

## ***NP*-Hard Problems in Hierarchical-Tree Clustering**

Mirko Křivánek<sup>1</sup> and Jaroslav Morávek<sup>2</sup>

<sup>1</sup> Research Institute of Mathematical Machines, Loretánské nám. 3, 11855 Praha 1, ČSSR

<sup>2</sup> Mathematical Institute, Czechoslovak Academy of Sciences, Žitná 25, 11567 Praha 1, ČSSR

**Summary.** We consider a class of optimization problems of hierarchical-tree clustering and prove that these problems are *NP*-hard. The sequence of polynomial reductions and/or transformations used in our proof is based on relatively laborious graph-theoretical constructions and starts in the *NP*-complete problem of 3-dimensional matching. Using our main result we establish the *NP*-completeness of a problem of the best approximation of a symmetric relation on a finite set by an equivalence relation, thus answering in the negative a question proposed implicitly by C.T. Zahn.

### **I. Introduction and Statement of the Main Result**

Within the last twenty years an enormous number of strategies for cluster analysis have been proposed [3, 6, 8, 15]. Though the main emphasis in these efforts has been concentrated on the creation or application of clustering techniques, relatively little attention has been paid to the study of computational complexity of clustering algorithms. Because of wide range of applications of cluster analysis [1, 13] there are many variations of problem formulation. Generally we can consider two types of goals:

(i) *Nonhierarchical clustering* where the goal is to partition a given finite set of objects into nonempty clusters (the number of clusters can be specified beforehand), those objects in the same cluster being considered as close or similar and those in different clusters as distant or dissimilar. The “quality” of this clustering is usually expressed by a real objective function defined on the family of all partitions of the set of objects.

(ii) *Hierarchical clustering (hierarchical-tree clustering)* where the goal is to construct a sequence of nested nonhierarchical clusterings which form a so called *hierarchical tree* and which have to be optimal with respect to a criterion [7, 9].

The problems of nonhierarchical clustering were studied partially from the point of view of computational complexity, see e.g. [2, 5]. On the other hand,

to our knowledge, no comparable results concerning the complexity of hierarchical clustering have been published. Our main aim is to present a result in this respect. For the terminology concerning the computational complexity (*NP*-theory) used in this paper see [4].

Let us review and formalize the main concept of hierarchical clustering in which lies our main interest, cf. e.g. [10]. Throughout this paper,  $n$  will denote an integer,  $n \geq 3$ ,  $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$  an  $n$ -element set, and  $D = (d_{i,j})$  a symmetric real  $n \times n$ -matrix such that

$$d_{i,j} > 0 \text{ if } i \neq j, \text{ and } d_{i,j} = 0 \text{ if } i = j \quad (i, j \in \{1, 2, \dots, n\}).$$

Elements of  $\Omega$  are called objects (these are to be clustered), and  $D$  is called a dissimilarity matrix. Within our context we interpret a “small” value of  $d_{i,j}$  as a close relationship between objects  $\omega_i$  and  $\omega_j$ .

A hierarchical tree  $T$  over  $\Omega$  is defined as a finite sequence of pairs  $T = ((P_1, l_1), (P_2, l_2), \dots, (P_q, l_q))$  where

- (i)  $P_1, P_2, \dots, P_q$  are partitions<sup>1</sup> of  $\Omega$ ;
- (ii)  $l_1, l_2, \dots, l_q$  are integers,

$$0 = l_1 < l_2 < \dots < l_q;$$

- (iii)  $P_k$  is proper refinement of  $P_{k+1}$  ( $1 \leq k \leq q-1$ );
- (iv)  $P_1 = \{\{\omega_1\}, \{\omega_2\}, \dots, \{\omega_n\}\}$  and  $P_q = \{\Omega\}$ .

The integer  $q$  is called the height of  $T$  and number  $l_k$ , the  $k$ -th level of partition  $P_k$  in  $T$ . It follows that  $2 \leq q \leq n$ .

Let  $\mathfrak{A}(\Omega)$  be the set of all hierarchical trees  $T$  over  $\Omega$  and  $\mathfrak{A}_q(\Omega)$ , where  $2 \leq q \leq n$ , the set of all hierarchical trees over  $\Omega$ , having the height  $q$ ; if  $\Omega$  is evident from the context we shall write  $\mathfrak{A}$  and  $\mathfrak{A}_q$  instead of  $\mathfrak{A}(\Omega)$  and  $\mathfrak{A}_q(\Omega)$ , respectively.

