

1 **Title**

2 Comparison of silhouette-based reallocation methods for vegetation classification

3 **Author names**

4 Attila Lengyel¹, David W. Roberts^{2,3} & Zoltán Botta-Dukát¹

5 **Author addresses**

6 Lengyel, A. (lengyel.attila@okologia.mta.hu, corresponding author)¹, Roberts, D.W.
7 (droberts@montana.edu)², Botta-Dukát, Z. (botta-dukat.zoltan@okologia.mta.hu)¹

8 ¹ Centre for Ecological Research, Institute of Ecology and Botany, Alkotmány u. 2-4. H-2163
9 Vácrátót, Hungary

10 ² Swiss Federal Research Institute WSL, CH-8903 Birmensdorf, Switzerland

11 ³ Ecology Department, Montana State University, Bozeman, Montana, USA, 59717-3460

12 **ORCIDs**

13 Lengyel, A.: <https://orcid.org/0000-0002-1712-6748>

14 Roberts, D.W.: <https://orcid.org/0000-0001-7128-6243>

15 Botta-Dukát, Z.: <https://orcid.org/0000-0002-9544-3474>

16

17 **Author contributions**

18 A.L. raised the idea, wrote the scripts, did data analysis, lead writing; D.W.R. did data
19 analysis, discussed results, contributed to the manuscript; Z.B.D. discussed results,
20 commented on the manuscript.

21

22 **Data availability**

23 The simulated and the Grassland data sets are available in the attachment. Bryce and
24 Shoshone data sets are available through the labdsv and optpart R packages respectively.

25

26 **Acknowledgements**

27 The work of A.L. was supported by the National Research, Development and Innovation
28 Office, Hungary (project number PD123997).

29

30 **Abstract**

31 *Aims:* To introduce REMOS, a new iterative reallocation method (with two variants) for
32 vegetation classification, and to compare its performance with OPTSIL. We test (1) how
33 effectively REMOS and OPTSIL maximize mean silhouette width and minimize the number
34 of negative silhouette widths when run on classifications with different structure; (2) how
35 these three methods differ in runtime with different sample sizes; and (3) if classifications by
36 the three reallocation methods differ in the number of diagnostic species, a surrogate for
37 interpretability.

38 *Study area:* Simulation; example data sets from grasslands in Hungary and forests in
39 Wyoming and Utah, USA.

40 *Methods:* We classified random subsets of simulated data with the flexible-beta algorithm for
41 different values of beta. These classifications were subsequently optimized by REMOS and
42 OPTSIL and compared for mean silhouette widths and proportion of negative silhouette
43 widths. Then, we classified three vegetation data sets of different sizes from two to ten
44 clusters, optimized them with the reallocation methods, and compared their runtimes, mean
45 silhouette widths, numbers of negative silhouette widths, and the number of diagnostic
46 species.

47 *Results:* In terms of mean silhouette width, OPTSIL performed the best when the initial
48 classifications already had high mean silhouette width. REMOS algorithms had slightly lower
49 mean silhouette width than what was maximally achievable with OPTSIL but their efficiency
50 was consistent across different initial classifications; thus REMOS was significantly superior
51 to OPTSIL when the initial classification had low mean silhouette width. REMOS resulted in
52 zero or a negligible number of negative silhouette widths across all classifications. OPTSIL
53 performed similarly when the initial classification was effective but could not reach as low
54 proportion of misclassified objects when the initial classification was inefficient. REMOS
55 algorithms were typically more than an order of magnitude faster to calculate than OPTSIL.
56 There was no clear difference between REMOS and OPTSIL in the number of diagnostic
57 species.

58 *Conclusions:* REMOS algorithms may be preferable to OPTSIL when (1) the primary
59 objective is to reduce or eliminate negative silhouette widths in a classification, (2) the initial
60 classification has low mean silhouette width, or (3) when the time efficiency of the algorithm
61 is important because of the size of the data set or the high number of clusters.

62

63 **Keywords**

64 Flexible-beta; classification; clustering; iterative; OPTIMCLASS; optimization; OPTSIL;
65 REMOS; silhouette; validation

66

67 **Abbreviations**

68 MSW = mean silhouette width; MR = misclassification rate

69

70 **Introduction**

71 Numerical classification methods are essential data analytical tools in vegetation ecology and
72 several other scientific fields, including genomics, psychology, or sociology. Basically,
73 classification algorithms can be divided into two groups. Hierarchical algorithms produce a
74 perfectly nested hierarchy of clusters of objects, while the output of non-hierarchical methods
75 is a partition in which each classified object is assigned exclusively to one cluster (or, in the
76 special case of fuzzy clustering methods, non-exclusively to several clusters using fuzzy
77 membership weights) at the same level. Hierarchical methods can be subdivided into
78 agglomerative and divisive methods based on whether they initiate the clustering algorithm
79 from treating each single object as a separate cluster, and then merge them until all objects are
80 included in a single cluster at the highest hierarchical level, or they proceed in the opposite
81 direction by dividing the entire sample iteratively into smaller and smaller subsets in a nested
82 way. The diversity of numerical classification methods is reviewed by several authors, e.g.
83 Kaufman & Rousseeuw (1990), Podani (2000), Peet & Roberts (2013), Legendre & Legendre
84 (2012).

85 The advantage of hierarchical methods is that they do not need a pre-defined cluster number;
86 however, if a single-level classification is the objective, as is generally the case, a hierarchical
87 classification requires a post-hoc assessment for choosing the ‘best’ number of clusters.
88 Moreover, a disadvantage of hierarchical methods is that earlier steps (either merging or
89 division) constrain further ones, hence the final solution may be suboptimal. In such a case the
90 *a posteriori* reallocation of misclassified objects might be necessary.

91 Recently Roberts (2015) introduced two reallocation-based methods which can be used for
92 improving already existing classifications by optimizing a pre-selected goodness-of-clustering
93 criterion. One of these two, called OPTSIL, optimizes the silhouette width which is a widely
94 used index for evaluating classifications and identifying ‘core’ and misclassified objects
95 individually (Rousseeuw 1987, Kaufman & Rousseeuw 1997). Let i be a focal object
96 belonging to cluster A . Let C be a cluster not containing i . $a(i)$ is defined as the average
97 dissimilarity between i and all other objects in A , while $c(i, C)$ is the average dissimilarity
98 between i and all objects in C .

$$b(i) = \min_{C \neq A} c(i, C)$$

99 That is, $b(i)$ is the average dissimilarity between i and the members of its closest neighbour
100 cluster. The silhouette width, $S(i)$, is defined as:

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

101 $S(i)$ ranges between -1 and +1. Values near +1 indicate that object i is much closer to other
102 objects in its assigned cluster than to objects of the closest other cluster, implying a correct
103 classification. If $S(i)$ is near 0, the correct classification of the focal object is doubtful, thus
104 suggesting intermediate position between two clusters. $S(i)$ values < 0 indicate poor fit, and
105 such objects are often considered ‘misclassified’ (Rousseeuw 1987). In each iteration,
106 OPTSIL evaluates how much the reallocation of any single object in the classification
107 increases the sample-wise mean of silhouette width. It is done by re-assigning each object
108 from its current cluster to every other cluster, and then re-calculating the silhouette widths for
109 all objects. The reallocation which causes the highest increase in the sample-wise mean
110 silhouette is accepted in each step, until no further improvement is possible. Roberts (2015)
111 concluded that OPTSIL is able to significantly improve the initial classification; however, it is
112 slow to converge, and thus recommended for ‘polishing’ of classifications made by other
113 methods.

