

F. A. SZCZOTKA (Warszawa)

## ON A METHOD OF ORDERING AND CLUSTERING OF OBJECTS

### 1. THE PROBLEM

Given  $p$  objects  $X_1, X_2, \dots, X_p$ , such as persons, teams, phenomena, characteristics etc., two problems are often encountered:

- (1) to order the objects, and
- (2) to divide them into groups of similar objects.

To solve these problems, a criterion is needed which would decide in a unique way which of two orderings or clusterings is the better one. Such a criterion determines also the best ordering or clustering. The number of objects is finite, thus the number of possible orderings or clusterings is also finite. Therefore it is possible, at least theoretically, to find an optimum solution. The number of possibilities is, however, vast and solving the problem by enumeration is practically impossible even with high-speed computers. That is the reason why, having determined an optimality criterion, one is seeking an algorithm which would give the optimum or near-optimum solution in a fast way. Often the optimality criterion is not formulated explicitly; it is determined by the proposed solution methods which — of course — should be objective and unique ones.

Usually, the mathematician encounters ordering and clustering problems only if the optimality criterion is not implied directly by the problem considered or if more general criteria are wanted. Similarity measures between objects are introduced then and on their basis optimality criteria are formulated.

Objects characterized by the standardized values of  $m$  quantitative characteristics are treated as points in  $m$ -dimensional Euclidean space whose coordinates are equal to the values of the characteristics, and similarity is determined by the distance between these points. Similarity between characteristics is measured by correlation coefficients or their squares. Similarity between plant groups may be determined by the index of similarity introduced by Marczewski and Steinhaus [3]. Other measures are, of course, also possible.

The present paper is composed of two parts. In the first part an optimality criterion for linear ordering is proposed and an algorithm allowing to find a near-optimum solution is given. The second part of the paper contains a discussion of optimality criteria for optimal clustering of characteristics and a method of clustering based on the optimal ordering is proposed.

## 2. ORDERING

**2.1. A goodness-of-ordering criterion.** The ordering of objects will be determined by giving a relation of neighbourhood for objects. The notation for two objects  $X_i$  and  $X_j$  being neighbours will be as follows:  $X_i - X_j$ . This relation has to satisfy two conditions:

- (a)  $X_i - X_j \Rightarrow X_j - X_i$ ,
- (b) every object has at least one neighbour.

The relation of neighbourhood may not be defined for some pairs of objects.

An ordering will be called *linear* if the relation  $-$  satisfies two additional conditions:

- (c) every object has at most two neighbours,
- (d) there exist exactly two objects, say  $X_1$  and  $X_p$ , which have only one neighbour each.

A relation satisfying conditions (a)-(d) ranges the objects either as  $X_1, \dots, X_p$  or as  $X_p, \dots, X_1$ . Both these orderings will be considered identical.

An example of non-linear ordering is given by the minimum-spanning-tree ordering [1].

The linear ordering  $X_{a_1}, X_{a_2}, \dots, X_{a_p}$  is determined uniquely by the permutation of the indices  $\mathbf{a} = (a_1, a_2, \dots, a_p)$ .

Assume that the distance between every pair of objects  $d(X_i, X_j) = d_{ij} \geq 0$  is defined (in any permissible way). By assumption, these distances satisfy the conditions

- (i)  $d_{ii} = 0$ ,
- (ii)  $d_{ij} = d_{ji}$ ,
- (iii)  $d_{ij} + d_{jk} \geq d_{ik}$ .

The distance matrix will be denoted by  $D(a_1, a_2, \dots, a_p) = (d_{a_i a_j})$ . Thus, every ordering has a different distance matrix. The set of elements is, however, for all matrices identical.