Further we define the function  $u\langle T \rangle: \Omega \times \Omega \rightarrow \mathbb{N}_0$  ( $\mathbb{N}_0$  is the set of all nonnegative integers), corresponding to a given hierarchical tree

$$T = ((P_1, l_1), (P_2, l_2), \dots, (P_q, l_q)) \in \mathfrak{A}(\Omega),$$

as follows:

$$u\langle T \rangle(\omega_i, \omega_j) \stackrel{\text{def}}{=} \min \{l_k \mid \text{there exists } M \in P_k \\ (1 \leq k \leq q) \text{ such that } \{\omega_i, \omega_j\} \subseteq M\}.$$

Remarks: 1) Function  $u\langle T \rangle$  is an ultrametric on  $\Omega$  (cf. [10]). From the point of view of graph theory, a hierarchical tree can be interpreted as a rooted tree. As the rooted tree is an upper semilattice we can interpret  $u\langle T \rangle(\omega_i, \omega_j)$  as the level  $l_k$  assigned to the least upper bound of elements  $\omega_i$  and  $\omega_j$ . In Fig. 1 we give the graphical representation of the hierarchical tree

$$T' = ((P'_1, l'_1), (P'_2, l'_2), (P'_3, l'_3), (P'_4, l'_4)),$$

<sup>1</sup> i.e. finite disjoint decompositions into nonempty classes

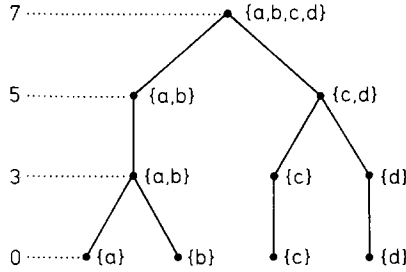


Fig. 1

where

- $(P'_1, l'_1) = (\{\{a\}, \{b\}, \{c\}, \{d\}\}, 0)$
- $(P'_2, l'_2) = (\{\{a, b\}, \{c\}, \{d\}\}, 3)$
- $(P'_3, l'_3) = (\{\{a, b\}, \{c, d\}\}, 5)$
- $(P'_4, l'_4) = (\{\{a, b, c, d\}\}, 7).$

The ultrametric  $u\langle T' \rangle$  corresponding to  $T'$  is given in the underlying tableau:

$u\langle T' \rangle$	a	b	c	d	
a	0	3	7	7	
b	3	0	7	7	
c	7	7	0	5	
d	7	7	5	0	//

For the evaluation of hierarchical clustering we use the objective function  $F: \mathfrak{A} \rightarrow \mathbb{R}_+$  (the set of all nonnegative reals), defined as follows (cf. [9]):

$$F(T) \stackrel{\text{df}}{=} \sum_{1 \leq i < j \leq n} |d_{i,j} - u\langle T \rangle(\omega_i, \omega_j)|.$$

Using  $F$  we introduce the following optimization problems: Problem HIC (hierarchical clustering):

**INSTANCE:** Set of objects  $\Omega$  and dissimilarity matrix  $D$  of the size  $n \times n$ , where  $n = \text{card}(\Omega)$ ;

**PROBLEM:** Determine a hierarchical tree  $T_* \in \mathfrak{A}$  such that

$$F(T_*) = \min \{F(T) \mid T \in \mathfrak{A}\}. \quad //$$

Problem  $\text{HIC}_q (q=2, 3, \dots)$ :

**INSTANCE:** Set of objects  $\Omega$  and dissimilarity matrix  $D$  of the size  $n \times n$ , where  $n = \text{card}(\Omega)$ ;

**PROBLEM:** Determine a hierarchical tree  $T_+ \in \mathfrak{A}_q$  such that

$$F(T_+) = \min \{F(T) \mid T \in \mathfrak{A}_q\}. \quad //$$

It is evident that problem  $\text{HIC}_2$  has a polynomial time complexity. Our main result can be summarized as follows:

**Theorem.** *Problems  $\text{HIC}$  and  $\text{HIC}_q$  for  $q \geq 3$  are  $NP$ -hard. This result was presented at  $\text{COMPSTAT}$  1984 (Prague), cf. [12].  $\square$*

## II. Proof of the Main Result

We shall obtain the proof of our theorem by proving several lemmas.

