

- [3] P. Billingsley, *Convergence of probability measures*, John Wiley & Sons Inc., New York, London, Sydney, Toronto 1968.
- [4] S. Kaczmarz, H. Steinhaus, *Theorie der Orthogonalreihen*, Chelsea Publishing Company 1951.

*Presented to the semester
 MATHEMATICAL STATISTICS
 September 15–December 18, 1976*

SERiation WITH APPLICATIONS IN PHILOLOGY

L. BONEVA

*Institute of Mathematics and Mechanics, Bulgarian Academy of Sciences,
 Sofia, Bulgaria*

1. Introduction

Seriation has turned out to be a common problem not only in archaeology but in philology and other fields as well. Generally speaking, the basic idea of the recently developed new mathematical, statistical and computing methods connected with seriation was the reconstruction of the “true” chronological order of a set of objects using only the available nonmetric information about the similarities (or dissimilarities) between pairs of objects. These methods have done a good service to all problems dealing with a great amount of data for numerous objects about which only a chronological ordering is needed. We are going to discuss here the SKK-method, which we call so in honour of the names of the three most famous men (Shepard–Kruskal–Kendall) who took part in creating the “main body” of this useful technique.

In fact, the seriation problem was formulated for the first time by the English archaeologist Flinders Petrie [19] at the very end of the last century. He was confronted with a very difficult problem—to find an approximate dating for 4000 prehistoric Egyptian graves, each containing pottery, jewellery and other objects permitting a final classification into types of varieties. Evidently, a chronological trend of these types is to be expected according to which the approximate dating of the graves might be done. Actually, Petrie managed to arrange 900 graves containing a total amount of 800 varieties. The weakest point of his laborious work is the “reverse connection” between graves and varieties, i.e. the varieties were classified according to the graves in which they were found, while the graves were ordered according to the varieties they contained. However, he is to be thanked for the so-called “Petrie’s Concentration Principle”, which shortly states that the more close together in temporal order two graves are the more likely they are to contain varieties of the same or similar types.

A second merit of Petrie’s work should not be omitted. It is he who gave the initial impulse (though it resounded about 50 years later) to many mathematicians, such as Robinson [20], Shepard [23], Kruskal [17], [18], Kendall [10]–[16], Sibson

[21], [22], Wilkinson [24], etc., for developing new techniques or for refining old ones.

The availability of high-speed computers has given lately a powerful push to the seriation problem not only in archaeology but also in psychology, in history, in classical and modern philology, in the reconstruction of maps, etc. A considerable amount of literature on this subject is contained and referred to in the excellent Edinburgh University Press edition of the *Mamaia Proceedings* [8].

2. Mathematical formulations and algorithms

The mathematical approach is mostly due to D. G. Kendall and could be divided into cases A and B.

Suppose we have n objects and k varieties, i.e. we have n k -dimensional vectors or an $(n \times k)$ -matrix $A = \{a_{ij}\}$ which has as many rows as there are objects and as many columns as there are varieties. The problem could be formulated shortly with several definitions and theorems.

Case A. Seriation from incidence matrices

The following four definitions could be summarized here:

D.1. We call A an *incidence matrix* when

$$a_{ij} = \begin{cases} 1 & \text{if the } i\text{th object contains the } j\text{th variety,} \\ 0 & \text{if the } i\text{th object does not contain the } j\text{th variety.} \end{cases}$$

D.2. We call A *Petrie* (or a *matrix with Pattern P*) if in each column there is only one sequence of consecutive 1's, provided such a sequence does exist.

D.3. We call A *petrifiable* if there exists a permutation matrix π such that πA is *Petrie* (Pattern P).

D.4. We say that a square symmetric matrix is in the *Robinson form* (Pattern R) if, when going to the left or down from any position of the main diagonal, the elements never increase.