114 We present two new silhouette-based reallocation algorithms, called REMOS (reallocation of
115 Misclassified Objects based on Silhouette width). Using artificial and real data sets, we
116 compare them with OPTSIL in terms of three criteria: optimization success, time efficiency,
117 and interpretability.

118

119 **Materials and Methods**

120 *The REMOS algorithms*

121 Instead of evaluating the effect of the reallocation of each object (typically sample unit) on the
122 mean silhouette width, REMOS algorithms simply reallocate one or all of the objects which
123 have negative silhouette width. According to how objects to reallocate are selected, we
124 introduce two versions of REMOS. REMOS1 reallocates only the object with the most
125 negative silhouette width (i.e., the ‘worst classified’ object), while REMOS2 reallocates all
126 objects with negative silhouette width (i.e., all misclassified objects). Both algorithms stop if
127 the lowest silhouette width reaches a threshold L , or if no further improvement is possible. By
128 default L is 0; however, using different values between -1 and 0 can control tolerance towards
129 misclassifications. The steps of the algorithms are presented below:

130 (1) Calculating the silhouette widths, $S(i)$, for the classified objects;

131 (2) Are there any objects with $S(i) < L$?

132 2a. If no, then go to (5)

133 2b. If yes, go to (3)

134 (3) Updating the classification by reallocating objects:

135 REMOS1: reallocate only the object with the most negative silhouette width to its
136 neighbour cluster;

137 REMOS2: reallocate all the objects with $S(i) < L$ to their respective neighbour clusters;

138 (4) Go to (1).

139 (5) End – no further optimization is possible

140 Our preliminary runs showed that both REMOS algorithms frequently converge into loops
141 where the iteration proceeds repeatedly over a finite number of suboptimal solutions without
142 finding any of them as a final solution. To break such a loop, the algorithm checks for
143 repetitions and stops if two identical solutions occur. In this case the solution with the lowest
144 number of negative silhouette widths is selected from the previous iterations. In case of tied
145 minimum of negative silhouette widths, the solution giving the higher absolute sum of
146 negative silhouette widths (a surrogate for smaller ‘classification error’) is chosen as final.
147 Not surprisingly, in most cases REMOS1 requires many more iterations than REMOS2.
148 According to our pilot analyses with differently sized data matrices and different initial
149 classifications, this can extend the computation time of REMOS1 in comparison with
150 REMOS2. It is possible to set an upper limit to the number of iterations; however, as there is
151 no standard value for this threshold, the default setting is infinity (that is, no limit).

152 An R script of the REMOS algorithms is provided in the Electronic Supplement.

153

154 *Data sets*

155 We compared the performance of the REMOS1, REMOS2 and the OPTSIL algorithms on
156 three real and one artificial data set. The Shoshone data set is a random subset comprising 150
157 plots selected from a larger forest inventory database. This data set represents coniferous
158 forests of Shoshone National Forests (WY, USA). In the plots vascular species were recorded
159 using an ordinal scale. The Bryce data set was sampled in the Bryce Canyon National Park
160 (UT, USA; Roberts 1992). It includes 160 circular plots of $\sim 404.7 \text{ m}^2$ (0.1 acre) where the
161 cover of 169 vascular species (except trees) were recorded on ordinal scale. The Grasslands
162 data set is a subset of a larger sample of mesic grasslands of northern Hungary (Lengyel et al.
163 2016). The size of the matrix is 55 plots by 269 species. Abundances are coded on a
164 percentage scale. As artificial data, we employed a simulated data set of 400 points in two
165 dimensions. The points are aggregated into eight fuzzy clusters (Fig. 1). For different test
166 scenarios, random subsets of different size were used.

167

168 *Data analysis*

169 The performance of the REMOS and OPTSIL algorithms was evaluated from three aspects:
170 optimization success on different initial classifications of artificial and real data, dependence
171 of computation time on sample size with artificial data, and interpretability of the optimized
172 classification of real data based on indicator species.

173 For testing optimization success, initial classifications of random subsamples of the artificial
174 data set containing 200 points were prepared using the flexible-beta classification algorithm
175 (Lance & Williams 1966). This method uses a parameter called beta which enables producing
176 classifications with different sensitivity of ‘chaining’ vs. ‘grouping’ effect. The beta is
177 adjustable between -1 and +1. With lower values the grouping effect is emphasized, while
178 higher beta gives more weight to chaining. With beta = -1 flexible-beta clustering is identical
179 with the complete linkage method, with beta = 0 it agrees with the average linkage
180 (UPGMA), with beta = +1 it is the same as single linkage. Several authors reported that the
181 flexible clustering method provides the most satisfactory classifications using beta = -0.25. In
182 this analysis, values of beta were changed between -1 and +1 in steps by 0.25 in between. The
183 hierarchical classifications were cut at the 8-cluster level. The procedure was repeated 5 times
184 resulting in $5 \times 9 = 45$ initial classifications. Each of them was optimized using the REMOS1,
185 REMOS2, and OPTSIL algorithms. We compared the change of mean silhouette widths
186 (MSW) and misclassification rate (that is, the proportion of negative silhouette widths; MR)
187 across beta values between the optimized classifications and the initial classification. In the
188 Electronic Supplement we show some exemplary classifications.

189 For comparing time efficiency, we drew subsamples containing 50, 100, 200, and 300 points
190 of the artificial data set in 20 repeats, and additionally, used also the entire sample of 400
191 points. Each of them were classified to 8 clusters using the flexible-beta algorithm with beta =
192 -0.25, resulting in 81 initial classifications. These were optimized using REMOS1, REMOS2,
193 and OPTSIL, and the time elapsed during the optimization process was compared between the
194 three algorithms.

195 Real data sets were classified to 2 to 20 clusters using the flexible-beta algorithm (Lance &
196 Williams 1966) where beta = -0.25. For all real data sets and both classifications, the
197 dissimilarity measure was Sørensen index. Each partition was optimized using the REMOS1,
198 REMOS2, and OPTSIL methods. To assess differences in optimization success, mean

199 silhouette width and misclassification rate were calculated and compared between reallocation
200 methods, the original classification, and across numbers of clusters.

201 Lötter et al. (2013) argued that species fidelity should be a leading criterion in the evaluation
202 of vegetation classifications. Therefore, we used the Optimclass 1 index as a proxy for
203 interpretability of classifications (Tichý et al. 2010) that is the total number of faithful species
204 across all clusters. Faithful species were determined using Fisher's exact test and a $p=0.001$
205 threshold for supporting the null hypothesis that the species shows random distribution across
206 clusters (Chytrý et al. 2002). Hence, we also compared flexible-beta classifications optimized
207 by REMOS1, REMOS2, and OPTSIL, as well as the initial classifications in terms of the
208 number of faithful species across number of clusters.