Let us introduce the notion of object apartness in an ordering. Neighbouring objects in an ordering will be said to be *apart 1 unit*, every

second object is from every other second object *apart 2 units*, etc. Generally, in the ordering  $a_1, a_2, \dots, a_p$  the objects numbered  $a_i$  and  $a_j$  will be said to be *apart*  $\delta_{ij} = |i - j|$  units. The apartness matrix  $\Delta$  is identical for all possible orderings and has the form

$$\Delta = \begin{pmatrix} 0 & 1 & 2 & 3 & \dots & p-1 \\ 1 & 0 & 1 & 2 & \dots & p-2 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ p-1 & \dots & \dots & 1 & 0 & \dots \end{pmatrix}.$$

Consider the function

$$(1) \quad Q^*(a_1, a_2, \dots, a_p) = \sum_{i < j} \delta_{ij} d_{a_i a_j} = \sum_{v=1}^{p-1} v \sum_{s=1}^{p-v} d_{a_s, a_{s+v}}.$$

Function (1) is the weighted sum of distances between every pair of objects, where the weights are given by the apartness of the objects in the ordering. The weights of distances of objects being near in the ordering are small, those of distant objects are great. Intuitively, one wants such an ordering in which objects with small distances would be neighbours, i.e. small distances would have small weights, and objects with great distances would be far apart, i.e. great distances would have great weights. The suggestion is that a better ordering should have a greater value of  $Q^*$ ; thus, the best ordering will be that one for which the function  $Q^*$  reaches its maximum. Therefore one may formulate the following criterion of goodness-of-ordering:

CRITERION  $Q^*$ . *The ordering  $(a_{11}, a_{12}, \dots, a_{1p})$  is better than the ordering  $(a_{21}, a_{22}, \dots, a_{2p})$  if*

$$Q^*(a_{11}, a_{12}, \dots, a_{1p}) > Q^*(a_{21}, a_{22}, \dots, a_{2p}).$$

Such a criterion makes it possible to test which of two given orderings is the better one.

The proposed criterion may be generalized by considering the function

$$Q_f^*(a_1, a_2, \dots, a_p) = \sum_{v=1}^{p-1} f(v) \sum_{s=1}^{p-v} d_{a_s, a_{s+v}},$$

where  $f(v)$  is monotone increasing in  $v$ .

Instead of the distance, a measure of similarity between objects is often defined. If the objects are random variables or characteristics, such a measure may be given by the correlation coefficients  $r_{ij}$  between them. In such a case it is assumed that  $r_{ij} \geq 0$  for all  $i, j = 1, 2, \dots, p$ . The goodness of ordering may now be defined by the function

$$(2) \quad Q(a_1, a_2, \dots, a_p) = \sum_{v=1}^{p-1} v \sum_{s=1}^{p-v} r_{a_s, a_{s+v}},$$

under the condition that goodness-of-ordering is now given by the following

**CRITERION  $Q$ .** *The ordering  $(a_{11}, a_{12}, \dots, a_{1p})$  is better than the ordering  $(a_{21}, a_{22}, \dots, a_{2p})$  if*

$$Q(a_{11}, a_{12}, \dots, a_{1p}) \leq Q(a_{21}, a_{22}, \dots, a_{2p}).$$

Under this criterion the best ordering is that one for which the function  $Q$  reaches its minimum.

**2.2. Quasi-optimal ordering.** The ordering  $(a_1, a_2, \dots, a_p)$  will be called *quasi-optimal* if no better one can be obtained by transposition of two objects.

**Definition.** The ordering  $(a_1, a_2, \dots, a_p)$  is called  $Q$ - (or  $Q^*$ -) *quasi-optimal* if  $Q(a_1, a_2, \dots, a_i, \dots, a_j, \dots, a_p) \leq Q(a_1, \dots, a_j, \dots, a_i, \dots, a_p)$  for every pair  $i, j = 1, 2, \dots, p$  (or if  $Q^*(a_1, \dots, a_i, \dots, a_j, \dots, a_p) \geq Q^*(a_1, \dots, a_j, \dots, a_i, \dots, a_p)$ , respectively).

It follows from the definition that an optimal ordering is also a quasi-optimal ordering. The opposite theorem is not true.

**2.3. Algorithm for quasi-optimal ordering.** Let the objects be ordered in the sequence  $(1, 2, \dots, s, \dots, t, \dots, p)$ . The value of the function  $Q$  will then be  $Q(1, 2, \dots, s, \dots, t, \dots, p)$ . A transposition of objects numbered  $s$  and  $t$  leads to the ordering  $(1, 2, \dots, t, \dots, s, \dots, p)$  and to  $Q(1, 2, \dots, t, \dots, s, \dots, p)$ . The quantity

$$q(s, t) = Q(1, 2, \dots, s, \dots, t, \dots, p) - Q(1, 2, \dots, t, \dots, s, \dots, p)$$

denotes the change of criterion  $Q$  caused by a transposition of objects  $X_s$  and  $X_t$ . Because of

$$Q(1, 2, \dots, t, \dots, s, \dots, p) = Q(1, 2, \dots, s, \dots, t, \dots, p) - q(s, t),$$

one obtains a better ordering by transposition of  $X_s$  and  $X_t$  if  $q(s, t) > 0$ .