**Lemma 1.** *For an arbitrary integer  $q \geq 2$  we have  $\text{HIC}_q \propto \text{HIC}_{q+1}$ .*

*Proof.* We apply the well-known method of local replacement, see e.g. [4]. To each instance  $(\Omega, D)$  of  $\text{HIC}_q$  let us assign an instance  $(\Omega', D')$  of  $\text{HIC}_{q+1}$  as follows:

$$\Omega' \stackrel{\text{df}}{=} \Omega \cup \{\omega_{n+1}\} = \{\omega_1, \omega_2, \dots, \omega_n, \omega_{n+1}\},$$

where  $\omega_{n+1}$  is a 'new' object (joined to  $\Omega$ ),

$$D' = (d'_{i,j}), \quad (1 \leq i, j \leq n+1),$$

where

$$d'_{i,j} \stackrel{\text{df}}{=} d_{i,j} \quad \text{for } 1 \leq i, j \leq n,$$

$$d'_{i,n+1} \stackrel{\text{df}}{=} d'_{n+1,j} \stackrel{\text{df}}{=} n^2(q + \max\{d'_{i',j'} \mid 1 \leq i', j' \leq n\}) + 1, \quad (1 \leq i, j \leq n),$$

and

$$d'_{n+1,n+1} \stackrel{\text{df}}{=} 0.$$

(Let us observe that the  $d'_{i,j}$  are computable from  $d_{i,j}$  in polynomial time.)

Now, the following equivalence is easily verified:

$$T' = ((P'_1, l'_1), (P'_2, l'_2), \dots, (P'_q, l'_q), (P'_{q+1}, l'_{q+1}))$$

is a solution of  $\text{HIC}_{q+1}$  if and only if

$$T \stackrel{\text{df}}{=} ((P'_1 \setminus \{\omega_{n+1}\}, 0), (P'_2 \setminus \{\omega_{n+1}\}, l'_2), \dots, (P'_q \setminus \{\omega_{n+1}\}, l'_q))$$

belongs to  $\mathfrak{A}_q(\Omega)$  and it is a solution of  $\text{HIC}_q$ .  $\square$

By virtue of Lemma 1 it is sufficient to prove the  $NP$ -hardness of  $\text{HIC}_3$ . In fact, we shall obtain a slightly stronger result. A dissimilarity matrix  $D = (d_{i,j})$  will be called *binary* if each off-diagonal element of  $D$  equals either 1 or 2, ( $d_{i,j} \in \{1, 2\}$  if  $i \neq j$ ). Let  ${}^b\text{HIC}$  and  ${}^b\text{HIC}_q$  denote the 'binary restrictions' of  $\text{HIC}$  and  $\text{HIC}_q$ , respectively, i.e. the computational problems defined in the precisely same way as  $\text{HIC}$  and  $\text{HIC}_q$  respectively, except that the instance  $D$  is a binary matrix. It follows immediately that

$${}^b\text{HIC} \propto \text{HIC} \quad \text{and} \quad {}^b\text{HIC} \propto \text{HIC}_q \quad \text{for } q \geq 2. \quad (1)$$

**Lemma 2.** *It holds that  ${}^b\text{HIC}_3 \propto {}^b\text{HIC}$ .*

*Proof.* It is sufficient to prove the following assertion: If  $T$  is a solution of  ${}^b\text{HIC}$  then  $T \in \mathfrak{A}_2 \cup \mathfrak{A}_3$ . Let us assume on the contrary, that there exists a solution

$$T = ((P_1, l_1), (P_2, l_2), \dots, (P_r, l_r))$$

of  ${}^b\text{HIC}$  with the property  $r > 3$ . Then it is easy to show that for the hierarchical tree

$$T^* = ((P_1, 0), (P_2, l_2), (P_r, l_2 + 1)) \in \mathfrak{A}_3$$

we have

$$u\langle T^* \rangle(\omega_i, \omega_j) = u\langle T \rangle(\omega_i, \omega_j) \quad \text{if } u\langle T \rangle(\omega_i, \omega_j) \leq 1 \quad (2)$$

and

$$\leq u\langle T \rangle(\omega_i, \omega_j) \quad \text{otherwise,} \quad (3)$$

$$\max(2, u\langle T^* \rangle(\omega_i, \omega_j)) < u\langle T \rangle(\omega_i, \omega_j) \quad (4)$$

for some pair  $(\omega_i, \omega_j) \in \Omega \times \Omega$ .