209 The data analysis was carried out in the R software environment (R Core Team 2017) using
210 the cluster (Maechler et al. 2018) package. Source code for REMOS1 and REMOS2 is
211 supplied in the Electronic Supplement S3. OPTSIL was calculated using the optpart package
212 (Roberts 2016).

213

214 **Results**

215 When comparing optimization success, the REMOS and OPTSIL algorithms differed
216 markedly in MSW values they reached at different values for beta (Fig. 2). With beta ≤ 0 the
217 mean silhouette width of the initial classification was already high (MSW > 0.60), yet all
218 three optimization methods achieved minor improvement. Within this range of beta, the
219 largest increment in MSW was made by OPTSIL (+0.0134), less by REMOS2 (+0.0105) and
220 REMOS1 (+0.0118) (Table 1). On average, OPTSIL was superior to all other methods in this
221 respect, although differences were very slight ($|0.0013|$ to $|0.0029|$) between OPTSIL,
222 REMOS1, and REMOS2. From beta = 0.25 and higher, initial classifications showed a
223 dramatic decline in MSW; with beta = 0.5 and higher, MSW dropped below 0. OPTSIL was
224 able to optimize these initial classifications only to a limited degree: MSW ranged between
225 0.45 and 0.69 with beta = 0.25, and between 0.12 and 0.45 with higher beta. On the contrary,
226 REMOS1 and REMOS2 performed well, achieving a lowest median MSW of 0.599 with beta
227 = 0.75; even the minima were near 0.5. A very similar pattern was detectable with
228 misclassification rates. MR was near 0 with beta ≤ 0 for both optimization methods (Fig. 3).
229 Within this range, REMOS1 reached the lowest MR on average but its advantage over
230 REMOS2 was minimal ($|0.0002|$ difference; Table 2). REMOS1 and REMOS2 had slightly
231 lower MR than OPTSIL ($|0.0034|$ and $|0.0032|$ differences, respectively). All optimization
232 methods decreased MR in comparison with the initial classification (REMOS1: -0.0187,
233 REMOS2: -0.0185, OPTSIL: -0.0153). With increasing beta, especially with beta ≥ 0.5 ,
234 REMOS1 and REMOS2 kept MR at the same level, while OPTSIL resulted in gradually
235 higher values reaching medians over 0.1.

236 The number of iterations for REMOS1 were between 2 and 234, for REMOS2 between 2 and
237 53, and for OPTSIL between 0 and 67. Not surprisingly, from less efficient initial
238 classifications more iterations were necessary to reach a final solution; however, the upper
239 limit of number of iterations was never reached.

240 Visual checking of the classifications showed that with beta = 0 or lower all classifications
241 mirrored the a priori point aggregations efficiently (Figures S4-1 to S4-4). Classifications
242 differed mostly in the assignments of transitional points. With beta > 0 initial classifications
243 tended not to distinguish point aggregations as separate clusters. OPTSIL classification tended
244 to delimit one (or a few) heterogeneous clusters including several aggregations of many points

245 in a single cluster and several clusters with very few points distant from each other (see Fig.
246 S4-7). Additionally, OPTSIL tended to eliminate clusters completely, thus often keeping only
247 2 to 7 clusters from the initial eight (see Fig. S4-5 to S4-7). In a few cases REMOS2 also
248 eliminated one or two clusters but REMOS algorithms were rather consistent in delineating
249 point aggregations rather independently of the beta value.

250 There was a significant difference in the relationship between sample size and computation
251 time among the three optimization methods (Fig. 4). Considering average runtimes with 50
252 points REMOS2 was the fastest (0.0006 s), followed closely by REMOS1 (0.0010), while
253 OPTSIL ran approximately three times longer (0.0254 s). For larger samples, REMOS2 was
254 the fastest, completing classifications in 0.0008 s, 0.0010 s, 0.0018 s, and 0.0010 s, with 100,
255 200, 300, and 400 objects, respectively. However, its advance over REMOS1 was minimal,
256 which needed 0.0015 s, 0.0025 s, 0.0051 s, 0.0030 s on average. The lag of OPTSIL was even
257 more significant at these sample sizes: average runtimes were 0.2556 s, 3.7027 s, 13.7790 s,
258 and 24.4539 s with 100, 200, 300, and 400 objects.

259 On the Grasslands data set OPTSIL reached the highest MSW at all but two examined cluster
260 levels (Fig. 5). With 6 and 10 clusters REMOS1 performed the best and it was only slightly
261 worse than OPTSIL in all other cases. Interestingly, REMOS2 gave the same MSW values
262 with 2 to 5 clusters (likely due to identical final solutions), but at finer resolutions it was much
263 poorer. With 6, 7 and 9 clusters REMOS2 even decreased the MSW of the initial
264 classification. Regarding misclassification rate, REMOS1 performed the best with no negative
265 silhouette width values over all runs. As with MSW, from 2 to 5 clusters REMOS2 gave the
266 same result, but the weak performance with 6 or more clusters was visible here, too. OPTSIL
267 solutions ranked in intermediate position between REMOS1 and the initial classification, the
268 latter being the worst in all but two cases. These differences were not observable with
269 diagnostic species. Between 2 to 7 clusters all methods (including the initial classification)
270 showed similar numbers of diagnostic species, while at finer resolutions REMOS2 was the
271 best. Nevertheless, the Grassland data set is small, thus at this level the sizes of clusters are so
272 small and the number of diagnostic species so low that these differences are probably not
273 relevant.

274 With the Bryce data set OPTSIL produced the highest MSW at most cluster levels (Fig. 5).
275 REMOS1 and REMOS2 had very similar, often identical performance. With a minimal
276 difference they outperformed OPTSIL at two clusters. At 3 and 4 clusters they were slightly
277 worse than OPTSIL but this difference increased with the number of clusters, and became
278 striking from 7 and more clusters. The initial classification had the lowest MSW across the
279 tested numbers of clusters. REMOS1 and REMOS2 provided solutions with the lowest MR,
280 most often with no negative silhouette widths at all. OPTSIL had MR between 0.02 and 0.07,
281 while the initial classification had the highest MR in at all cluster numbers (MR between
282 0.048 and 0.15). OPTSIL performed the best in terms of diagnostic species at 3, as well as at
283 6 and more clusters. Interestingly, at 4 clusters the initial classification had the most
284 diagnostic species, while at 2 and 5 clusters REMOS algorithms reached the highest values.

285 On the Shoshone data set, OPTSIL reached the highest MSW across all cluster numbers,
286 REMOS1 was the second best, showing similar (in a few cases identical) MSW values with
287 REMOS2, and the worst was the initial classification (Fig. 6). REMOS1 had the lowest MR
288 again. This position was shared with REMOS2 between 2 and 5 clusters when both
289 algorithms provided no misclassifications. OPTSIL had MR between 0.04 and 0.07, which
290 positioned it behind REMOS2 in all but two cluster numbers. The initial classification had the
291 highest MR (between 0.15 and 0.27). Regarding the number of diagnostic species, the picture
292 was different. REMOS1 gained the highest numbers, again in a few cases together with

293 REMOS2, while OPTSIL was always inferior. The initial classification was again the worst in
294 all cases, except for the 10-cluster level, where REMOS2 had the fewest diagnostic species.