The value of  $q(s, t)$  is easy to calculate even without a computer. It equals

$$(3) \quad q(s, t) = -h \sum_{v=1-s}^{-1} (r_{s, s+v} - r_{s+h, s+v}) + \sum_{v=1}^{h-1} (-h + 2v)(r_{s, s+v} - r_{s+h, s+v}) + h \sum_{v=h+1}^{p-s} (r_{s, s+v} - r_{s+h, s+v}),$$

where  $h = t - s$  and  $t > s$ .

The following formula may be of use. If  $s < t < u < w$ , we have

$$\begin{aligned} & Q(1, \dots, t, \dots, s, \dots, w, \dots, u, \dots, p) \\ & = Q(1, \dots, s, \dots, t, \dots, u, \dots, w, \dots, p) - q(s, t) - q(u, w), \end{aligned}$$

where  $q(s, t)$  and  $q(u, w)$  are calculated after formula (3).

From the definition of quasi-optimal ordering it follows directly

**COROLLARY.** *If  $(a_1, \dots, a_p)$  is a quasi-optimal ordering,  $q(a_s, a_t) \leq 0$  holds for every pair of  $s, t = 1, 2, \dots, p, s < t$ .*

Now we may proceed to formulate the algorithm for finding a quasi-optimal ordering.

Let us have the ordering  $\mathbf{a}_k = (a_{k1}, a_{k2}, \dots, a_{kp})$ .

1. Calculate  $q(a_{ks}, a_{kt})$  for all  $s, t = 1, 2, \dots, p, s < t$ .

If for all pairs  $s$  and  $t, s < t$ , holds  $q(a_{ks}, a_{kt}) \leq 0$ , the ordering  $\mathbf{a}_k$  is quasi-optimal. Go then to 5. If there exists a  $q(a_{ks}, a_{kt}) > 0$ , go to 2.

2. Find  $q(a_{ks}^*, a_{kt}^*) = \max_{s,t} q(a_{ks}, a_{kt})$ .

3. Transpose the objects numbered  $a_{ks}^*$  and  $a_{kt}^*$ .

4. Denote the new permutation of indices by  $\mathbf{a}_{k+1} = (a_{k+1,1}, a_{k+1,2}, \dots, a_{k+1,p})$  and go back to 1.

5. Calculate the value of  $Q(a_{K1}, a_{K2}, \dots, a_{Kp})$ .

Since  $Q(\mathbf{a}_k) > Q(\mathbf{a}_{k+1})$  for every  $k$ , the procedure ends with an ordering  $\mathbf{a}_K = (a_{K1}, a_{K2}, \dots, a_{Kp})$  such that  $q(a_{Ks}, a_{Kt}) < 0$  for all  $s < t$ . That ordering is a quasi-optimal one.

**2.4. An example.** The presented algorithm has been used to find the quasi-optimal ordering of 14 characteristics of motor ability. (The data were kindly provided by Mr. W. Czworonóg from the Academy of Physical Education, Warsaw). The characteristics were the following: 1. 100 m run, 2. shot-put, 3. high jump, 4. 30 m run (velocity test), 5. 250 m run (endurance test), 6. vertical jump, 7. standing broad jump, 8. triple jump, 9. knee bends, 10. weight lifting, 11. 60 m hurdles, 12. discus throwing, 13. javelin throwing, 14. broad jump.

The matrix of sample correlation coefficients is given in Table 1. The algorithm leads to the following quasi-optimal ordering:

1. 250 m run (endurance test), 2. 100 m run, 3. 60 m hurdles, 4. 30 m run (velocity test), 5. broad jump, 6. vertical jump, 7. triple jump, 8. standing broad jump, 9. high jump, 10. knee bends, 11. shot-put, 12. weight lifting, 13. discus throwing, 14. javelin throwing.

