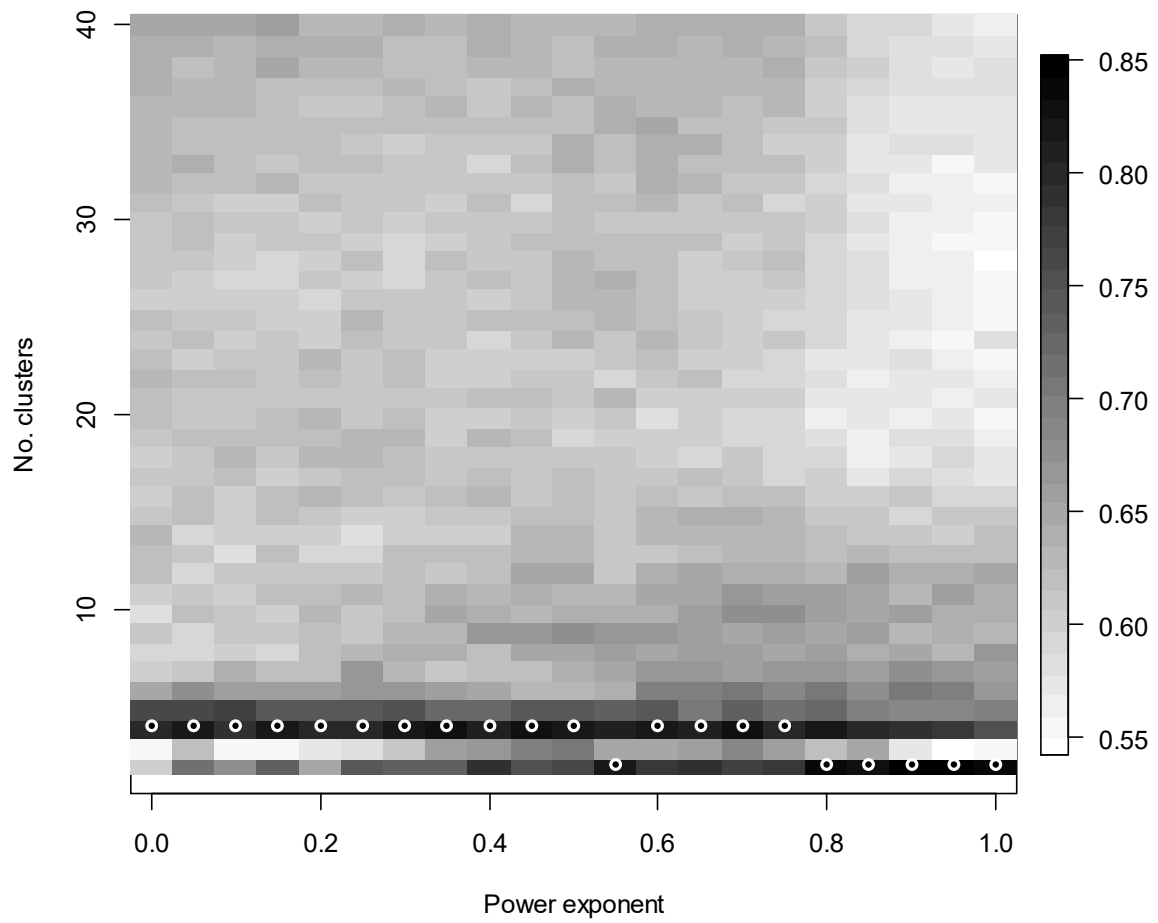
700

701  Fig. 2. Mean and standard deviation as error bars of the marginal of the heat map of the

702  Grasslands data set.

703

704

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

708

709

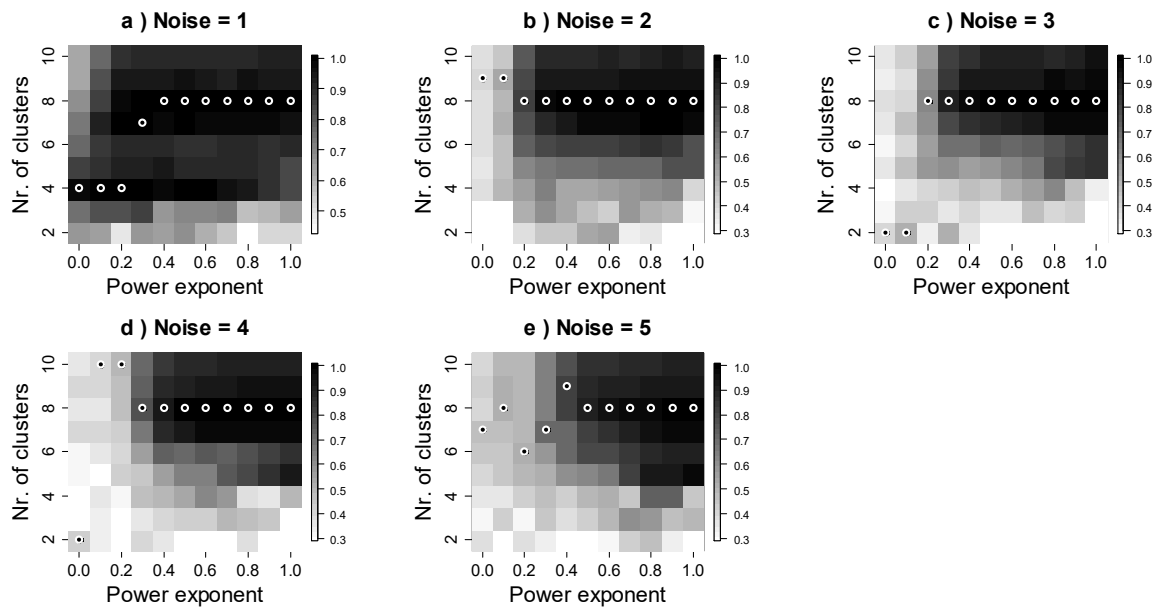Fig. 4. Analysis of the Kwongan data set showing the heat map of the classification stability