It has been proved firstly that to decide whether there exists a row permutation that will bunch together the 1's in each column of A it is enough to know $V = A'A$ and, secondly, that if A has such a property then it could be rearranged if we knew $G = AA'$, thus, $V = A'A$ and $G = AA'$ contain, respectively: (i) information about the possibility of such rearrangement of A , and (ii) sufficient information for constructing a sorting algorithm, provided A is petrifiable.

Now we could formulate, in the language of D.1, ..., D.4, the two main theorems in this case, the first due to Fulkerson and Gross [7], and the second due to Kendall [11].

THEOREM (F and G). *If A and B are two incidence matrices with the same number of rows and columns, and if $V = A'A = B'B$, then B is Petrie if and only if A is Petrie.*

Consequently, the knowledge of the matrix V suffices for answering the question whether A is petrifiable or not. There is a nicely formulated graph technique in [7] permitting us to identify a petrifiable matrix A only from its V matrix.

Evidently, a more important seriation question is whether there exists a possibility of applying a sorting algorithm. To answer it Kendall showed that it is enough to know $G = AA'$ (called *similarity matrix* with elements, s_{ij} , the experimentally obtained similarities between objects i and j) by proving the following

THEOREM (K.a). *If an incidence matrix A is petrifiable, then the same row permutations which petrify A will, when applied both to rows and to columns of $G = AA'$, turn G into Pattern R.*

So, $A'A$ shows us whether there exists a seriation solution, while AA' gives us all the relevant information for finding it.

Case B. Seriation from abundance matrices

Here the four definitions differ as follows:

1. If A is again an $(n \times k)$ -matrix but with arbitrary real numbers a_{ij} (or frequencies, say p_{ij} , for which $\sum_{j=1}^k p_{ij} = 1$) as elements, we call A an *abundance matrix*.

2. Instead of Pattern P we deal with Pattern Q here, i.e. we say that an abundance matrix A for which there exists at least one row permutation π turning πA into

Pattern Q: the elements of each column are unimodal functions of i for each j , is liable to seriation.

3. Instead of *petrifiable* we say *queutrifable matrices*.

4. The similarity matrix, which has to have Pattern R, is $S = A \circ A'$ here with elements

$$s_{ij} = (A \circ A')_{ij} = \sum_k w_k \min(a_{ik}, a_{jk}), \quad \text{where } w_k > 0 \text{ but arbitrary.}$$

It is shown by Wilkinson [24] that no relevant information is lost if we deal with the "promoted" matrix $S \circ S$ with elements

$$(S \circ S)_{ij} = \sum_k w_k \min(s_{ik}, s_{jk})$$

as well as with some of its "degrees" $(S \circ S) \circ (S \circ S)$, etc.

Kendall's theorem could now be formulated as

THEOREM (K.b). *If an abundance matrix A is queutrifable, then the same row permutations which queutrifify A will, when applied both to rows and to columns of $S = A \circ A'$, turn S into Pattern R.*

We know so far that from V we can learn whether there exists a solution; but if so, how to find it?

A natural approach will be to rearrange the rows and columns of the similarity matrix until we get Pattern R. But this would be practically impossible if n is large enough (remember Petrie's case).

Then other seriation techniques could be adopted, such as MDSCAL, proposed by Kruskal, or HORSHU, proposed by Kendall. Both procedures use G or S , or some of their degrees, as an input and produce a t -dimensional configuration, $t \geq 2$, as an output. We shall call the joint technique a *KK seriation algorithm* because it could be said that MDSCAL + ($S = A \circ A'$) = HORSHU.

