# A PROCEDURE FOR RECOGNITION OF UNCOMMON SPECIES COMBINATIONS IN SETS OF VEGETATION SAMPLES

by

DAVID W. GOODALL[1])

Department of Population and Environmental Biology,
University of California, Irvine

## INTRODUCTION

An argument sometimes raised against the use of objective methods in the study of vegetation, and particularly against non-selective sampling (e.g. GOUNOT 1961, IVIMEY-COOK & PROCTOR 1966), has been that the less common types of communities are likely to be ignored. Either they will not be recorded at all, or the records of each will be so few that they will not be recognised as distinct communities. It will be shown, however, that provided sampling is extensive enough to cover an uncommon community, even if only once, there is no bar to its identification by objective methods if it is sufficiently distinct from the other community types sampled.

In a collection of vegetation samples representing more than one community type, groups of species characteristic (in the general sense) of one or another community are correlated among themselves, and such groups of species enable the communities in question to be recognized. These interspecific correlations underlie the Zürich-Montpellier system of classifying vegetation, and have also been used as the basis for objective classification (e.g. GOODALL 1953, WILLIAMS & LAMBERT 1959, 1960) or ordination (e.g. GOODALL 1954, ORLOCI 1966) techniques. In many such cases, however, species occurring in few samples have been excluded from consideration, thus reducing the chances of distinguishing communities represented by a few samples only, whose characteristic species were also likely to be recorded only rarely. A modification of the same approach, however, enables the information provided by such uncommon characteristic species to be taken into account while still depending on interspecific correlation.

## PROCEDURE

If in a series of vegetation records a particular species is recorded

once only, that information in isolation leads nowhere. The records may constitute a uniform and homogeneous collection from a single vegetation type in which this species is uncommon enough for only a single record to be expected in a chance sample. On the other hand, the records may be heterogeneous and include one sample from a peculiar vegetation type of which this species is characteristic, and with which it would always be associated. Without further information one cannot distinguish between these very different possibilities.

Let us now suppose that a second species is also recorded once only, and that this single occurrence was in the same sample as the first species. Already the second hypothesis gains greatly in credibility, for the chances may be quite small that these two uncommon species should occur together. If there are one hundred samples, the probability that both these species occur in the same sample by chance is just .01. On the other hand, if there are additional species occurring once only, the hypothesis of a homogeneous set of samples may not be so untenable; for instance, if there are ten species each occurring once in a hundred samples, the chance that two or more of them should have their single occurrence in the same sample is

$$1 - \frac{99!}{90! \; 100^9} = .3718$$

If, however, *three* out of the ten species occur together, the probability is much less (.012), and if *four* do so rejection of the hypothesis of homogeneity becomes practically inescapable, for the chance of such an occurrence is only .0002. In this event, then, it is clear that this sample distinguished by the presence together of four unique species stands apart from the rest, and that these four species may provisionally be regarded as characteristic of a community type of which only a single sample has been recorded. If other species recorded only two or three times in the collection as a whole have also been recorded in the sample thus distinguished, this provides further confirmatory evidence of its distinctness (cf. MOORE 1962).

Apart from its value in recognizing the infrequent record of the uncommon community, this line of argument can serve the converse purpose of enabling one to recognize deviant quadrats in a set, and to remove them with a view to making the set internally more uniform.

The procedure described here, being capable of objective and precise definition, is well-suited to the electronic computer. A programme for this purpose has been written (in FORTRAN IV for the IBM 360), and is available on request.

| | QUADRAT NUMBER |
|---|---|
| | 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 |

1. Acacia bidwillii Benth.
2. Acacia dictyophleba F. Muell.
3. Alloteropsis semialata Hitchcock
4. Alphitonia excelsa Reissek
5. Amyliana gracilense N. Br.

6. Adisomeles salvifolia N. Br.
7. Acorualla acaulis Domin.
8. Aristida sp. *
9. Aristidochloa rubens F. Bauer
10. Aristidochloa thompifil F. Muell.

11. Arundinella spp. *
12. Atylosia marmorata Benke.
13. Blumea diffusa N. Br.
14. Borreria sp. *
15. Brachiaria holosericea D.K. Hughes

16. Bracyia atticlata Muel Arg.
17. Buchnera asiaticifolia R. Br.
18. Cadacua calgo Philip.
19. Calliris columellaris F. Muell.
20. Canthiigmodom parviflorus Stapf.

21. Capparis canescens G. Don.
22. Cassia absus L.
23. Cassia mimosoides DC.
24. Curculis trifolia (L.) Domin.
25. Chetlanthen caudata

26. Cheilanthes hirsuta
27. Cheilanthes velley F.v.M.
28. Chrysopogon pallidus Trin.
29. Cissus opaca F. Muell.
30. Coelorhachis polichollades L.

31. Coelospermum reticulatum Benth.
32. Commelina undulata Benth.
33. Crotalaria linifolia Linn. f.
34. Crotalaria trifoliastrum Willd.
35. Cucurligo recurvata N. Br.

36. Cycas media R. Br.
37. Cymbopogon bombycinus Domin.
38. Cymbopogon spp. *
39. Desmodium filiforme Hook. f.
40. Desmodium rhytidophyllus F. Muell.

41. Dianella sp.
42. Dichanthium sp.
43. Dicrania spp. *
44. Dolichandrone alternifolia
45. Enneapogon sp. *

46. Eragrostis spp. *
47. Eriachne armittii F. Muell.
48. Eriachne spp. *
49. Erythrochlamys chlorostachys Anil.
50. Eucalyptus alba Reinw.

51. Eucalyptus citriodora Hook.
52. Eucalyptus dichromophloia F. Muell.
53. Eucalyptus intermedia F.V. Baker
54. Eucalyptus leptophleba F. Muell.
55. Eucalyptus papuana F. Muell.