obtained using different parameters for number of clusters and power exponent. For the
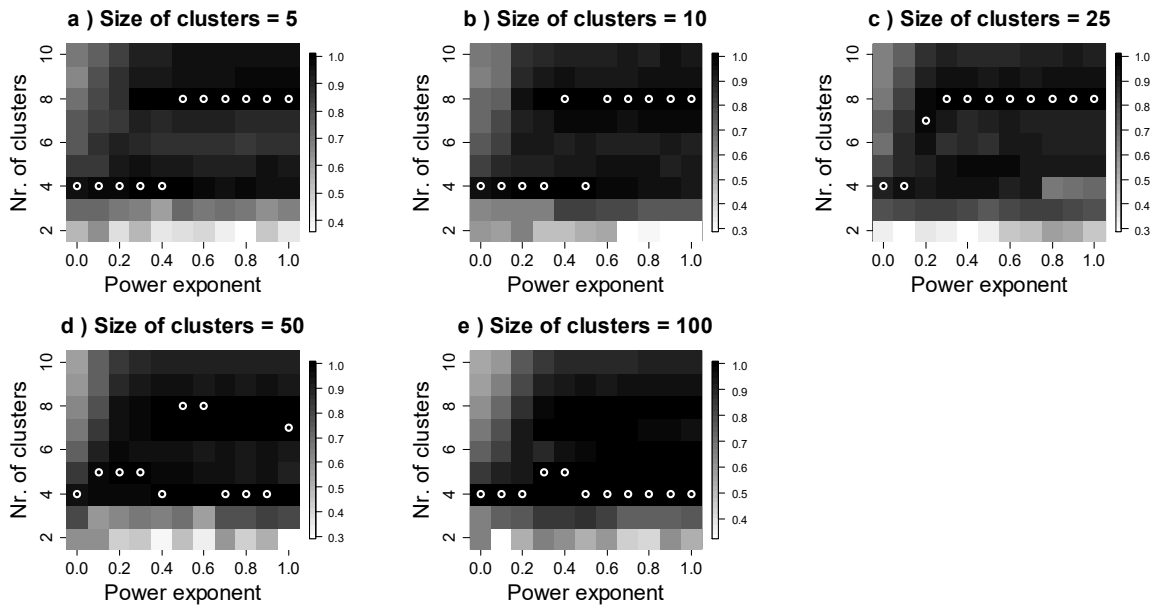
meaning of shading and other symbols see Fig. 1.

713

714

715

Fig. 5. Simulated data with different noise levels showing the heat maps of classification stability obtained with different parameters for number of clusters and power exponent. For the meaning of shading and other symbols see Fig. 1. The abundance-based numbers of clusters is eight, and the species-pool-based number of clusters is four.
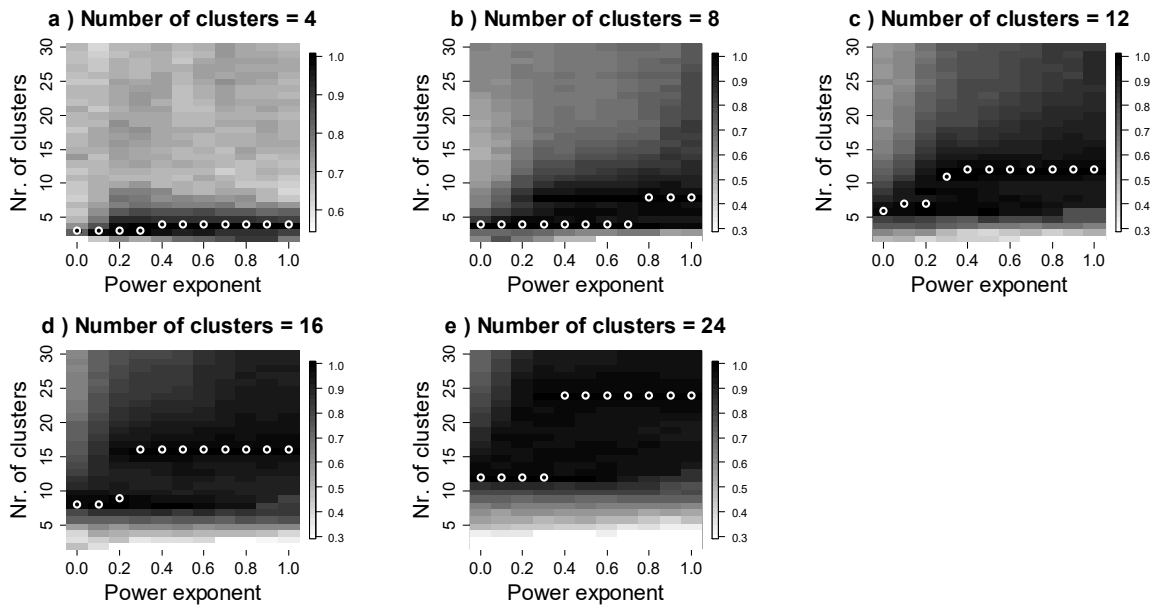
722

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

723
724
725
726
727

728

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