295

296 **Discussion**

297 In this paper we introduced the REMOS algorithms which can be used for improving already
298 existing classifications by reallocating misclassified objects using the silhouette criterion.
299 Two versions are available: REMOS1 reallocates only the single object with the lowest
300 silhouette width, while REMOS2 re-assigns all objects with negative silhouette width to their
301 respective closest neighbour cluster. We provide evidence on the high optimization success
302 and time efficiency of the new algorithms.

303 Our tests showed that the efficiency of the tested reallocation algorithms (REMOS1,
304 REMOS2 and OPTSIL) has different degrees of dependence on the initial classification.
305 Regarding MSW, the optimization success of OPTSIL is higher than REMOS algorithms'
306 when the initial classification already has rather high MSW; although, the difference is
307 usually small. Since mean silhouette width prefers spherical cluster shapes (Rousseeuw 1987),
308 it is typically high for classifications produced by group-forming methods, e.g. flexible-beta
309 with $\beta \leq 0$, and similar behaviour can be expected when applied to average linkage,
310 complete linkage, Ward's method, K-means, or PAM classifications. However, with chaining
311 algorithms, e.g. $\beta > 0$, REMOS1 and REMOS2 outperform OPTSIL. Chaining algorithms
312 optimize on criteria emphasising nearest neighbour distances which are not well reflected by
313 the traditional form of silhouette width also applied here (but see Lengyel & Botta-Dukát in
314 press), resulting in non-spherical clusters and low MSW. Using such classifications as input,
315 OPTSIL frequently converges into local optima, while REMOS algorithms provide more
316 robust optimization and reach high MSW. As it was shown by our examples, OPTSIL
317 solutions in these situations often fail to mirror the original cluster structure of the data set. In
318 concurrence with Roberts (2015), we suggest classifying the data set by a grouping method
319 first, and then optimizing the result with OPTSIL in order to reach the highest possible MSW.
320 Alternatively, REMOS algorithms seem more effective with other types of initial
321 classifications, although, their final MSW might be slightly lower than what is maximally
322 possible with OPTSIL.

323 Regarding misclassification rate, REMOS1 performed the best. In many cases REMOS2 led
324 to exactly the same solution containing no negative silhouette width values at all; however,
325 with the real data examples and higher number of clusters REMOS2 tended not to reach such
326 efficiency. The sensitivity to the initial classification of OPTSIL was visible also on the
327 presence of negative silhouette widths: OPTSIL had significantly higher MR than REMOS
328 algorithms when initiated from classifications with a chained structure. It must be noted that
329 different algorithms may reach the same value for MR, while their final solutions are not
330 necessarily identical. It occurred in some times with REMOS1 and REMOS2 that their final
331 solutions contained no, or only very few misclassified objects, while their classifications were
332 different. Even the number of clusters can differ between REMOS1 and REMOS2 despite
333 equal MR (e.g., Fig. S4-4). Such agreement in MSW is less probable due to its continuous
334 scale.

335 In general, optimizing a single criterion results in trade-offs for other criteria, and OPTSIL
336 and REMOS demonstrate this clearly. It is not surprising that OPTSIL reached the highest
337 MSW values, while REMOS outperformed OPTSIL in terms of MR. When comparing the
338 optimization success of OPTSIL and REMOS on MSW and MR, it must be noted that
339 OPTSIL directly maximizes MSW, a 'global' criterion of classification efficiency. REMOS,

340 on the other hand, has a more local perspective on classification efficiency, and focuses on
341 neighbourhoods of adjacent clusters. Surely, MSW and MR correlate strongly, and in general
342 optimizing MR will lead to high, although not necessarily optimal, MSW. In addition,
343 REMOS implicitly minimizes the absolute value of the sum of negative silhouette widths. In
344 our tests, this criterion behaved very similarly to MR, thus we present its results only in the
345 Electronic Supplement 5.

346 OPTSIL employs an anticipatory algorithm that tentatively reallocates an object to another
347 cluster, but then calculates the consequences of doing so before making the reallocation
348 effective. As a result, the trace of the optimization criterion is strictly monotonic increasing.
349 REMOS, on the other hand, identifies candidate objects to reallocate and makes the
350 reallocation effective immediately. In some cases this causes objects in the target cluster to
351 exhibit newly negative silhouette widths in the next iteration, and subsequent reallocations
352 must undo the negative consequences of a previous reallocation. As a result, the trace of the
353 optimization criterion shows non-monotonic behaviour, and in some cases oscillates or
354 exhibits cycles. While in general this behaviour is undesirable it may help avoid local optima
355 in a manner similar to genetic algorithms.

356 The difference between ‘global’ vs ‘local’ perspective can be seen on the classifications of the
357 artificial data sets (see the Electronic Supplement). OPTSIL solutions initiated from less
358 efficient classifications often contained one or more clusters with a single object, or a few
359 objects which were distant from each other (e.g., Figure S4-7). Such solutions are presumed
360 to have the highest possible MSW from the respective initial classification with the cost of a
361 few very heterogeneous or overlapping clusters and misclassified objects.

362 From the perspective of optimizing silhouette width, it is not correct to say that an object with
363 negative silhouette width is misclassified if reallocating it to its nearest neighbour cluster
364 decreases MSW. Rather, a misclassification is an assignment that lowers mean silhouette
365 width. However, as noted above, MSW cannot be high with many negative silhouette widths.
366 Alternatively, the viewpoint that correct classification reflects strictly positive silhouette
367 widths for as many objects as possible might be more straightforward than an ‘on-average
368 correct’ solution. This requires a decision from the investigator before choosing between these
369 methods.

370 An important property of REMOS2 and OPTSIL is that they are able to eliminate complete
371 clusters from the initial classification, thus the final number of clusters becomes lower than
372 the initial. This can be useful if the initial classification has more clusters than is optimal.
373 However, our simulation examples showed that the number of clusters by these methods, but
374 especially OPTSIL, can decrease even if the initial classification is not effective, despite its
375 cluster number corresponds to the number of point aggregations.

376 We found clear difference in computation time between the three methods. REMOS
377 algorithms were magnitudes faster than OPTSIL. This is not surprising considering that in
378 every iteration of the OPTSIL algorithm all possible reallocations of all objects to each cluster
379 are recalculated and only the one bringing the highest increment in MSW is accepted. In our
380 tests for computation time we used rather small data sets (i.e., containing max. 400 objects)
381 with clear cluster structure, and optimized initial classifications with relatively high MSW.
382 Presumably, such classifications would be faster to optimize than real data sets. Therefore,
383 our measured runtimes are likely to be shorter than what we can expect for larger and more
384 complicated data sets, less efficient classifications or more clusters. If time efficiency of the
385 analysis is crucial and the small difference in optimization success can be neglected,
386 REMOS1 or REMOS2 should be considered instead of OPTSIL.