56. Eucalyptus polycarpa F. Muell.
57. Eucalyptus tessellaris F. Muell.
58. Eucalyptus spp. *
59. Euphorbia mitchelliana Boiss.
60. Eustrephus latifolius R. Br.

61. Evolvulus alsinoides L.
62. Ficus opposita Miq.
63. Fimbristylis cinnamometorum Hance
64. Fimbristylis dichotoma Vahl
65. Fimbristylis microcarpa F. Muell.

66. Fimbristylis monostachya Hassk.
67. Fimbristylis recta F. M. Bailey
68. Galactia tenuiflora Willd.
69. Glycine clandestina Wendl.
70. Glycine labialba Benth.

71. Glycine tomentosa L.
72. Grevillea glauca Knigh
73. Grevillea parallela Knight
74. Grewia latifolia F. Muell.
75. Grewia retusifolia Kurz.

76. Haloa perakhana F. Muell.
77. Heteropogon contortus Beuv.
78. Heteropogon triticeus Pers.
79. Hibbertia longifolia F. Muell.
80. Hibiscus rhodopetalus F. Muell.

81. Hybanthus enneaspermus F. Muell.
82. Indigofera colutea Merrill
83. Indigofera enneaphylla L.
84. Indigofera trifoliata L.
85. Indigofera spp. *

86. Justicia procumbens Wall.
87. Lepidosiloma fruticosum Benth.
88. Leonardia hexifolia Lauxil.
89. Leptanena curningianii (F. Muell.) Lane.
90. Melaleuca mundifolia F. Muell.

91. Melaleuca 'etunostecnva'
92. Melaleuca viridiflora Soland.
93. Merremia dentata
94. Mnesaria parviflora (Benth.) Kuntze
95. Panicum spp. *

96. Panpalidium distans D. K. Hughes
97. Perotis rara R. Br.
98. Petaloctylon laxmiki F. Muell.
99. Petaloctylon pubescens Domin.
100. Petaloctylon pubescens Domin.

101. Phyllanthus simplex Retz.
102. Piselia cornucopiae Vahl.
103. Planchonia careya (F. Muell.) Knuth
104. Polygala chinensis L.
105. Pseudopogonatherum contortum
      (Brogn.) A. Camus

106. Pterocaulon redolens Boer.
107. Pterocaulon sphacelatum Benth. &
      Hook. f.
108. Rhynchelytrum roseum (Willd.)
      C. E. Hubb
109. Rhynchosia minima DC.
110. Rhynchospora pterygasta F. Muell.

111. Rottboellia formosa R. Br.
112. Sarcilla linearioria R. Br.
113. Schizachyrium sp.
114. Scleria browmii Kunth
115. Securinega virosa Baill. Adananonia

116. Setaria glauca Beauv.
117. Sida corvifolia L.
118. Sida rhombifolia L.
119. Sorghum plumosum Beauv.
120. Tephrosia filipes Benth.

121. Tephrosia juncea Benth.
122. Themeda australis Stapf.
123. Triacogyne spp. *
124. Tristania spp. *
125. Uvaria cylindracea Benth.

126. Vernonia cinerea Less.
127. Vitadinia brachycomoides F. Muell
128. Waltzbergia spp. *
129. Xedulia juliflorales F. Muell.
130. Xanthorrhoea sp.

131. Xornia diphylla Pers.

---

Additional species present in three quadrats or fewer (the quadrat number is followed by the cover percent in parentheses):

* Certain groups of species could not be consistently distinguished in the field, and accordingly have been combined for the purpose of this table. They are:

8. Aristida browniana Henrard, A. pluraris Henrard, A. holathera Domin., A. hygrometrica R. Br., A. ingrata Domin., A. pruinicosa Domin., A. queenslandica Henrard, A. superpondens Domin., A. vallida F. V. Bailey.

11. Arundinella repalensis Trin., A. setosa Trin.

14. Borreria brachystoma Valcnton, B. laevigata Mart. & Gal.

37. Cymbopogon exaltatus A. Camus, C. refractus A. Camus.

42. Dichtanis adscorten (N.B. & K.) Henrard, D. sibbosa Beauv., D. sesostorionsum F.M. Bailey.

45. Enneapogon pallidus Desv., E. spp.

46. Eragrostis australasiensis Domin., E. bella Domin., E. brownii Nees., E. sp. aff. brownei, E. cumingii Steud., E. elongata Jacq., E. pubescens Steud.

48. Eriachne mucronata N. Br., E. obtusa R. Br.

52. Eucalyptus crebra F. Muell., E. decompophylla F. Muell.

82. Indigofera australis Willd., I. pratensis F. Muell.

95. Panicum mindana Fluege, P. decompositum N. Br., P. sp. aff. effusum N. Br., P. seminudum Domin., P. simile Domin., P. spp.

123. Triacogyne angens N. Br., T. elatior N. Br. T. spp.

128. Waltzbergia zolliflora F. Muell., W. grandiflora Cheel, W. quaroulens Oa.

## EXAMPLES

The procedure described may be made clearer by some examples; as a beginning we take some data collected between 1959 and 1961 in North Queensland.

a) SAVANNAH WOODLAND IN NORTH QUEENSLAND

The data in question were collected on Springmount Station between Mareeba and Dimbulah (17°0′ S 145°20′ E). A full account of the study will be published in due course, and a few brief particulars only need be given here. The vegetation of the area is in the main savannah woodland with an open canopy of *Eucalyptus* spp. (especially *E. leptophleba*) above an understorey of tall grasses dominated by *Themeda australis* and *Heteropogon contortus*. Records were taken on 67 quadrats of 50 × 5 metres, one placed at random in each square kilometre of country. Cover was estimated for all species of vascular plants present, and these data are recorded in Table I.