This KK-algorithm, which has been worked out and continuously improved during the last 15 years, could be schematized as a geometrical representation of n objects by n points. More precisely, it consists of finding n points in a t -dimensional space (a configuration) in such a way that the distances between the pairs of points correspond (in some sense) to the similarities (dissimilarities) between the pairs of objects. The basic hypothesis is that distances and similarities are monotonically related. Some more details now:

Suppose that n and $t \geq 2$ are fixed. Let us call the arbitrarily chosen points P_1, \dots, P_n . Let d_{ij} denote the distance from P_i to P_j and let $(x_{i1}, x_{i2}, \dots, x_{it})$ be the orthogonal coordinates of P_i . We can use for the distances the l -metric

$$d_{ij} = \left(\sum_{s=1}^t |x_{is} - x_{js}|^r \right)^{1/r} \quad \text{where } r \geq 1,$$

but we mostly use the Euclidean metric, e.g. when $r = 2$. To evaluate how well the distances match the similarities we produce first of all a scatter diagram of all (i, j) points with coordinates (s_{ij}, d_{ij}) for all $i < j$, $i = 1, \dots, n-1$ and $j = 2, \dots, n$. At each step of the algorithm, i.e. for any configuration, we perform a monotone regression of similarity upon distance. The best configuration (the output) is supposed to represent the i th object with the i th point in such a way that for two different points, say (i, j) and (l, k) , we have $s_{ij} > s_{lk}$ whenever $P_i P_j < P_k P_l$, i.e. the smaller the distances the more similar the objects. If the objects are liable to seriation, then intuition tells us that the output-points should be on a straight line but, as Kendall has shown experimentally, they are plotted along a horseshoe curve (or something like an arched Milky Way, when "noise" is available) which, as proved by Wilkinson, has turned out to be the shortest Hamiltonian circuit (or the shortest path in the travelling salesman problem). But let us go on with the rough description of the algorithm. If we ignore s_{il} and suppose that $s_{ij} = s_{ji}$, and also that there are no ties (equal similarities), then we have $M = n(n-1)/2$ similarities, which we could arrange in a strictly ascending (descending) order $s_{i_1 j_1} < \dots < s_{i_M j_M}$. Usually we do not expect the corresponding M distances to be in ascending (descending) order because the starting configuration is absolutely arbitrary. The aim is to find numbers d_{ij}^* which are supposed to be as near to the d_{ij} as possible and monotonically related to them, i.e. satisfying the (Mon) condition

$$d_{i_1 j_1}^* \leq d_{i_2 j_2}^* \leq \dots \leq d_{i_M j_M}^*.$$

In fact, d_{ij}^* are the ordinates of the M points when we move them, (say, with step-size α) so as to make the curve which is passing through the M new points with coordinates (s_{ij}, d_{ij}^*) a monotone ascending one. The KK-algorithm enables us to find the best configuration with a reasonably small number of iterations (usually no more than 50), each time finding a new d_{ij}^* and a new monotonically ascending curve belonging to a family of "nonparametric" curves. A natural measure for the "goodness of fit" seems to be the number

$$\text{stress } S = \sqrt{\frac{\sum_{i < j} (d_{ij} - d_{ij}^*)^2}{\sum_{i < j} (d_{ij} - \bar{d})^2}} \geq 0, \quad \text{where } \bar{d} = \frac{1}{M} \sum_{i < j} d_{ij}.$$

For finding the best fitting configuration we have to minimize S in two ways: at each step, fitting each time new d_{ij}^* 's, i.e.

$$\text{stress of fixed configuration} = S(P_1, \dots, P_n) = \min_{d_{i_1 j_1}^* \leq \dots \leq d_{i_M j_M}^*} S,$$

and as a whole, i.e.

$$\text{stress in } t\text{-dimensions} = \min_{\text{all } t\text{-dimensional configurations}} S(P_1, \dots, P_n).$$

For this purpose the well-known numerical method of the steepest descent was used. It should be noted here that to ensure that the resulting minimum is not only a local one (even if it happens to be very close to zero, say ≤ 0.025) it is advisable to repeat the whole procedure several times, starting from different initial configuration each time.