The corresponding matrix of correlation coefficients is given in Table 2. The value of  $Q$  is equal to 193,860.

### 3. CLUSTERING

**3.1. The problem.** The clustering problem may be stated as follows. Given  $p$  objects, they are to divide into  $k$  clusters in such a way as to have some optimality conditions satisfied. Intuitively, the clustering is a good one if the objects within clusters are similar and non-similar objects belong to different clusters.

Let us treat the objects as points in Euclidean space of appropriate dimension and let us know the distances between them. Denote by  $X_{ij}$ ,  $j = 1, 2, \dots, p_i$ , the point  $j$  in cluster  $i$ , and by  $S_i$  the centroid of this cluster. Let  $d(X_{ij}, S_i)$  denote the distance of point  $X_{ij}$  from  $S_i$ . A possible optimality criterion for clustering is the criterion of minimum sum of squares of distances of the points from their centroid.

CRITERION *K*. Among all possible divisions of the points  $X_1, \dots, X_p$  into  $k$  disjoint non-empty clusters the optimum one is that for which the quantity

$$K = \sum_{i=1}^k \sum_{j=1}^{p_i} d^2(X_{ij}, S_i)$$

reaches its minimum.

If characteristics form the objects which are to be divided into clusters, one uses an optimality criterion based on Wilks' statistics which is used to verify the hypothesis that  $k$  groups of random variables with normal distributions are mutually uncorrelated [6]. Let  $C_i$  denote the matrix of sums of products of deviations for the characteristics which belong to group  $i$ , and  $C$  the appropriate matrix for all characteristics. Wilks' statistics is given by the formula

$$(4) \quad W = |C| / \prod_{i=1}^k |C_i|.$$

If every characteristics of group  $i$  is uncorrelated with every characteristics of group  $j$ , and if this condition is satisfied for all pairs of indices such that  $i \neq j$ , then  $|C| = \prod_{i=1}^k |C_i|$  and  $W = 1$ . In all other cases  $W$  is smaller than 1. The optimality criterion is formulated as follows:

CRITERION *W*. Among all divisions of  $p$  characteristics into  $k$  disjoint and non-empty clusters the optimum one is that for which (4) reaches its minimum.

The idea of this criterion is clear; optimal is such a division into clusters for which the clusters are correlated as small as possible.

Given a fixed number of clusters, the problem of finding an optimal clustering is theoretically solvable for both criteria  $K$  and  $W$ . Up to now, however, algorithms for fast clustering are not known. The number of possible divisions is vast enough and even high-speed computers do not allow to solve the problem by enumeration. Therefore, instead of searching for the optimum solution, "good" solutions must satisfy.

The methods of finding good solutions may be divided in two groups: agglomeration and division methods. Agglomeration methods begin with treating each cluster as consisting of one element. The number of clusters

is then successively diminished by one through merging of clusters after given rules. The methods of Ward [5] and King [2] belong to that type of methods.

Division methods begin with treating all objects as one cluster. The number of clusters is then successively increased by 1 through division of one of the existing clusters after given rules. Such a method is the division of the minimum spanning tree in the Wrocław taxonomy [1].

### 3.2. Optimality criterion for division of characteristics into clusters.

Division of a set of characteristics into clusters of characteristics mutually strongly correlated or into clusters with correlations as weak as possible between them is an intuitive demand. Some formalizations of it, however, may lead to unsatisfactory results. King [2] observes that an application of the  $W$ -criterion to division into  $k = 2$  clusters resulted in optimum divisions into a one-element and a  $(p - 1)$ -element clusters. I have tried several simple division criteria and obtained also paradoxical results. Thus, the criterion of maximum sum of the sums of intra-cluster correlation coefficients and that of maximum sum of mean intra-cluster correlation coefficients (summed together with the ones on the main diagonal) lead to  $k - 1$  clusters of one element each and one cluster of  $p - k + 1$  elements. The criterion of maximum sum of mean intra-cluster correlation coefficients gives division into clusters of nearly equal numbers of elements. These experiments learned me to formulate the criterion of division of characteristics after correlation coefficients in a rather elaborate way.