Combining (2), (3) and (4) we obtain  $F(T^*) < F(T)$ , which completes the proof.  $\square$

In the next lemma we investigate the levels of a hierarchical tree solving  ${}^b\text{HIC}_3$ .

**Lemma 3.** *If  $T = ((P_1, 0), (P_2, l_2), (P_3, l_3)) \in \mathfrak{A}_3$  is a solution of  ${}^b\text{HIC}_3$  then  $l_2 = 1$  and  $l_3 = 2$ .*

*Proof.* It is sufficient to prove that for each

$$T = ((P_1, 0), (P_2, l_2), (P_3, l_3)) \in \mathfrak{A}_3$$

with the property  $l_3 \geq 3$  there exists

$$T' = ((P'_1, 0), (P'_2, l'_2), (P'_3, l'_3)) \in \mathfrak{A}_3$$

such that

$$l'_3 < l_3 \quad \text{and} \quad F(T') < F(T). \quad (5)$$

Indeed, let us consider following two cases for  $(0, l_2, l_3)$ :

$$(\alpha) \quad l_3 > 3 \quad \text{or} \quad (l_2, l_3) = (1, 3); \quad (\beta) \quad (l_2, l_3) = (2, 3).$$

In the  $(\alpha)$  case we can put evidently  $P'_j = P_j$  ( $j = 1, 2, 3$ ),  $l'_2 = \min(2, l_2)$  and  $l'_3 = l_3 - 1$ .

In the  $(\beta)$  case we put  $l'_2 = 1$ ,  $l'_3 = 2$  and

$$P'_2 = \{\{\omega_1, \omega_2\}, \{\omega_3\}, \{\omega_4\}, \dots, \{\omega_f\}\}.$$

Since  $l'_3 < l_3$  we have to verify inequality  $F(T') < F(T)$ . Indeed,

$$F(T') = (|d_{1,2} - 1| - |d_{1,2} - 2|) + \sum_{1 \leq i < j \leq n} |d_{i,j} - 2|$$

i.e.

$$F(T') \leq 1 + \text{card}\{(i, j) \mid 1 \leq i < j \leq n, d_{i,j} = 1\}. \quad (6)$$

On the other hand, let  $P_2 = \{I_1, I_2, \dots, I_r\}$ , hence  $2 \leq r \leq n - 1$ . Then we have

$$\begin{aligned} F(T) &= \sum_{\rho=1}^r \sum_{\{i,j\} \subset I_\rho} |d_{i,j} - 2| + \sum_{1 \leq \rho' < \rho'' \leq r} \sum_{i \in I_{\rho'}} \sum_{j \in I_{\rho''}} |d_{i,j} - 3| \\ &= \sum_{\rho=1}^r \sum_{\{i,j\} \subset I_\rho} |d_{i,j} - 2| + \sum_{1 \leq \rho' < \rho'' \leq r} \sum_{i \in I_{\rho'}} \sum_{j \in I_{\rho''}} |d_{i,j} - 2| \\ &\quad + \sum_{1 \leq \rho' < \rho'' \leq r} \sum_{i \in I_{\rho'}} \sum_{j \in I_{\rho''}} (|d_{i,j} - 3| - |d_{i,j} - 2|) \\ &= \sum_{1 \leq i < j \leq n} |d_{i,j} - 2| + \sum_{1 \leq \rho' < \rho'' \leq r} \text{card}(I_{\rho'}) \text{card}(I_{\rho''}), \end{aligned}$$

i.e.

$$F(T) \geq \text{card} \{ \{i,j\} \mid 1 \leq i < j \leq n, d_{i,j} = 1 \} + (n - 1). \tag{7}$$

Combining (6) and (7) we obtain (5) (since  $n \geq 3$ ). This completes the proof.  $\square$

In the next lemma, problem  ${}^b\text{HIC}_3$  is restated using the following terminology. For an arbitrary partition  $\{I_1, I_2, \dots, I_r\}$  of  $\Omega$  let us set  $i_\rho = \text{card}(I_\rho)$  and

$$j_\rho = \text{card} \{ \{i,j\} \subseteq I_\rho \mid d_{i,j} = 1 \} \quad (\rho = 1, 2, \dots, r).$$