A total of 321 species were recorded, and of these 101 occurred in a single quadrat only. If these 101 species were distributed at random among the 67 quadrats, the numbers in each quadrat would be distributed approximately as a Poisson variate, with parameter $\frac{101}{67}$. The observed distribution is:

| No. of unique species | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 ... 18 |
|---|---|---|---|---|---|---|---|---|---|
| No. of quadrats | | 32 | 13 | 10 | 5 | 1 | 3 | 0 | 0 | 2 ... 1 |

Agreement with expectation may be tested by $\chi^2 = \dfrac{\Sigma(x - \bar{x})^2}{\bar{x}} = 324,22.$ which with 66 degrees of freedom is very highly significant.

An alternative test is possible by considering only the quadrat differing most markedly from the remainder. It is reasonably evident that Quadrat 49 must be regarded as exceptional. It contains 18 out of the 101 unique species. The chance that any one out of 67 samples should contain as many as 18 out of 101 species, each occurring in one of them, could in principle be calculated as follows. The probability that $N$ unique species independent of one another in their distribution will be partitioned among $M$ samples so that $n_0$ of the quadrats have none of them, $n_1$ have one only, and so forth, is given by

$$\frac{M! \quad N!}{M^N \prod\limits_{i=0}^{N} \left\{ n_i! (i!)^{n_i} \right\}}$$

If the values of this expression are summed for all sets of $n_i$ where $\sum_{i=R}^{N} n_i > 0$, the sum will give the probability $P_R$ that at least one sample will contain at least $R$ of the unique species.

A computer programme [1]) has been written to perform this calculation; but for the particular set of data presented the calculation is too long even for the computer. Luckily, a convenient approximation is available. The probability that $R$ or more out of $N$ unique and independent species should occur in a specified quadrat (out of $M$) is given by summing the appropriate terms of the binomial expansion [2])

$$p_1 = \sum_{i=R}^{N} {}^{N}C_i \left(\frac{M-1}{M}\right)^{N-i} \left(\frac{1}{M}\right)^{i}$$

The probability that one or more of the $M$ quadrats should have such a concentration of these species is then approximated by

$$P = 1 - (1-p_1)^{M}$$

The closeness of this approximation has been tested with some simpler combinations of values for $M$ and $N$ (Table II). It will be

TABLE II

Probability — $P_R$ that at least $R$ out of $N$ species will occur together in at least one of $M$ samples.

|  | $N = 6$  $M = 20$ | | $N = 30$  $M = 20$ | | $N = 70$  $M = 10$ | |
|---|---|---|---|---|---|---|
|  | Exact | Approximation | Exact | Approximation | Exact | Approximation |
| $P_2$ | .5640 | .4865 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| $P_3$ | $.4454 \times 10^{-1}$ | $.4366 \times 10^{-1}$ | .9947 | .9844 | 1.0000 | .9999 |
| $P_4$ | $.1728 \times 10^{-2}$ | $.1727 \times 10^{-2}$ | .7918 | .7146 | 1.0000 | .9798 |
| $P_5$ | $.3594 \times 10^{-4}$ | $.3594 \times 10^{-4}$ | .2880 | 2703 | .8725 | .7599 |
| $P_6$ | $.3125 \times 10^{-6}$ | $.3125 \times 10^{-6}$ | $.6490 \times 10^{-1}$ | $.6364 \times 10^{-1}$ | .3973 | .3569 |
| $P_7$ | | | $.1146 \times 10^{-1}$ | $.1141 \times 10^{-1}$ | .1114 | .1070 |
| $P_8$ | | | $.1693 \times 10^{-2}$ | $.1692 \times 10^{-2}$ | $.2384 \times 10^{-1}$ | $.2361 \times 10^{-1}$ |
| $P_9$ | | | $.2135 \times 10^{-3}$ | $.2135 \times 10^{-3}$ | $.4156 \times 10^{-2}$ | $.4149 \times 10^{-2}$ |
| $P_{10}$ | | | $.2323 \times 10^{-4}$ | $.2323 \times 10^{-4}$ | $.5986 \times 10^{-3}$ | $.5984 \times 10^{-3}$ |
| $P_{11}$ | | | $.2194 \times 10^{-5}$ | $.2194 \times 10^{-5}$ | $.7151 \times 10^{-4}$ | $.7151 \times 10^{-4}$ |

All approximations below this line are accurate to four figures at least.

---

[1]) This programme (in FORTRAN IV, for the IBM 360/50) is available to enquirers.

[2]) The notation ${}^{n}C_m$, which may be unfamiliar to some readers, stands for the number of combination of $n$ things taken $m$ at a time.

seen that, for the smaller probabilities, the approximation is very close. In the present case, this approximation gives a value of $1.11 \times 10^{-12}$ for $P$.

We may now go further and consider those species present in two quadrats, of which there are 57.

Again, the hypothesis of random distribution cannot be sustained, for there are 2211 pairs of quadrats $\left(\dfrac{67 \cdot 66}{2}\right)$ which could have contained these species, distributed as follows:

No. of species                0    1    2    3
No. of quadrat pairs   2160   40   7    1

The value of $\chi^2$, with 2210 degrees of freedom, is 2930 with a probability of about $10^{-25}$.

Again, we note that Quadrat 49 figures prominently, for of the 57 species occurring in two quadrats seven occur in Quadrat 49. The chances of such a concentration of this new class of species in this same quadrat may be calculated by an argument similar to that already used for the unique species. It is

$$p_2 = \sum_{r=8}^{57} {}^{57}C_r \left(\frac{65}{67}\right)^{57-r} \left(\frac{2}{67}\right)^r$$

Similarly, there are 30 species occurring in three quadrats of which seven are in Quadrat 49, giving

$$p_3 = \sum_{r=8}^{30} {}^{30}C_r \left(\frac{64}{67}\right)^{30-r} \left(\frac{3}{67}\right)^r$$

and so forth.