Assume that each of the characteristics  $X_1, X_2, \dots, X_p$  has been obtained for  $n$  individuals. Denote the normalized values of those characteristics by  $x_{ia}$ , where  $a$  is the individual number. The set of values  $(x_{i1}, x_{i2}, \dots, x_{in})$  may be treated as a point in  $n$ -dimensional Euclidean space. The square of the distance between two characteristics is equal to the square of the distance between two such points. It is equal to

$$d^2(X_i, X_j) = 2(1 - r_{ij}).$$

Elementary algebraic derivations lead to the conclusion that — with such a distance definition — the  $K$ -criterion is equivalent to the following

CRITERION  $K'$ . *Among all divisions of  $p$  characteristics into a fixed number of  $k$  disjoint and non-empty clusters the optimum division is that one for which*

$$K' = \sum_{i=1}^k \frac{1}{p_i} \sum_{s,t=1}^{n_i} r_{i,st}$$

reaches its maximum, where  $r_{i,st}$  is the correlation coefficient between characteristics numbered  $s$  and  $t$  in cluster number  $i$ . The sum  $\sum_{s,t=1}^{p_i} r_{i,st}$  represents the sum of all correlation matrix elements, thus also of the ones on the main diagonal.

Besides of the geometrical interpretation, the  $K'$ -criterion has also a "statistical" interpretation. Assume that division into clusters of characteristics is made for the purpose of reducing the number of characteristics, i.e. to replace the characteristics of one cluster by the weighted mean

$$Y_i = \frac{1}{\sqrt{p_i}} (X_{i1} + X_{i2} + \dots + X_{ip_i}),$$

where  $X_{ij}$ ,  $j = 1, 2, \dots, p_i$ , are the characteristics of cluster  $i$ , normalized to mean zero and dispersion 1. The quantity  $Y_i$  is the distance between the projection of point  $X_i = (X_{i1}, X_{i2}, \dots, X_{ip_i})$  on the main diagonal of the coordinate system and the origin of the system. This quantity is called the *first centroid component* [4]. It is easy to calculate that

$$\mathbf{E}Y_i = 0, \quad \text{Var } Y_i = \mathbf{E}Y_i^2 = \frac{1}{p_i} \sum_{s,t=1}^{p_i} r_{i,st},$$

where the symbol  $\mathbf{E}$  denotes the expected value. It is also easy to verify that

$$\text{Var } Y_i \leq p_i^2/p_i = p_i.$$

The variance of  $Y_i$  is thus the greater the stronger the correlation between characteristics; it assumes its maximum equal to  $p_i$  if all correlation coefficients are equal to 1. For  $K'$  to be maximum is equivalent to requiring the maximum of

$$\sum_{i=1}^k (\text{Var } Y_i - p_i).$$

That condition is easily generalized. Its interesting interpretation in factor analysis will be published elsewhere.

**3.3. Division by using the optimum ordering.** It is right to assume that division of optimally (i.e. so that consecutive objects are similar) ordered objects leads to a clustering which is satisfactory in the sense of having in clusters similar objects. Therefore I want to propose the following principle of clustering, where we assume that the ordering  $X_1 X_2 \dots X_p$  is  $Q$ -optimal:

Definition. A division into  $k$  clusters is called *admissible* if it is performed in the following manner:

$$\begin{aligned} G_1 &= \{X_i: i = 1, 2, \dots, p_1\}, \\ G_2 &= \{X_i: i = p_1+1, p_1+2, \dots, p_1+p_2\}, \\ &\dots\dots\dots \\ G_k &= \{X_i: i = p_1+\dots+p_{k-1}+1, \dots, p\}. \end{aligned}$$

CRITERION  $AK'$ . An admissible division into  $k$  clusters is optimal if the quantity  $K'$  reaches its maximum.

To find the optimal division into  $k$  clusters, it is necessary to verify  $\binom{p-1}{k}$  admissible divisions. I believe that in practice one obtains a good clustering applying the agglomeration principle to admissible divisions. One should therefore begin with  $k = p$  clusters by considering each characteristics as one cluster. The number of clusters is then successively diminished through

- (a) merging of neighbouring characteristics,
- (b) merging of a neighbouring characteristics with a cluster of several characteristics,
- (c) merging of two clusters consisting of neighbouring characteristics.