**Lemma 4.** *Problem  ${}^b\text{HIC}_3$  can be stated equivalently as follows: Find a partition  $\{I_1, I_2, \dots, I_r\}$  of  $\Omega$  such that*

$$\sum_{\rho=1}^r \left( \binom{i_\rho}{2} - 2j_\rho \right)$$

*is minimum.*

*Proof.* For every hierarchical tree

$$T = ((P_1, 0), (\{I_1, I_2, \dots, I_r\}, 1), (\{\Omega\}, 2)) \in \mathfrak{A}_3$$

we have

$$\begin{aligned} F(T) &= \sum_{\rho=1}^r \sum_{\{i,j\} \subset I_\rho} |d_{i,j} - 1| + \sum_{1 \leq \rho' < \rho'' \leq r} \sum_{i \in I_{\rho'}} \sum_{j \in I_{\rho''}} |d_{i,j} - 2| \\ &= \sum_{\rho=1}^r \left( \binom{i_\rho}{2} - j_\rho \right) + \text{card} \{ \{i,j\} \mid d_{i,j} = 1 \} - \sum_{\rho=1}^r j_\rho \\ &= \text{card} \{ \{i,j\} \mid d_{i,j} = 1 \} + \sum_{\rho=1}^r \left( \binom{i_\rho}{2} - 2j_\rho \right), \end{aligned}$$

which concludes the proof.  $\square$

For proving the NP-hardness of  ${}^b\text{HIC}_3$  we shall use the following decision problem: Problem EC3 (exact cover by ordered 3-tuples):

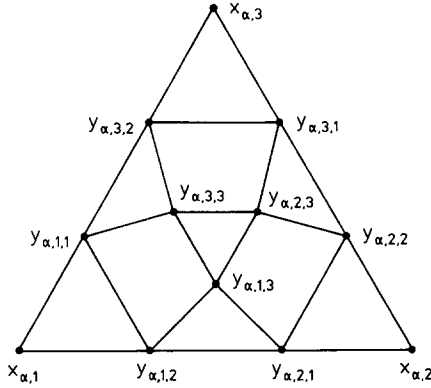


Fig. 2

- INSTANCE:** 1) finite set  $X$  with  $\text{card}(X)=3m$  for some positive integer  $m$ ;  
 2) finite indexed family

$$\mathcal{C} = ((x_{\alpha,1}, x_{\alpha,2}, x_{\alpha,3}))_{\alpha=1}^a$$

of ordered 3-tuples of elements of  $X$  with the property that each element of  $X$  occurs at least in one 3-tuple of  $\mathcal{C}$ , i.e.

$$\bigcup_{\alpha=1}^a \{x_{\alpha,1}, x_{\alpha,2}, x_{\alpha,3}\} = X;$$

**QUERY:** Decide whether  $\mathcal{C}$  contains a subfamily  $\mathcal{C}'$  such that each element of  $X$  occurs in exactly one 3-tuple of  $\mathcal{C}'$ . //

$\mathcal{C}'$  is called an exact cover for  $X$ .

**Lemma 5.** *Problem EC 3 is NP-complete, (cf. [11]).*  $\square$

Our next aim is to reduce polynomially EC 3 to  ${}^b\text{HIC}_3$ . For this aim we assign to each instance  $(X, \mathcal{C})$  of EC 3 an instance  $(\Omega, D)$  of  ${}^b\text{HIC}_3$ . First, let us put

$$n \stackrel{\text{df}}{=} 3m + 9 \text{card}(\mathcal{C}) = 3m + 9a \tag{8}$$

and

$$\Omega \stackrel{\text{df}}{=} X \cup \{y_{\alpha,\beta,\gamma} | \alpha \in \{1, 2, \dots, a\}; \beta, \gamma \in \{1, 2, 3\}\}, \tag{9}$$

where  $y_{\alpha,\beta,\gamma}$  are 9.a ‘new’ objects joined to  $X$ . (It is appropriate to index these objects by the triple subscripts.)