The question now arises how these probabilities from classes of species differing in rarity are to be combined. If a number $n$ of independent tests give results with probability $p_i$, their joint probability may be determined (FISHER 1963) by calculating

$$-2 \sum_{i=1}^{n} \ln p_i,$$

this sum being distributed as $\chi^2$ with $2n$ degrees of freedom. In the present instance, the groups of species with different frequency may be regarded as independent tests of the hypothesis that a particular quadrat falls into the same population as the rest. If the probabilities for Quadrat 49 are combined in this way for all species occurring in 34 quadrats or less, the $\chi^2$ variate obtained is 158.21 with 68 degrees of freedom, which has a probability of less than $10^{-10}$.

The probability that so large a $\chi^2$ value would occur in one or more out of 67 quadrats is about $2 \times 10^{-9}$, which is accordingly the significance one may ascribe to this deviation from expectation,

given the observed numbers of species of differing degrees of rarity.

It is thus clear that this quadrat cannot be regarded as a sample of the same vegetation type as that presumptively represented by the other 66 quadrats. Let Quadrat 49 therefore be separated from the rest, and let them be reconsidered in the same way.

With the removal of Quadrat 49, there are now 91 species confined to a single quadrat; six of these species are in Quadrat 47. This quadrat also has 11 of the 56 species confined to two quadrats, five of the 26 confined to three, and so forth. Proceeding in the same way as before, we can calculate the combined probability that species of the various degrees of rarity indicated will be concentrated in any one out of 66 remaining quadrats as: $2 \times 10^{-9}$.

This is in fact the smallest such probability for any quadrat. Accordingly we proceed to remove Quadrat 47 from the set, and start again. By repeated application of this process, we find that 15 other quadrats (viz., Nos. 36, 63, 64, 48, 51, 67, 50, 65, 37, 62, 66, 54, 30, 58, and 15) also include significantly large concentrations of uncommon species and so should be removed.

If there is in a set of quadrat data a single quadrat from a community deviating from the rest, characterized by a number of characteristic species of reasonably high frequency and presence, it is to be expected that such a quadrat will be identified and excluded by the process described. If, however, there are two such quadrats from the same deviant community, neither of them considered in isolation may qualify for rejection. To cover such a possibility, those species occurring in two quadrats only are considered further. If the same pair of quadrats contain more than one of the species recorded in two quadrats only, they may well represent a community type differing from the rest, and if this conjunction is such as could not easily happen by random assortment of the species records, the two quadrats in question may merit removal, as a further step in restricting the set to samples of a single community type.

In the present case, after the 17 quadrats listed have been removed on account of unacceptable individual deviations from expectation, there are 35 species occurring in two quadrats only. Only two of these occur in the same pair of quadrats, and this does not represent a significant degree of concentration on these species. Thus, the residue of 50 quadrats may be regarded of these criteria as uniform, and not including any quadrats deviating unacceptably from the norm of the community.

If the same tests are repeated on the 17 quadrats which have been removed, it is found that Quadrat 49 stands apart, with 21 species not shared with any of these others (the total number of species unique to one or other of these aberrant quadrats is 108). As far as

the present tests go, the other 16 quadrats removed from the original set may be regarded as samples from the same population.

b) DATA OF LAMBERT AND WILLIAMS FOR TUMULUS HEATH

These data, already analyzed by LAMBERT & WILLIAMS's (1962) "nodal analysis" technique, cover the floristic composition of 20 stands within an area of heathland intersected by valley bogs in Hampshire, England. Full data were made available by the authors' courtesy. A total of 80 species were recorded.

Analysis by the method described led to the rejection first of Stand 8 (10 unique species out of 24), followed by Stands 20, 19, 13, 3, and 11. After these six stands have been removed, no other single stand deviates significantly from the rest. In this case, however, a test for pairs of stands deviating from the rest, on the lines suggested in the preceding section, leads to further rejections.

At this point, there are six species occurring in two only out of the 14 remaining stands. These are:

|  | In stands |
| --- | --- |
| Carex panicea | 9, 14 |
| Juncus acutiflorus | 9, 14 |
| Narthecium ossifragum | 9, 14 |
| Pedicularis sylvatica | 9, 17 |
| Potentilla erecta | 7, 17 |
| Pteridium aquilinum | 7, 10 |

It will be noted that three out of the six occur in the same pair of stands. Since there are 91 possible pairs of stands among 14, the chance that three or more will occur in the same pair by random assortment is .0024. Consequently, these two stands may be recognized as differing from the other 12.

When they are removed there is still no single stand calling for rejection, but there are five species occurring in two of the remaining 12 stands:

|  | In stands |
| --- | --- |
| Drosera intermedia | 2, 12 |
| Drosera rotundifolia | 2, 12 |
| Eriophorum angustifolium | 2, 12 |
| Potentilla erecta | 7, 17 |
| Pteridium aquilinum | 7, 10 |

The probability that three out of these five species should occur in the same two stands is about .0022. Again, the rejection of these two stands is indicated. The remaining ten show no further concen-

trations of uncommon species, and may be regarded as samples from
the same population.

The ten stands removed in this way were examined afresh. The
six quadrats rejected individually still call for rejection, though the
remaining four (2, 9, 12, and 14) do not differ significantly by these
criteria. When the six quadrats still rejected are examined again,
four are again rejected, leaving only Nos. 3 and 11. The rejected
four, however, have similar numbers of unique species, and do not
show any significant pairing among themselves on these criteria.
Thus, the data have been divided into subsets of ten, four, two, and
four at the usual significance level of 5 %.