Among the in such sense admissible clusters it is necessary to choose in every step that one for which  $K'$  is maximum.

**3.4. An example.** The set of motor ability characteristics (see 2.4) has been clustered after the  $AK'$ -criterion. Table 3 presents the optimum divisions for  $k = 13, 12, 11, \dots, 2$ . The characteristics are numbered as in the quasi-optimal ordering. Characteristics belonging to one cluster are parenthesized. However, for one-element clusters no parentheses were used. The last column of Table 3 gives the values of  $K'$ .

The method known as *Wrocław taxonomy* [1] divides into clusters as follows. A minimum spanning tree over all objects is formed. For division into  $k$  clusters,  $k-1$  greatest distances are deleted in the tree. This method has been applied to the mentioned set of characteristics using as distance the values  $1-r_{ij}$  in the minimum spanning tree. Table 4 presents the clusters obtained in such a way. Besides of the division into 13 clusters, the division of the quasi-optimal ordering lead to a greater  $K'$  than the division of the minimum spanning tree. For  $k = 13$  the Wrocław taxonomy merged characteristics nos. 2 and 4 in one cluster which was an inadmissible clustering in the first method.

**References**

[1] K. Florek, J. Łukaszewicz, J. Perkal, H. Steinhaus et S. Zubrzycki, *Sur la liaison et la division des points d'un ensemble fini*, Coll. Math. 2 (1951), p. 282-285.

- [2] B. King, *Step-wise clustering procedures*, J. Amer. Statist. Assn. 62 (1967), p. 86-101.
- [3] E. Marczewski i H. Steinhaus, *O odległości systematycznej biotopów*, Zastosow. Matem. 4 (1959), p. 195-203.
- [4] F. Szczotka, *Wskaźniki przyrodnicze i składowe zespołu cech*, Listy Biometryczne 12-15 (1966), p. 3-54.
- [5] J. H. Ward, *Hierarchical grouping to optimize an objective function*, J. Amer. Statist. Assn. 58 (1963), p. 236-244.
- [6] S. S. Wilks, *On the independence of  $k$  sets of normally distributed statistical variables*, Econometrica 3 (1935), p. 309-326.

Received on 15. 4. 1971

T A B L E 1. Correlation matrix for 14 tests of motor ability

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
1	1													
2	.269	1												
3	.491	.422	1											
4	.842	.347	.563	1										
5	.738	.266	.493	.745	1									
6	.607	.431	.582	.685	.478	1								
7	.603	.483	.547	.678	.510	.708	1							
8	.639	.482	.606	.715	.554	.724	.792	1						
9	.454	.648	.483	.540	.395	.510	.561	.551	1					
10	.269	.611	.257	.357	.250	.365	.420	.403	.656	1				
11	.760	.333	.578	.787	.703	.571	.616	.658	.463	.287	1			
12	.249	.740	.353	.324	.265	.315	.454	.407	.534	.538	.298	1		
13	.238	.588	.323	.304	.311	.280	.343	.349	.401	.375	.304	.561	1	
14	.741	.413	.667	.816	.745	.641	.654	.709	.545	.346	.765	.378	.366	1

T A B L E 2. Correlation matrix for 14 tests of motor ability. The variables are in quasi-optimal order

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
1	1													
2	.738	1												
3	.703	.760	1											
4	.745	.842	.787	1										
5	.745	.741	.765	.816	1									
6	.478	.607	.571	.685	.641	1								
7	.554	.639	.658	.715	.709	.724	1							
8	.510	.603	.616	.678	.654	.708	.792	1						
9	.493	.491	.578	.563	.667	.582	.606	.547	1					
10	.395	.454	.463	.540	.545	.510	.551	.561	.483	1				
11	.266	.269	.304	.347	.413	.431	.482	.483	.422	.648	1			
12	.250	.269	.287	.357	.346	.365	.403	.420	.257	.656	.611	1		
13	.265	.249	.298	.324	.378	.315	.407	.454	.353	.534	.740	.538	1	
14	.311	.238	.333	.304	.366	.280	.349	.343	.323	.401	.588	.375	.561	1