387 Although we found OPTSIL sensitive to the initial classification, we must note that OPTSIL
388 performed poorly in situations which are scarcely realistic, since chaining algorithms are
389 rarely used in practice. If high silhouette width is a desired outcome it makes little sense to
390 begin with a classification emphasizing connectivity (e.g. single linkage or flexible-beta with
391 $\beta > 0$), and classifications emphasizing cluster disjunction (e.g. complete linkage or
392 flexible-beta with $\beta < 0$) should be preferred. In vegetation science, group-forming
393 methods are much more popular and straightforward, thus these drawbacks of OPTSIL may
394 not obtain in practice.

395 Tests on real data showed that OPTSIL combined with flexible-beta ($\beta = -0.25$) is more
396 efficient than REMOS algorithms in terms of MSW, although, the difference is often small.
397 As a contrast, with respect to minimizing the proportion of negative silhouette widths
398 REMOS1 provided consistently the best classifications. However, these differences may not
399 affect interpretability the same way since we could not detect consistent difference between
400 OPTSIL and REMOS algorithms in the number of diagnostic species. We suggest considering
401 which cluster validity measure fits the research question the best, and then decide between the
402 methods discussed above.

403

404 **Conclusions**

405 We present REMOS1 and REMOS2 as new reallocation methods for the optimization of
406 classifications and compare them with the related OPTSIL algorithm. When the initial
407 classification is already relatively efficient, most frequently OPTSIL gives the highest final
408 mean silhouette width; however, REMOS solutions are often only slightly worse. When the
409 initial classification has low mean silhouette width, OPTSIL performs poorly, while REMOS
410 algorithms are similarly straightforward as with more efficient initial classifications. With
411 respect to the proportion of misclassified objects, REMOS algorithms, especially REMOS1,
412 provided better classifications than OPTSIL, and this difference increased toward less
413 efficient initial classifications. REMOS algorithms are much time efficient to compute than
414 OPTSIL. We found no systematic difference in the number of diagnostic species between
415 vegetation classifications obtained by OPTSIL and REMOS algorithms.

416

417 **Acknowledgements**

418 The work of A.L. was supported by the National Research, Development and Innovation
419 Office, Hungary (PD-123997).

420

421 **References**

- 422 Chytrý M, Tichý L, Holt J, Botta-Dukát Z. (2002) Determination of diagnostic species with
423 statistical fidelity measures. *Journal of Vegetation Science* 13: 79-90.
- 424 Fleishman E (2015) *Vegetation structure and composition in the Shoshone Mountains and*
425 *Toiyabe, Toiyabe and Monitor ranges, Nevada*. 2nd Edition. Fort Collins, CO: Forest Service
426 Research Data Archive. <https://doi.org/10.2737/RDS-2013-0007-2>
- 427 Kaufman L, Rousseeuw PJ (1990) *Finding groups in data*. Wiley, New York
- 428 Lance GN, WT Williams (1966) A General Theory of Classificatory Sorting Strategies, I.
429 Hierarchical Systems. *Computer Journal* 9: 373-380.

- 430 Legendre P, Legendre L (2012) Numerical ecology, 3rd edn. Elsevier, Amsterdam
- 431 Lengyel A, Botta-Dukát, Z. (in press) Silhouette width using generalized mean – a flexible
432 method for assessing clustering efficiency. Ecology and Evolution, accepted
- 433 Lengyel A, Illyés E, Bauer N, Csiky J, Király G, Purger D, Botta-Dukát Z (2016)
434 Classification and syntaxonomical revision of mesic and semi-dry grasslands in Hungary.
435 Preslia 88: 201-228.
- 436 Lötter, MC, Mucina L, Witkowski ETF (2013) The classification conundrum: species fidelity
437 as leading criterion in search of a rigorous method to classify a complex forest data set.
438 Community Ecology 14(1): 121-132.
- 439 Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K (2018) cluster: Cluster Analysis
440 Basics and Extensions. R package version 2.0.7-1.
- 441 Peet RK, Roberts DW (2013) Classification of natural and semi-natural vegetation. In: van
442 der Maarel E, Franklin J (eds) Vegetation ecology, 2nd edn. Wiley-Blackwell, Oxford, pp 26–
443 62
- 444 Podani J (2000) Introduction to the exploration of multivariate biological data. Backhuys,
445 Leiden, NL.
- 446 R Core Team (2017) R: A language and environment for statistical computing. R Foundation
447 for Statistical Computing, Vienna, Austria. URL: <https://www.R-project.org/>
- 448 Roberts DW (1992) Plant Community Distribution and Dynamics in Bryce Canyon National
449 Park: Final Report for Project PX 1200-7-0966
- 450 Roberts DW (2015) Vegetation classification by two new iterative reallocation optimization
451 algorithms. Plant Ecology 216(5): 714-758.
- 452 Roberts DW (2016) optpart: Optimal Partitioning of Similarity Relations. R package version
453 2.3-0. <https://CRAN.R-project.org/package=optpart>
- 454 Tichý L, Chytrý M, Hájek M, Talbot SS, Botta-Dukát Z (2010) OptimClass: Using
455 species-to-cluster fidelity to determine the optimal partition in classification of ecological
456 communities. Journal of Vegetation Science 21: 287-299.

457

458 **Electronic Supplements**

- 459 Supplement S1: the R code of REMOS
- 460 Supplement S2: the R code for the simulated data set
- 461 Supplement S3: the Grassland data set (in txt format directly readable to R)
- 462 Supplement S4: Exemplary classifications of the simulated data set
- 463 Supplement S5: Comparison of reallocation methods on real data sets using misclassification
464 rate and the absolute sum of negative silhouette widths

465

466

467

468 **Table 1.** Differences between the (unoptimized) initial classification, REMOS1, REMOS2,
469 and OPTSIL solutions when the beta = 0 or lower in the flexible-beta classification. In the
470 cells are averages of differences calculated for each run as [MSW by the method in the row] –
471 [MSW by the method in the column].

	Initial	REMOS1	REMOS2	OPTSIL
Initial	0	-0.0105	-0.0118	-0.0134
REMOS1	0.0105	0	-0.0013	-0.0029
REMOS2	0.0118	0.0013	0	-0.0016
OPTSIL	0.0134	0.0029	0.0016	0

472

473

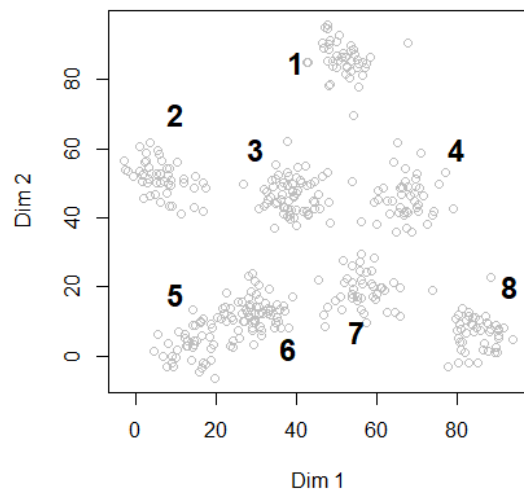
474 **Table 2.** Differences between the (unoptimized) initial classification, REMOS1, REMOS2,
475 and OPTSIL solutions when the beta = 0 or lower in the flexible-beta classification. In the
476 cells are averages of differences calculated for each run as [MR by the method in the row] –
477 [MR by the method in the column].