There is broad agreement between these results and those of
LAMBERT & WILLIAMS's "nodal analysis" — which likewise is based
on presence or absence only. The four stands 2, 9, 12, and 14
constitute LAMBERT & WILLIAMS's Community 3, described as
"wet-heath". Stands 8, 13, 19, and 20 (the other set of four),
together with stand 3, constitute their Community 1, of moderately
grazed heath on soil with slight podzolization. The undifferentiated
residue of ten quadrats, together with stand 11, constitute their
Communities 2 and 4, which appear not to be distinguishable on the
criteria used here.

c) DATA OF FIJAŁKOWSKI FOR COMMUNITIES WITH ADONIS VERNALIS

FIJAŁKOWSKI (1961) published a table (No. 2) giving complete
floristic data for a series of 29 stands in which *Adonis vernalis* was
abundant. A total of 373 species was recorded in these stands.
Analysis leads to rejection, one after the other, of no fewer than 17
stands, the residue consisting of those numbered 3, 4, 6, 9, 10, 14,
15, 17, 18, 19, 20, and 23. If the 17 rejected stands are then analyzed
again, six are rejected — numbers 1, 2, 12, 13, 26, and 28. Re-
analysis of these six leads to the rejection of 1, 2, 12, and 13, and
reanalysis of the last four to the rejection of 1 and 2. Thus, this
procedure leads to the division of the 29 quadrats into a group of 12,
a group of 11 and three groups of two.

FIJAŁKOWSKI did not claim that these 29 samples were from the
same type of vegetation — in fact, he divided them into six com-
munities. There is a certain measure of agreement between the
groups recognized in the present analysis and those into which
FIJAŁKOWSKI divides them, though he was using the usual wide
range of floristic characters and not merely the presence and absence
of less common species. Stands 1 and 2 are the samples he allots to
the Corylo-Peucedanetum cervariae, and stand 26 is the
only representative of the *Thalictro-Salvietum*. The whole of the
uniform residue of 12 stands falls into the *Brachypodium pinnatum-
Teucrium chamaedrys* association or the Carici-Inuletum. But

stands 5, 12, 13, 15, 21, and 22 also are listed in these communities
although on the criteria here they are distinct; and the Cynosurion
(stand 29) and most samples from the *Festuca sulcata-Koeleria
gracilis* association are not here distinguished from those of the two
most abundantly represented communities.

If our procedure is applied to the two main communities sep-
arately, it is found that the *Brachypodium pinnatum-Teucrium cha-
maedrys* association is far from uniform, for stands 11, 8, and 5 (in
that order) are again rejected; the same applies to the Carici-
Inuletum, with stands 13, 12, 22, and 21 being rejected. FIJAŁ-
KOWSKI does in fact separate stands 21 and 22 as a distinct facies,
but stands 12 and 13 are regarded as typical of the association,
along with stands 14—19.

d) Data Selected for Uniformity

The sets of data mentioned above have in each case been re-
cognized as covering several distinct communities. Where, on the
other hand, the data in question have been selected or sorted to
represent a single type of community, the results are quite different.
The tests have, for instance, been applied to tables for single as-
sociations published by McVean & Ratcliffe (1962), and by
Szynal (1962), and have failed to detect any stand as deviating
significantly from the association norm. As an example of this may
be mentioned the table for the Erico-Pinetum silvestris, pub-
lished by Braun-Blanquet et al. (1954) in their account of the
vegetation of the Swiss National Park in Graubünden. This table
includes 25 relevés, in which a total of 155 species were recorded.

Applying the technique described here showed that the 25 re-
levés constituted a uniform group at the 5 % significance level.
When, however, the first two relevés from Table III for the Pino-
Caricetum humilis (another association listed under the Pino-
Ericion alliance) were included, they were both rejected with high
significance (probabilities of .00001 and .002 respectively).

## APPLICATION TO TESTS OF THE DIFFERENCE BETWEEN TWO QUADRATS

Since this technique is non-parametric, calling for no knowledge
of the distribution of the data studied, it may be applied to a
comparison of two quadrats only, to provide a valid though in-
sensitive test of the hypothesis that the two quadrats could reason-
ably be regarded as samples from the same population. All that is
required is to compare the number of species recorded in one and
not the other. If the partition of the number of species recorded
once only between the two communities is compared with the

binomial distribution for $p = q = 1/2$, then the two-tail pro-
bability [1]) gives immediately the chance that the two samples could
have been drawn from the same population. As an example, one
may refer to MATUSZKIEWICZ's (1956) data for Quercetalia
pubescentis in Poland (Tables 1 and 2); community 15 has 15
species which do not occur in community 17, whereas community 17
has 5 which do not occur in community 15. The significance level
of this difference is

$$2 \cdot 2^{-20} \cdot \sum_{r=0}^{5} \frac{20!}{r! \ (20-r)!} = 0.0413$$

Thus a difference of this magnitude is unlikely, and MATUSZ-
KIEWICZ's allocation of the two communities to the same association
may be questioned [2]). If, further, we test each pair of the nine
communities (Nos. 10 to 18) included in the Querco-Potentil-
letum albae association, we find that ten of the 36 comparisons
show a difference exceeding that corresponding with the 5 % sig-
nificance level, and five even exceed the 1 % level. This could
happen by chance only very rarely, though it is difficult to estimate
the combined probability in view of the lack of independence
between the various comparisons.

If the data are examined further, one finds that most of the im-
probable deviations in the numbers of unmatched records occur in
the lists for communities 17 (five with $P < 0.05$) and 16 (three
with $P < 0.05$). If these two are excluded, the remaining seven
show only two comparisons out of 21 with a difference exceeding
the 5 % point, which could well be a chance effect. There is thus
good evidence for regarding the communities 16 and 17 (which do
not differ significantly from one another) as somewhat apart from
the others (10 to 15, and 18) allotted to the same association.