And one last remark. If ties are available, then we have two possibilities: (i) Primary treatment of ties (PTT) — we do not care if $s_{ij} = s_{kl}$ and deal as if there were no ties, and (ii) secondary treatment of ties (STT) — we accept that whenever $s_{ij} = s_{kl}$ then d_{ij} must be equal to d_{kl} , and we diminish the configuration even if $d_{ij} \neq d_{kl}$. Therefore we put the additional restriction $d_{ij}^* = d_{kl}^*$ whenever $s_{ij} = s_{kl}$. PTT and STT should be taken into account when dealing with the algorithm for fitting d_{ij}^* [17].

Up to now there have been some additional facilities for improving the final configuration like CIRCLEUP (which replaces S by $S \circ S$) or an option which permits a choice between PTT and STT, etc., proposed by Kendall. Of course, there are other algorithms, proposed by other authors, but we would like to refer to [8] once more.

3. Applications in philology

We came to the idea of applying SKK in classical philology when looking through a work of Cox and Brandwood [6]. They give a linear model for finding the chronological order of Plato's last 6 dialogues with the help of the maximum likelihood test. About 10 years later Atkinson [1] attacked them, challenging the linearity of their model, dealing again with the same 6 works of Plato.

However, Plato has written 45 different dialogues and it seemed worthwhile to try to arrange all of them, considering the "Republic" as 10 and the "Laws" as 12 separate books (because it was so in fact). For more than 110 years many philologists and philosophers argued frantically about the chronology of parts of them. Some managed to arrange 2 or 3 but none more than 6 books. For ordering all the 45 dialogues we used their clausula frequency distributions, done in 1904 by the German philologist [9] Kaluscha. The clausula chosen by him consists of 5 syllables of two types — long and short one — or altogether $2^5 = 32$ classes of clausulae. According to our experience the 5 syllable clausula seems to be the most convenient one. The results of applying SKK to all the 45 of Plato's works may be seen in [2].

Later on the same seriation technique but a slightly different phonetico-synthetic "clausula criterion" (the clausulae consisting again of 5 syllables but this time of the stressed and the unstressed types) have been applied to obtain the chronology of the works (but with a known "true" order now) of a modern author. Briefly, a sample of 16 short stories, taken at random from the total number of short stories written by the Bulgarian novelist Y. Yovkov, have been studied. To our great satisfaction the order obtained happened to be the same as the true one. The results may be seen in [3]. Here we give only the coordinates of the final configuration (Table 1) because that has been omitted in [3].

To work in Sofia and to use a computer in Cambridge turned out to be quite a difficult task. That is why an algorithm for "seriation at hand" when n is small

Table 1

No	Final configuration		Order of publishing	Order of writing	True order obtained
1	0.214	-1.606	1	1	1
2	-1.322	-0.904	2	2	2
3	-0.894	-0.542	3	4	4
4	-0.660	1.921	4	7	7
5	-0.061	0.714	5	8	8
6	-0.643	0.572	6	5	5
7	0.436	0.997	7	10	10
8	-0.068	0.398	8	9	9
9	0.992	-0.681	9	12	12
10	0.905	0.232	10	11	11
11	-1.347	-0.086	11	3	3
12	2.487	-1.596	12	13/14	14
13	1.135	0.687	13	13/14	13
14	-1.593	0.976	14	6	6
15	0.324	-0.830	15	16	16
16	0.096	-0.253	16	15	15

enough, say $n < 10$, has been published [4] in the meantime. Of course the algorithm, called *LMC* (*local-maximum-chain*), works only if the starting matrix A is queutrififiable. So, instead of rearranging the rows and columns of S we may try *LMC*. If this does not work, then a Pattern R is not to be expected in S .