	Initial	REMOS1	REMOS2	OPTSIL
Initial	0	0.0187	0.0185	0.0153
REMOS1	-0.0187	0	-0.0002	-0.0034
REMOS2	-0.0185	0.0002	0	-0.0032
OPTSIL	-0.0153	0.0034	0.0032	0

478

479

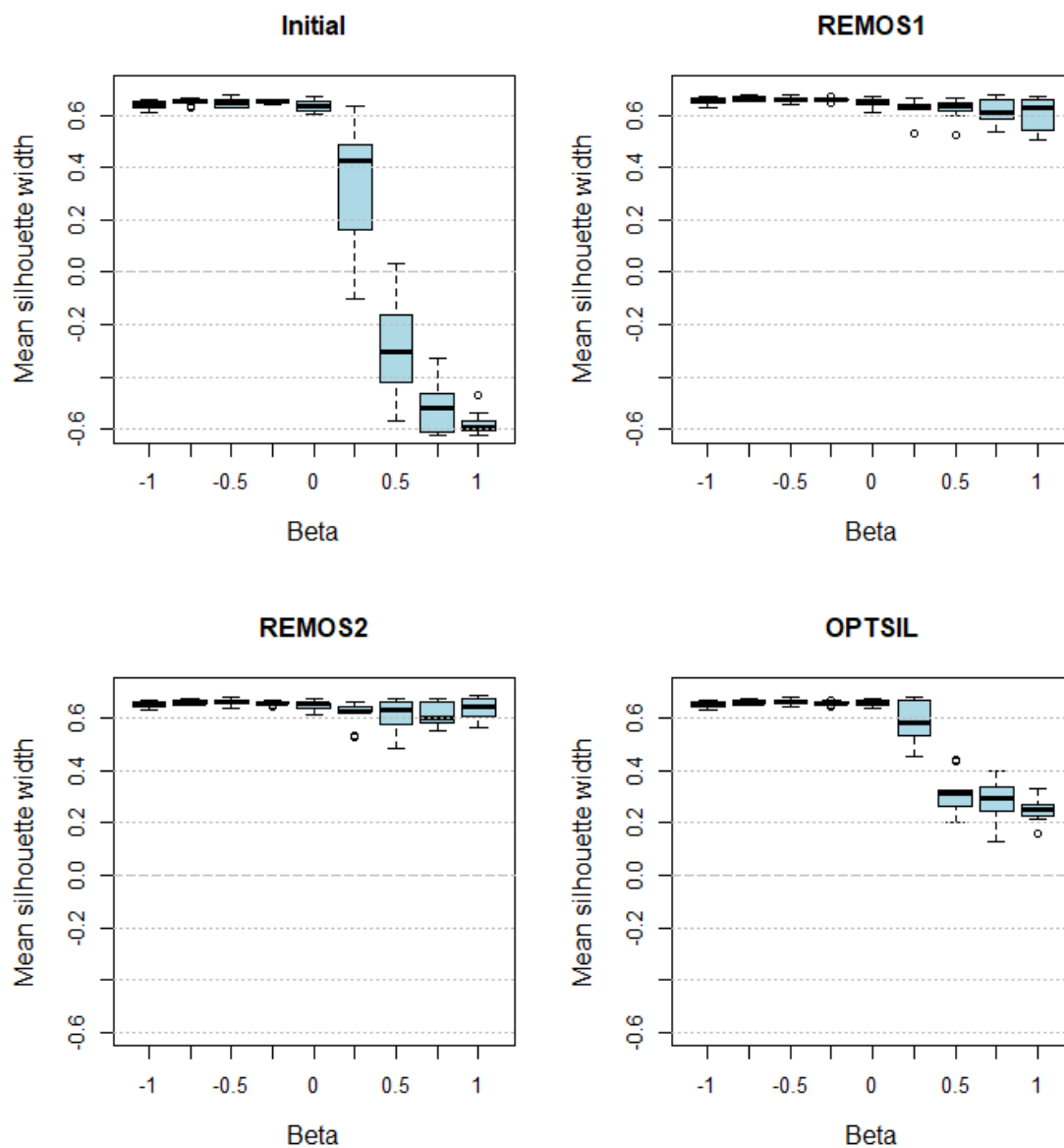
480 **Fig. 1.** The simulated data set containing 400 points in eight aggregations



481

482

483 **Fig. 2.** Comparison of the initial classification (without optimization), REMOS1, REMOS2,
484 and OPTSIL across different beta values of the flexible-beta classification based on mean
485 silhouette width.