The proposed test seems to be the only way in which conclusions
at specified confidence levels can be drawn when the total inform-
ation available consists of two species lists. Where more extensive
material is to hand, it would of course be foolish to disregard it.
Communities often differ in respects to which this test would be
completely insensitive; two communities similar in floristic diversity
might not be distinguished by this test, no matter how different they
were in other respects — they might not have a single species in
common. On the other hand, if this test shows that the communities
are distinct — if they could not have been drawn from the same

---

[1]) A short table of limiting values at different significance levels is given in the
Appendix.

[2]) The $2 \times 2$ $\chi^2$ test, which might at first glance be thought appropriate, is
inapplicable to counts of species varying enormously in frequency, if only because
the number of species absent from both quadrats is indeterminate.

population without invoking the long arm of chance — this is a conclusion which more extensive data are unlikely to negative. But the distinction between the two may often be rather trivial — perhaps only in aspect. If two desert areas are compared, one before and one after rain, they would clearly — and properly — be recognized as different, though the difference might only reflect the immediately preceding conditions. The method provides, however, a way of drawing valid conclusions from data which might otherwise be regarded as having only descriptive value.

This test is analogous to the various coefficients of community which have been proposed from the time of JACCARD (1901; see DAGNELIE 1960, SOKAL & SNEATH 1963). These compare the number of species common to the two communities with the total recorded in both, but significance tests have not usually been applied, and are possible only where the two communities are considered as samples from a larger set (GOODALL 1968).

## DISCUSSION

The procedure described was originally envisaged rather as a means of "purifying" a set of quadrat records by removing any clearly aberrant records than of partitioning it into a number of discrete subsets. Where such a set contains a single record from another community characterized by quite distinct floristic composition, it is unlikely that it will fail to be recognized; the same is true if a sample is heterogeneous, part only belonging to a floristically distinct community. If, on the other hand, the set is divided fairly equally between two communities, it is unlikely that this procedure will be able to separate them. Even a single deviant quadrat may not be recognized if its distinguishing features are quantitative rather than qualitative, or if it is floristically poor, so that the number of unique species it contains is small.

The method used here in selecting individual quadrats for rejection is closely related to the deviant index published elsewhere (GOODALL 1966a), for use in the classification of a set of individuals characterized by a large number of attributes. This is an expression of the probability with which a random assortment of the observed attribute values would differ from the norm of the whole set as much as a particular individual does. In the present case, the individuals are the quadrats, and the attributes (all binary) are the presence or absence of those species occurring in not more than half the quadrats; the deviant index is then the probability of the $\chi^2$ value obtained by combining probabilities for the different species categories.

For the second stage of the process here described, a similarity

index (see GOODALL 1964, 1966b) could also have been used.
Instead of considering only species confined to two quadrats, one
could take into account all species present in both quadrats, and
compute the overall probability that the two would be as alike by
random assortment of species records. If the minimum probability
for any pair was less than that which could be expected as a
minimum at the significance level selected, then that pair would be
rejected. Such a procedure would take much fuller account of the
information available, but would be much more time-consuming
than that proposed here.

Though this procedure may show beyond a reasonable doubt that
a particular quadrat differs from the rest of the set, the question
will always remain as to what sort of difference it is. Perhaps it is a
sample from a different community; or perhaps it is merely a sample
taken from the same community in a different aspect — following a
localized rainstorm, for instance, which led to germination of
annuals elsewhere present only as seeds and consequently not re-
corded. If the quadrat size is not large enough to cover the full
pattern of the community, whether dependent on local differences
in climatic or edaphic conditions or on cyclical regeneration, dif-
ferences may be detected and quadrats representing the less com-
mon elements of the pattern may be rejected. In all these cases,
the results of the procedure call for interpretation by ecological
insight; but even where the distinction is not one between com-
munities, analysis of the "purified" set of records may well provide
a clearer picture than if the deviant quadrats are included.

It should be noted that in some cases the list of deviant quadrats
removed might not be the same if pairs had been considered before
single quadrats. In the Queensland data, for instance, quadrats 37
and 38 show many similarities. Since, however, quadrat 37 dif-
fered significantly from the norm of the set, it was rejected as an
individual deviant, and there were then insufficient grounds for
removing quadrat 38. If, on the other hand, at the point where No.
37 was rejected pairs had been considered, it would have been found
that this pair differed significantly from the rest (with four of the 38
species occurring in two of the remaining quadrats), and with a
probability $(6 \times 10^{-6})$ less than that $(6 \times 10^{-4})$ on which the re-
jection of No. 37 as a single quadrat is based.

This suggests that the "purification" of the residual quadrat set
might be more complete if pairs (or perhaps even larger sub-sets)
were considered for rejection before single quadrats.

Although the 17 quadrats removed from the Queensland data
were rejected on account of their deviation from the norm, not
because of any relationship with one another, nevertheless an
examination of species recorded in two or three quadrats show that

they are in fact not unrelated. Of the 57 species occurring in two quadrats, 23 have both their occurrences among the deviant quadrats (if the quadrats in which they occurred were random pairs out of the whole set of 67, the expected number would be 3.6) and of the 30 occurring in three quadrats, 12 have all three of their occurrences among these deviant quadrats (the expected number here is 0.5). Consequently, the removal of these quadrats one by one from the set has resulted unintentionally in the separation of one or more related groups.

Thus, the procedure described can lead to the recognition of distinct communities represented by a minority of the samples in a set, even though this is not its primary purpose. One can, however, envisage an extension with this particularly in mind. As described, the procedure begins by considering individual deviant quadrats, and pairs, as candidates for exclusion. There is no reason why this should not extend to sub-sets larger than two, which would cover the possibility that more than one community should be distinguished in the vegetation sampled.