Table 2

No	Final configuration		Order of publishing	Order of writing	Order obtained
1	-0.758	1.156	2	2	2
2	0.331	0.973	3	3	3
3	0.768	0.967	6	5	5
4	1.012	-0.998	10	11	11
5	0.326	-0.352	12	14	14
6	0.641	-0.007	9	9	9
7	-0.121	-0.697	13	17	17
8	0.160	0.048	14	15	15
9	-0.657	-0.298	15	19	19
10	-0.413	-0.337	16	18	18
11	0.468	0.861	17	4	4
12	-1.970	-0.490	20	21/22	21
13	-0.944	0.067	21	21/22	22
14	1.320	-0.583	24	10	10
15	-0.285	0.301	25	26	26
16	-0.040	0.187	26	25	25
17	0.711	0.191	5	7	7
18	1.472	0.527	7	6	6
19	0.057	-0.249	8	16	16
20	0.606	0.408	4	8	8
21	-0.429	-1.473	18	12	12
22	-0.555	0.323	11	23	23
23	-0.223	0.588	1	1	1
24	0.182	-0.866	22	13	13
25	-0.050	-0.162	23	24	24
26	-0.385	-0.969	19	20	20
27	-1.073	0.482	27	27	27
28	-0.650	0.839	28	28	28

Not long ago a FORTRAN version of HORSHU was also prepared for our computer. The first thing to do was to try to arrange a bigger sample of Yovkov's stories ($n = 28$, $k = 32$, and the same "clausula criterion"). Note that two of the stories were not published at all (No 23 and No 27) but we took them just for comparison with their published "twins", No 6 and No 28. The results are given in Table 2, where the first column gives the input order, the second gives the coordinates of the final configuration, the third gives their order of publishing, the fourth gives the order in which, according to our philologists, they were written, the fifth gives the order obtained. An illustration of one of the several final configurations (similar