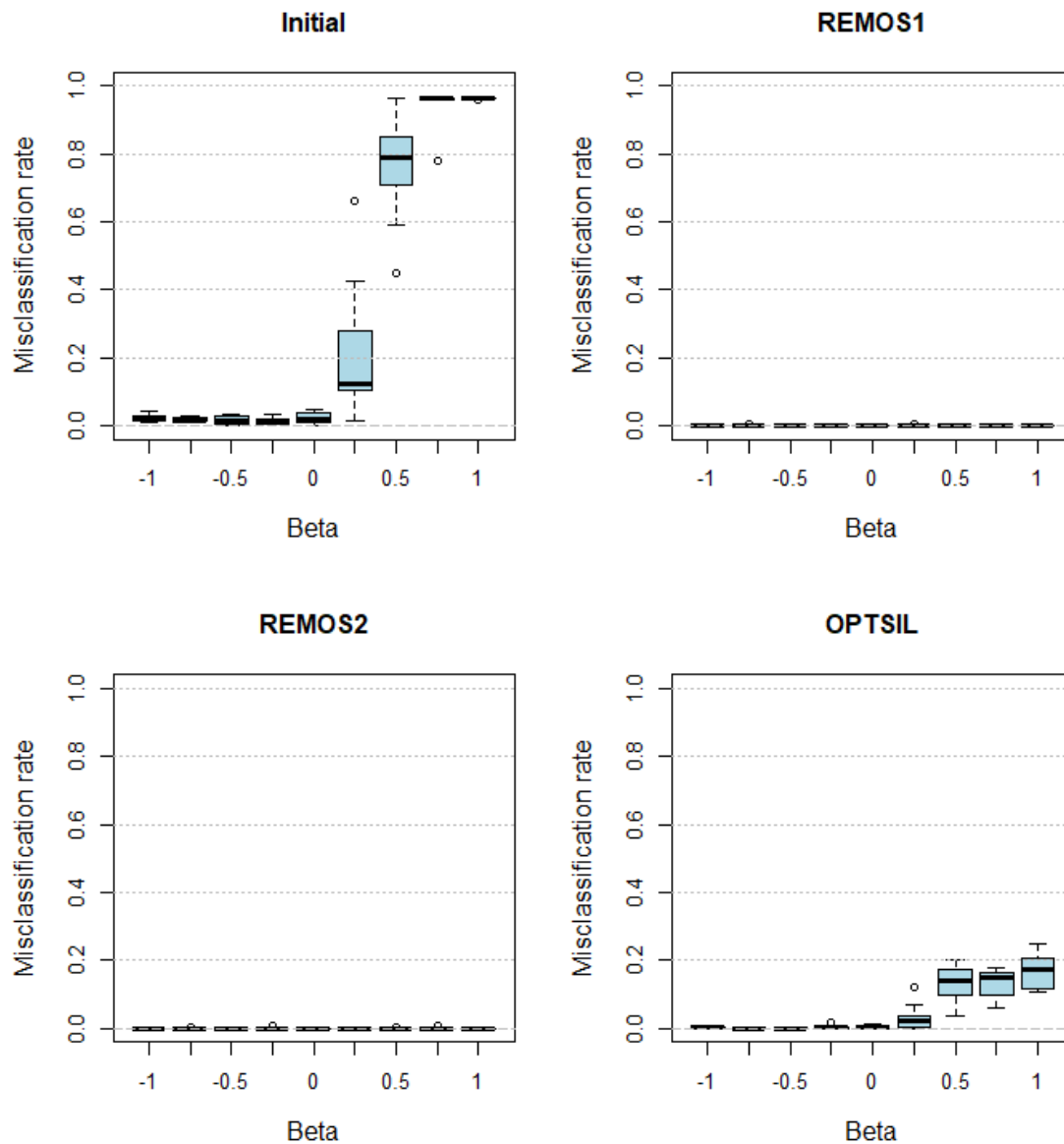
486

487

488

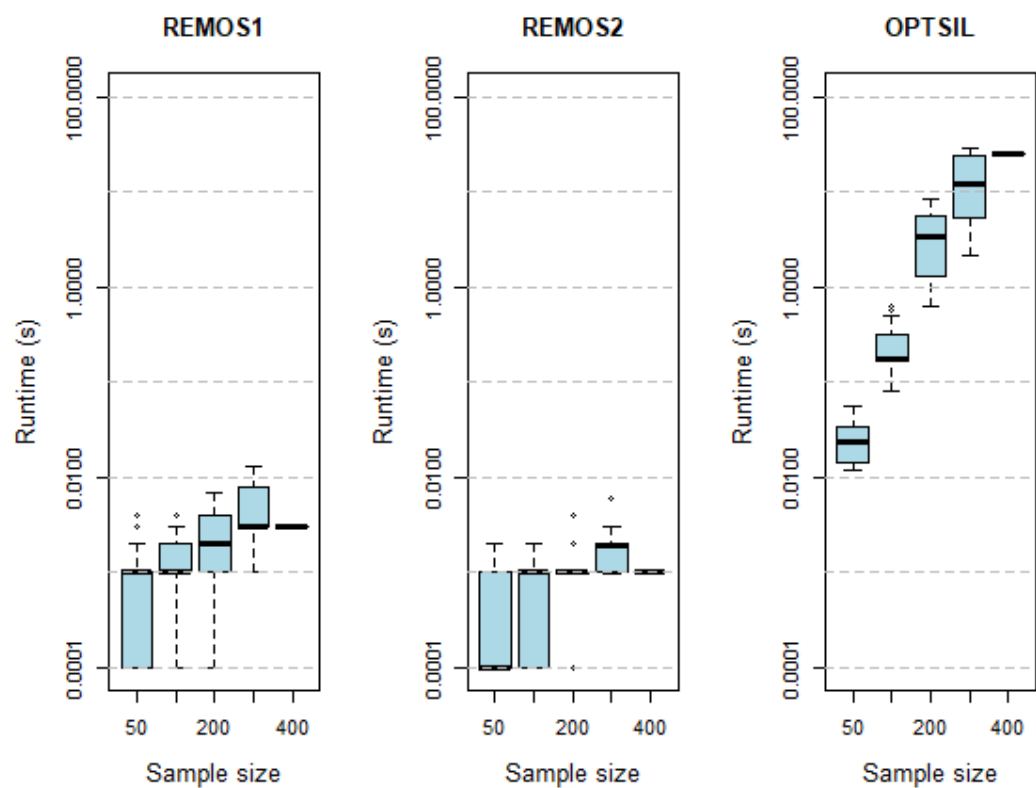
489 **Fig. 3.** Comparison of the initial classification (without optimization), REMOS1, REMOS2,
490 and OPTSIL across different beta values of the flexible-beta classification based on
491 misclassification rate.

492



493

494 **Fig. 4.** Computation times with different sample sizes by REMOS1, REMOS2, and OPTSIL.
495 Shortest computation times are truncated and replaced by 0.0001 s.

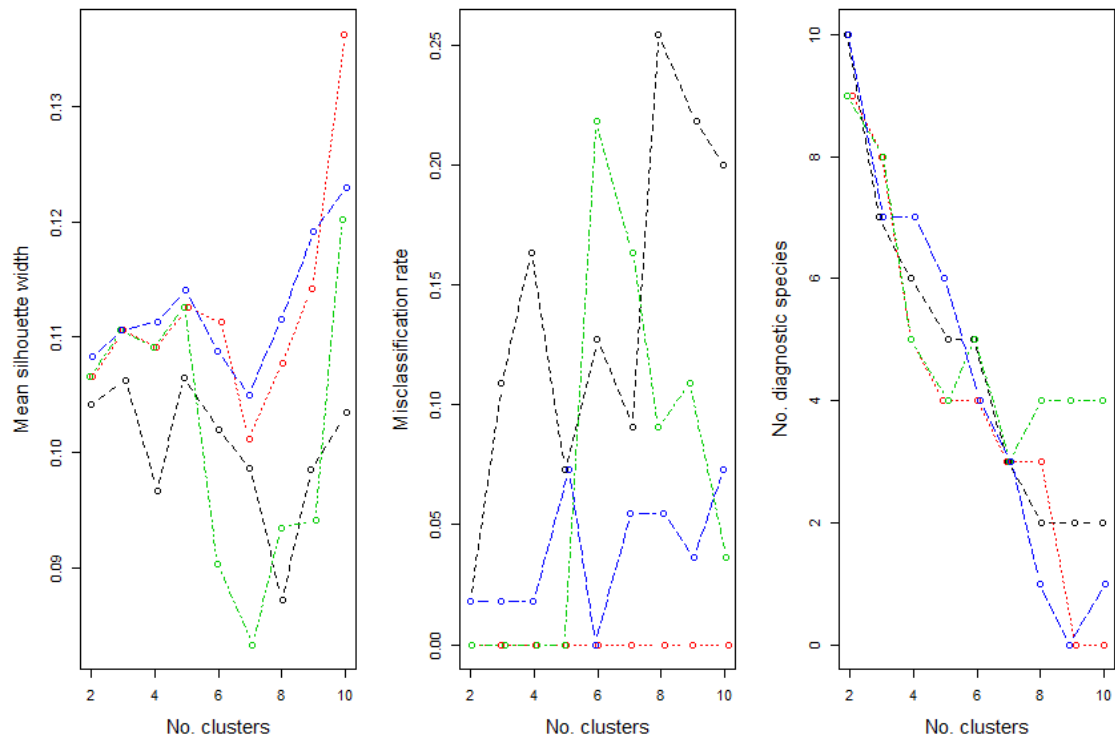


496

497

498

499 **Fig. 5.** Comparison of the initial classification (without optimization), REMOS1, REMOS2,
500 and OPTSIL solutions in terms of the change of mean silhouette width and number of
501 diagnostic species across the number of clusters on the Grassland data set. The initial
502 classification was produced by the flexible-beta method (beta = -0.25). To avoid overlap,
503 points are jittered in the horizontal direction on the graph. Colour code: red – REMOS1, green
504 – REMOS2, blue – OPTSIL, black – initial classification.