This extension of the procedure has not yet been worked out in any detail, but the possibilities may be illustrated by reference to the data of DABROWSKI (1956) for beech forest. In the 20 quadrats included in his Table 2 there are 24 species occurring in a single quadrat only; five of these are in Quadrat 16, three each in Quadrats 1, 3, and 18, and the Poisson $\chi^2$ is 32.67 with 19 degrees of freedom (P = .02). For the present purpose, however, these unique species are less helpful than those occurring in two or more quadrats, which give evidence on the similarity in some respect of the quadrats where they are recorded.

Starting with a particular quadrat (perhaps one set apart by the number of single records it contains) one may then compare it with each other in turn, counting the number of species occurring in both and comparing this with the expected number, given the possible pairwise combinations of records for each species. For this purpose, the infrequent species provide the most powerful test.

Quadrat 16, in addition to the five single records mentioned, also includes three species out of the 14 occurring in two quadrats. The other quadrats involved in these pairwise occurrences are 1, 2, and 15. One further notes that another species occurring in two quadrats only (*Melica nutans*) was in Quadrats 1 and 2, while *Vicia* sp. was in 2 and 15. We thus have five out of the 14 species recorded twice only occurring in four out of the 20 quadrats. The chance that any particular quadrat should occur three times among these 28 records is about 0.07, and the chance that the three species occurring with it should also occur with one another is 0.0062.

These four quadrats (1, 2, 15, and 16) are also linked by some of

the other less common species. Of the species recorded in three quadrats, none was recorded in one only of these quadrats, though two were recorded in two of them. Likewise, those with four, five, or six records tended either to be absent from these four quadrats, or to be recorded in more than one of them.

Quadrat 18 included three single records. Of the 11 species recorded in three quadrats, it included four, and all of these four also occurred in Quadrats 19 and 20. The chance that the same three quadrats should be those in which four of the 11 three-record species occur is about $1 \times 10^{-8}$, so DABROWSKI's separation of these three quadrats into Variant C of the association seems fully justified.

This illustration merely serves to show the potentialities of extension of the procedure described to the recognition of distinct communities represented by several quadrats in a set. It is clear that this is by no means the only, or necessarily the best, method of using uncommon species for this purpose, but may encourage the testing of alternative techniques for recognizing and distinguishing clusters of quadrats.

The analogy between these procedures and those of the Zürich-Montpellier school is clear. Where the adherents of this school rearrange their tables of relevés to bring together, on the one hand, vegetation samples which can be grouped into an association, and on the other, a group of species which can be regarded as discriminant or characteristic for this association, they are doing by intuitive methods what we have here been doing more formally. Conceptually, the differences are less important than the resemblances; perhaps the most important point of difference is that the procedure as described here is armed with significance tests based on the null hypothesis that the samples are taken from a common population, and that consequently the occurrence or non-occurrence of each species in any particular sample of the set is random. Only when this null hypothesis has to be abandoned, through a demonstration that the observed distribution of the species is one which would occur only with a probability below the chosen limit, does an alternative hypothesis that more than one community is present become acceptable.

## SUMMARY

In a set of vegetation samples the presence in a single sample of a number of unique (or rare) species may indicate that the sample is in some sense peculiar, and that its removal is likely to render the residue more uniform. By repeated test and removal — a process which can be controlled at a fixed significance level — one may restrict the samples to a uniform subset, within which the presence

or absence of different species may acceptably be ascribed to chance.

This procedure is illustrated with original data for savannah woodland in North Queensland, and with published data for several European vegetation types. It is shown that deviant samples can often be detected in this way, although this may not happen when the floristic differences are mainly quantitative. The procedure is more likely to be successful in floristically rich vegetation than if the number of species present is limited.

The possible development of methods of objective classification based on the same principle is discussed.

## ZUSAMMENFASSUNG

Die Anwesenheit einer Anzahl nur hier vorkommender (oder seltener) Arten in einer Aufnahme, die zu einer ganzen Serie von Vegetationsaufnahmen gehört, kann darauf hinweisen, daß diese Aufnahme irgendwie merkwürdig ist, und daß es wahrscheinlich ist, daß der Rest der Aufnahmen durch ihre Ausschaltung gleichförmiger wird.

Durch wiederholte Untersuchung und Ausschaltung — ein Prozess, der bei einer festgestellten statistischen Signifikanz kontrolliert werden kann — kann man die Aufnahmen auf eine uniforme Untereinheit einschränken, bei der angenommen werden kann, daß die Anwesenheit oder Abwesenheit verschiedener Arten zufällig ist.

Dieses Verfahren wird erläutert durch Originalangaben für Savannenwald in Nord-Queensland, und für veröffentlichte Angaben für verschiedene europäische Vegetationstypen.

Es wird gezeigt, daß oft auf diese Weise abweichende Aufnahmen gefunden werden können, es sei denn, daß die floristischen Unterschiede vor allem quantitativ sind. Das Verfahren dürfte in floristisch reicher Vegetation bessere Resultate geben, als wenn die Anzahl der anwesenden Arten beschränkt ist.

Die mögliche Entwicklung von Methoden objektiver Klassifizierung nach diesem Prinzip wird besprochen.

## ACKNOWLEDGMENTS

## REFERENCES

BRAUN-BLANQUET, J., PALLMANN, H. & BACH, R. 1954 — Pflanzensoziologische und bodenkundliche Untersuchungen im schweizerischen Nationalpark und seinen Nachbargebieten, Vegetation und Böden der Wald- und Zwerg-strauchgesellschaften (Vaccinio-Piceetalia). *Ergebn. wiss. Untersuch. schweiz. Nationalparks*, N.F. 4, 28: 1—200.