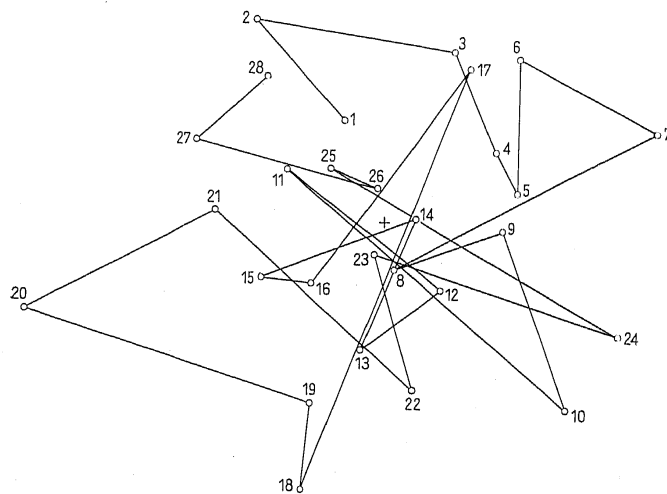


Fig. 1a

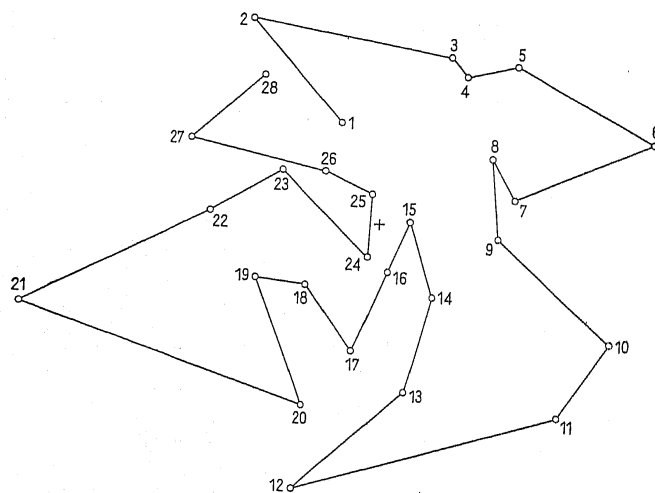


Fig. 1b

to rotation) that have been obtained is shown on Fig. 1. One thing to be noticed here is that in connecting the stories in their order of publishing the lines crossed each other in all the resulting final configurations, while, in connecting them in their writing order, the lines never crossed, and Figs. 1a and 1b actually reflect that result. In fact, the "horse-shoe" shape is a rather scrambled one. One of the reasons for that might be the great number of zero's among the frequencies of our distributions. Anyway, we hope to be able to change this shape — unsatisfactory for a "horse-shoe" — very soon either by diminishing the zero's or by combining similar classes of clausula or else by working out a new more appropriate for our case similarity measure. Meanwhile we are trying to straighten the configuration by using CIRCLE-UP, PTT, STT, and combinations of them. This is worth trying because Fig. 1 shows a good sense in the chronology obtained. There are five distinguishable groups (see Fig. 1.b): (1) The stories from 2 up to 8 (including even 1, or No 23) were written during the first world war; (2) The stories from 9 up to 13 are "peculiar" in some way. They had been written, rewritten and seriously changed before being published, e.g. 9, which is the "twin" of 1, or No 23, was written in 1910, rewritten several times, but published only in 1926. Consequently, it comes immediately after the war-period stories; (3) The intermediate period of publishing, actually 1925–1928 (the author published between 1913 and 1937) consists of the stories from 14 to 19; (4) Next comes the 1931–1935 period, consisting of the stories from 20 to 24. The only exception here is the story 22 published first in 1925 and then in 1927, but rewritten (according to the author himself) in such a way as to be different in style from the former version. He claimed the same for 11, which actually comes in the "peculiar" group; (5) Here only the works written in 1936 occur, including the one written twice 27 and 28 but published once (the second version).

Finishing, we repeat our hopes that after some modifications we will be able to obtain a better final configuration, because we do believe that the changes of one's style go "hand in hand" with the temporal changes of one's personality and, what is more, that they could be traced chiefly in sentence endings.

References

- [1] A. C. Atkinson, *A method for discriminating between models*, J. Roy. Statist. Soc. (B) 32 (1970), pp. 323–353.
- [2] L. Boneva, *A new approach to a problem of chronological seriation associated with the works of Plato*, in [8] (1971), pp. 173–185.
- [3] —, *Chronological seriation applied in literature*, Proc. of the 39th Session of Int. Statist. Inst. (1973), pp. 199–209.
- [4] —, *Chronological seriation of the works of an author by means of computer*, Proc. of the First Bulg. Conf. on the Application of Math. Models in Linguistics (1975), pp. 319–322.
- [5] R. M. Clark, *A survey of statistical problems in archaeological dating*, J. of Multivar. Anal. (4) 3 (1977), pp. 308–326.
- [6] D. R. Cox and L. Brandwood, *On a discriminatory problem connected with the works of Plato*, J. Roy. Statist. Soc. (B) 21 (1959), pp. 195–200.