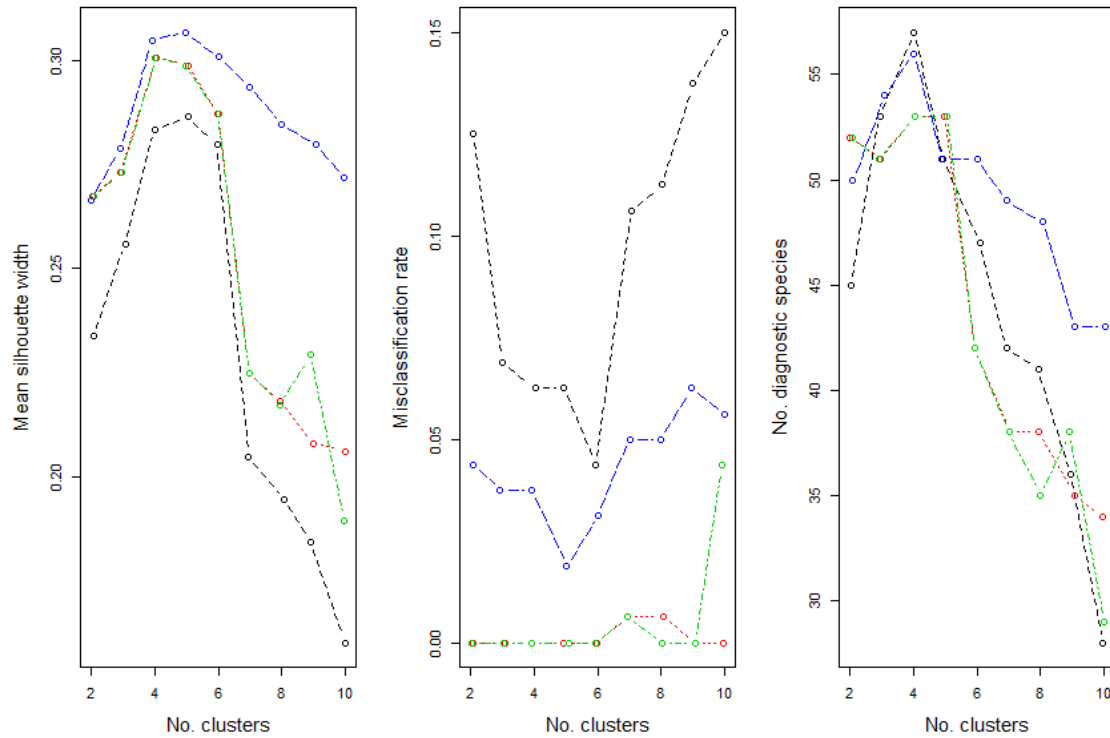
505

506

507

508

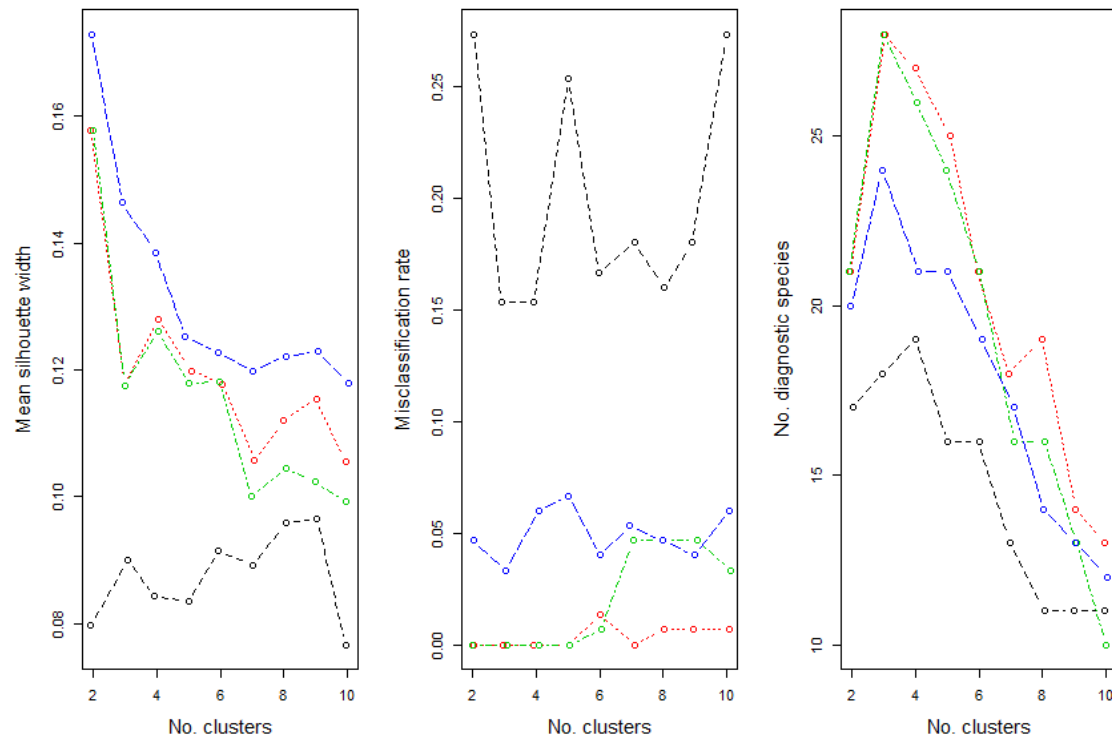
509 Fig. 5. Comparison of the initial classification (without optimization), REMOS1, REMOS2,
510 and OPTSIL solutions in terms of the change of mean silhouette width and number of
511 diagnostic species across the number of clusters on the Bryce data set. The initial
512 classification was produced by the flexible-beta method (beta = -0.25). To avoid overlap,
513 points are jittered in horizontal direction on the graph. Colour code: red – REMOS1, green –
514 REMOS2, blue – OPTSIL, black – initial classification.



515

516

517 Fig. 6. Comparison of the initial (without optimization) classification, REMOS1, REMOS2,
518 and OPTSIL solutions in terms of the change of mean silhouette width and number of
519 diagnostic species across the number of clusters on the Shoshone data set. The initial
520 classification was produced by the flexible-beta method (beta = -0.25). To avoid overlap,
521 points are jittered in horizontal direction on the graph. Colour code: red – REMOS1, green –
522 REMOS2, blue – OPTSIL, black – initial classification.



523