DABROWSKI, M. J. 1956 — Rozklad ilosciowy oraz frekwencje gatunkow v warstwie runa (Numerical distribution and occurrence of the species comprising the ground vegetation). *Ekol. Polska*, Ser. A. 4: 349—376.

DAGNELIE, P. 1960 — Contribution à l'étude des communautés végétales par l'analyse factorielle. *Bull. Serv. Carte phytogéogr.*, Ser. B. 5: 7—71, 43—195.

FIJAŁKOWSKI, D. 1961 — Miłek wiosenny (*Adonis vernalis* L.) w województwie lubelskim. (*Adonis vernalis* L. in Wojewodschaft Lublin). *Ann. Univ. Mariae Curie-Skłodowska*, Sect. C. 16: 49—76.

FISHER, R. A. 1963 — Statistical Methods for Research Workers. 13th edition. Oliver & Boyd, Edinburgh and London.

GOODALL, D. W. 1953 — Objective methods for the classification of vegetation. I. The use of positive interspecific correlation. *Aust. J. Bot.* 1: 39—63.

GOODALL, D. W. 1954 — Objective methods for the classification of vegetation. III. An essay in the use of factor analysis. *Aust. J. Bot.* 2: 304—324.

GOODALL, D. W. 1964 — A probabilistic similarity index. *Nature* 203: 1098.

GOODALL, D. W. 1966a — Deviant index — a new tool for numerical taxonomy. *Nature* 210: 216.

GOODALL, D. W. 1966b — A new similarity index based on probability. *Biometrics* 22: 882—907.

GOODALL, D. W. 1968 — The distribution of the matching coefficient. *Biometrics* 23: 647—656.

GOUNOT, M. 1961 — Les méthodes d'inventaire de la végétation. *Bull. Serv. Carte phytogéogr.*, Ser. B. 4: 147—177.

IVIMEY-COOK, R. B. & PROCTOR, M. C. F. 1966 — The application of association-analysis to phytosociology. *J. Ecol.* 54: 179—192.

JACCARD, P. 1901 — Distribution de la flore alpine dans le Bassin des Dranses et dans quelques régions voisines. *Bull. Soc. vaud. Sci. nat.* 37: 241—272.

LAMBERT, J. M. & WILLIAMS, W. T. 1962 — Multivariate methods in plant ecology. IV. Nodal analysis. *J. Ecol.* 50: 775—802.

McVEAN, D. N. & RATCLIFFE, D. A. 1962 — Plant Communities of the Scottish Highlands (A study of Scottish mountain, moorland, and forest vegetation). Her Majesty's Stationery Office, London.

MATUSZKIEWICZ, W. & MATUSZKIEWICZ, A. 1956 — Materiały do fitosocjologiczny systematyki ciepłolubnych dabrow w Polsce (Zur Systematik der Querce-talia pubescentis-Gesellschaften in Polen). *Acta Soc. Bot. Polon.* 25: 27—72.

MOORE, J. J. 1962 — The Braun-Blanquet system: a reassessment. *J. Ecol.* 50: 761—769.

ORLOCI, L. 1966 — Geometric models in ecology. I. The theory and application of some ordination methods. *J. Ecol.* 54: 193—215.

SOKAL, R. R. & SNEATH, P. H. A. 1963 — Principles of Numerical Taxonomy. W. H. Freeman and Company, San Francisco and London.

SZYNAL, T. 1962 — Ogólna analiza florystyczno-ekologiczna zespołów roślinnych Nadleśnictwa Kosobudy na Roztoczu Środkowym. (A general floristic and ecological analysis of plant associations of the forest district Kosobudy in Central Roztocze.) *Ann. Univ. Mariae Curie-Skłodowska* 17: 363—418.

WILLIAMS, W. T. & LAMBERT, J. M. 1959 — Multivariate methods in plant ecology. I. Association analysis in plant communities. *J. Ecol.* 47: 83—101.

WILLIAMS, W. T. & LAMBERT, J. M. 1960 — Multivariate methods in plant ecolo-
gy. II. The use of an electronic digital computer for association analysis.
*J. Ecol.* 48: 689—710.

## APPENDIX

Significance levels for Numbers of Species recorded only once in Two Species
Lists

N.B.  $n_1$  is in each case the larger number of species present in only one of the two
lists,  $n_2$  the smaller number.

Maximum value of  $n_2$  for a significance level of

| $n_1$ | 0.05 | 0.01 | 0.001 |
|---|---|---|---|
| 6 | 0 | — | — |
| 7 | 8 | — | — |
| 8 | 1 | 0 | — |
| 9 | 1 | 0 | — |
| 10 | 2 | 0 | — |
| 11 | 2 | 1 | 0 |
| 12 | 3 | 1 | 0 |
| 13 | 4 | 2 | 0 |
| 14 | 4 | 2 | 1 |
| 15 | 5 | 3 | 1 |
| 16 | 5 | 3 | 1 |
| 17 | 6 | 4 | 2 |
| 18 | 7 | 4 | 2 |
| 19 | 7 | 5 | 3 |
| 20 | 8 | 6 | 3 |
| 21 | 9 | 6 | 4 |
| 22 | 9 | 7 | 4 |
| 23 | 10 | 7 | 5 |
| 24 | 11 | 8 | 5 |
| 25 | 12 | 9 | 6 |
| 26 | 12 | 9 | 6 |
| 27 | 13 | 10 | 7 |
| 28 | 14 | 11 | 7 |
| 29 | 15 | 11 | 8 |
| 30 | 15 | 12 | 8 |
| 31 | 16 | 13 | 9 |
| 32 | 17 | 13 | 10 |