T A B L E 3. Optimal clusters for tests of motor ability

$k$															$K'$
13	1	2	3	[4	5]	6	7	8	9	10	11	12	13	14	13.816
12	1	2	3	[4	5]	6	[7	8]	9	10	11	12	13	14	13.608
11	1	[2	3	4	5]	6	7	8	9	10	11	12	13	14	13.807
10	1	[2	3	4	5]	6	[7	8]	9	10	11	12	13	14	13.599
9	1	[2	3	4	5]	[6	7	8]	9	10	11	12	13	14	13.290
8	1	[2	3	4	5]	[6	7	8]	9	[10	11]	12	13	14	12.938
7	1	[2	3	4	5]	[6	7	8]	9	[10	11	12]	13	14	12.567
6	1	[2	3	4	5]	[6	7	8]	9	[10	11	12]	[13	14]	12.128
5	1	[2	3	4	5]	[6	7	8	9]	[10	11	12]	[13	14]	11.625
4	1	[2	3	4	5]	[6	7	8	9	10]	[11	12	13	14]	11.143
3	[1	2	3	4	5]	[6	7	8	9	10]	[11	12	13	14]	10.393
2	[1	2	3	4	5	6	7	8	9	10]	[11	12	13	14]	9.530

T A B L E 4. Clustering of the motor ability tests by the Wrocław taxonomy

$k$															$K'$
13	1	3	[2	4]	5	6	7	8	9	10	11	12	13	14	13.842
12	1	3	[2	4	5]	6	7	8	9	10	11	12	13	14	13.599
11	1	3	[2	4	5]	6	[7	8]	9	10	11	12	13	14	13.391
10	1	[2	3	4	5]	6	[7	8]	9	10	11	12	13	14	13.148
9	[1	2	3	4	5]	6	[7	8]	9	10	11	12	13	14	12.849
8	[1	2	3	4	5]	6	[7	8]	9	10	[11	13]	12	14	12.589
7	[1	2	3	4	5	9]	6	[7	8]	10	[11	13]	12	14	12.010
6	[1	2	3	4	5	9]	[6	7	8]	10	[11	13]	12	14	11.701
5	[1	2	3	4	5	6	7	8	9]	10	[11	13]	12	14	11.009
4	[1	2	3	4	5	6	7	8	9]	[10	12]	[11	13]	14	10.665
3	[1	2	3	4	5	6	7	8	9]	[10	11	12	13]	14	10.133
2	[1	2	3	4	5	6	7	8	9]	[10	11	12	13	14]	9.530

F. A. SZCZOTKA (Warszawa)

## O PEWNEJ METODZIE PORZĄDKOWANIA I GRUPOWANIA OBIEKTÓW

## S T R E S Z C Z E N I E

Niech dla obiektów  $X_1, X_2, \dots, X_p$  będą określone odległości  $d(X_i, X_j) = d_{ij}$ . Autor proponuje — jako kryterium dobroci uporządkowania  $X_{a_1} X_{a_2} \dots X_{a_p}$  — funkcję (1), przy czym większe wartości  $Q^*$  wskazują na lepsze uporządkowanie. Jeżeli obiektami są cechy  $X_1, \dots, X_p$  o znanych współczynnikach korelacji  $r_{ij} \geq 0$ , to proponowanym kryterium jest funkcja (2), przy czym na lepsze uporządkowanie wskazują mniejsze wartości  $Q$ . Podaje się algorytm dla znalezienia uporządkowania quasi-optimalnego, tzn. takiego, którego nie można poprawić przez transpozycję dwóch obiektów.

W drugiej części autor proponuje i dyskutuje następujące kryterium optymalności podziału cech  $X_1, \dots, X_p$  na  $k$  grup. Za optymalny uznaje się ten podział, dla którego

$$K' = \sum_{i=1}^k \frac{1}{p_i} \sum_{s,t=1}^{p_i} r_{i.st}$$

jest maksymalne, gdzie  $p_i$  jest ilością cech w grupie o numerze  $i$ , a  $\sum_{s,t=1}^{p_i} r_{i.st}$  jest sumą wszystkich wyrazów macierzy korelacji tych cech. Sugeruje się, jako praktycznie dobry sposób podziału na grupy, pocięcie optymalnego (lub quasi-optymalnego) uporządkowania na  $k$  odcinków w ten sposób, by  $K'$  osiągało maksimum.

---