- [7] D. R. Fulkerson and O. A. Gross, *Incidence matrices and interval graphs*, Pacific J. Math. 15 (1965), pp. 835–855.
- [8] F. R. Hodson, D. G. Kendall and P. Tautu (eds.), *Proc. of the Anglo-Romanian Conference of Mathematics in the Archaeological and Historical Sciences*, Edinburgh University Press, Edinburgh 1971.
- [9] W. Kaluscha, *Zur Chronologie der Platonischen Dialoge*, Wiener Studien 25, 26 and 27 (1904).
- [10] D. G. Kendall, *A statistical approach to Flinders Petrie's sequence dating*, Bull. Int. Statist. Inst. 40 (1963), pp. 657–680.
- [11] —, *Incidence matrices, interval graphs, and seriation in archaeology*, Pacific J. Math. 28 (1969), pp. 565–570.
- [12] —, *Some problems and methods in statistical archaeology*, World Archaeology 1 (1969), pp. 68–76.
- [13] —, *Abundance matrices and seriation in archaeology*, Zeitschrift f. Wahrscheinlichkeitstheorie 17 (1971), pp. 104–112.
- [14] —, *A mathematical approach to seriation*, Philos. Trans. Roy. Soc. London (A) 269 (1971), pp. 125–134.
- [15] —, *Seriation from abundance matrices*, in [8] (1971), pp. 213–254.
- [16] —, *Some data-analytic problems in archaeology and history*, Sci. and Hum. North Holland Publ. Comp. (1972), pp. 1371–1376.
- [17] J. B. Kruskal, *Multidimensional scaling, I and II*, Psychometrika 29 (1964), pp. 1–27 and 28–42.
- [18] —, *Multidimensional scaling in archaeology*, in [8] (1971), pp. 119–132.
- [19] W. M. F. Petrie, *Sequences in prehistoric remains*, J. Anthropol. Inst. 29 (1899), pp. 259–301.
- [20] W. S. Robinson, *A method for chronologically ordering archaeological deposits*, American Antiquity 16 (1951), pp. 293–301.
- [21] R. Sibson, *Some thoughts on sequencing methods*, in [8] (1971), pp. 263–266.
- [22] —, *Local order multidimensional scaling*, Proc. 39 Session of I.S.I. (1973) (inv. paper).
- [23] R. N. Shepard, *The analysis of proximities: Multidimensional scaling with an unknown distance function, I and II*, Psychometrika 27 (1962), pp. 125–139 and 219–246.
- [24] E. M. Wilkinson, *Archaeological seriation and the travelling salesman problem*, in [8] (1971), pp. 276–283.

Presented to the semester
 MATHEMATICAL STATISTICS
 September 15–December 18, 1976

USE OF MATRIX APPROXIMATION IN STATISTICS

L. C. A. CORSTEN

Department of Mathematics, Agricultural University, Wageningen, The Netherlands

In trying to approximate an n by p array Y of data by a matrix C of rank k , one may want to minimize the approximation error matrix E in some sense. Minimization of $|Ex|/|x|$ for all x , or of $|E'y|/|y|$ for all y suggests minimization of all eigenvalues of $E'E$ or EE' simultaneously. This minimization can be attained by the canonical decomposition YAU' of Y , where $U'U = V'V = I_r$, r is the rank of Y , and λ_j is the positive square root of the j th largest characteristic value of $Y'Y$ or YY' ($j = 1, \dots, r$).

The required approximation C of $Y = YAU' = \sum_{j=1}^r \lambda_j v_{*j} u'_{*j}$ obtained by suppressing the last $r-k$ terms in this sum equals $V_k A_k U'_k$. In this way, $Y'Y = U A^2 U'$ will be approximated by $(U_k A_k)(U_k A_k)'$, and any symmetric matrix S by $U_k A_k^* U'_k$ where A_k^* contains k characteristic values of S in non-increasing order of their absolute value. The approximation C of Y may be written as AB' where $A = V_k = YU_k A_k^{-1}$ and $B = U_k A_k$, and the rank k approximation of $Y'Y$ equals BB' . When each row of Y contains a multivariate observation at a corresponding individual and each column corresponds to a component of such a multivariate vector, each column a_{*j} of A may be conceived of as the set of n values of a new characteristic (a factor), each row a_{i*} of A as a set of factor scores for the i th individual, and each row b_{j*} of B as the set of factor loadings for the j th component of the multivariate observations.

As the columns of A are orthonormal the structure of the columns of C approximating those of Y may be visualized by means of their coordinates b_{j*} in k -space, the inner products between those columns approximating those between y_{*j} . Each row a_{i*} of factor scores may, likewise in k -space, visualize the mutual position of the individuals, and the inner product between a_{i*} and b_{j*} is the approximation of c_{ij} .

Approximation of Y by a rank k matrix C plus $1\beta'_1$ and (or) $\beta_2 1'$ where β_1 is a p -vector and β_2 an n -vector is a useful modification of the first situation. Not only the situation that $n^{-1}Y'Y$ is a covariance matrix \mathcal{Z} is covered now, but it also leads to an exact test on the presence of a multiplicative term in a two-way analysis of variance table in addition possibly to row effect and (or) a column effect.