The dissimilarity matrix  $D$  will be introduced using certain graphs (we consider in this paper finite undirected graphs without loops and parallel edges). For each  $\alpha \in \{1, 2, \dots, a\}$  let us consider the graph  $G_\alpha = (V_\alpha, E_\alpha)$ , see Fig. 2, where

- (i) The vertex-set is:

$$V_\alpha \stackrel{\text{df}}{=} \{x_{\alpha,1}, x_{\alpha,2}, x_{\alpha,3}\} \cup \{y_{\alpha,\beta,\gamma} | 1 \leq \beta, \gamma \leq 3\};$$

(ii) The edge-set is (edges are defined as 2-element subsets of the vertex-set):

$$E_\alpha \stackrel{\text{df}}{=} \{ \{x_{\alpha,\beta}, y_{\alpha,\beta,\gamma}\} \mid 1 \leq \beta \leq 3, \gamma = 1, 2 \} \\ \cup \{ \{y_{\alpha,\beta,3}, y_{\alpha,\beta',3}\} \mid 1 \leq \beta \neq \beta' \leq 3 \} \\ \cup \{ \{y_{\alpha,\beta,2}, y_{\alpha,\beta,\gamma}\} \mid 1 \leq \beta \leq 3, \gamma = 1, 3 \} \\ \cup \{ \{y_{\alpha,\beta,1}, y_{\alpha,\beta',\gamma}\} \mid 1 \leq \beta \leq 3, \gamma = 2, 3 \text{ and } (\beta - \beta' = 1) \vee (\beta' - \beta) = 2 \}.$$

Using graphs  $G_\alpha$  we introduce the graph  $\mathbf{G} \stackrel{\text{df}}{=} (\Omega, \mathbf{E})$ , where

$$\mathbf{E} = \bigcup_{\alpha=1}^a E_\alpha.$$

The dissimilarity matrix  $D = (d_{i,j})$  will be now defined as follows: We consider an arbitrary fixed numbering  $\omega_1, \omega_2, \dots, \omega_n$  of elements of  $\Omega$  and put

$$d_{i,j} = 0 \quad \text{if } i = j, \tag{10}$$

$$d_{i,j} = 1 \quad \text{if } i \neq j \text{ and } \{\omega_i, \omega_j\} \in \mathbf{E}, \tag{11}$$

$$d_{i,j} = 2 \quad \text{otherwise.} \tag{12}$$

In the sequel, we shall use some additional graph-theoretical definitions and notations: The term *subgraph* will denote an induced subgraph; the subgraph of  $\mathbf{G}$  induced by a nonempty subset  $I \subseteq \Omega$  will be denoted by  $G(I)$ . A subgraph of  $\mathbf{G}$  with 3 vertices and 3 edges will be called a *triangle*.

Let  $E = \{(W_i, H_i)\}$  be a finite set of triangles;  $W_i$  is the vertex-set and  $H_i$  is the edge-set of the triangle  $(W_i, H_i)$ .  $E$  will be called a vertex-partition of the graph  $\mathbf{G}$  into triangles if  $\bigcup_i W_i = \Omega$  and  $\bigcup_i H_i \subseteq \mathbf{E}$ .

**Lemma 6.** *A solution of EC3 exists if and only if there exists a vertex-partition of  $\mathbf{G}$  into triangles.*

*Proof.* Let  $\mathcal{C}$  be an exact cover for  $X$  with respect to EC3. Let  $E$  be defined as the minimum, with respect to the cardinality, set of triangles such that the following two conditions are fulfilled ( $\alpha \in \{1, 2, \dots, \text{card}(\mathcal{C})\}$ ):

(i) If  $(x_{\alpha,1}, x_{\alpha,2}, x_{\alpha,3}) \in \mathcal{C}$  then

$$G(\{y_{\alpha,1,3}, y_{\alpha,2,3}, y_{\alpha,3,3}\}) \in E$$

and

$$G(\{x_{\alpha,\beta}, y_{\alpha,\beta,1}, y_{\alpha,\beta,2}\}) \in E$$

for all  $\beta = 1, 2, 3$ .

(ii) If  $(x_{\alpha,1}, x_{\alpha,2}, x_{\alpha,3}) \in \mathcal{C} - \mathcal{C}'$  then

$$G(\{y_{\alpha,\beta',1}, y_{\alpha,\beta,2}, y_{\alpha,\beta,3}\} \mid (\beta' - \beta) = 1 \vee (\beta - \beta') = 2, 1 \leq \beta, \beta' \leq 3).$$

The set  $E$  contains  $4m + 3(a - m) = m + 3a$  triangles, and it is easy to see that  $E$  is a vertex-partition of  $\mathbf{G}$ .

Conversely, given a vertex-partition  $E$  of  $\mathbf{G}$  into triangles we define  $\mathcal{C}'$  as the family of all  $(x_{\alpha,1}, x_{\alpha,2}, x_{\alpha,3}) \in \mathcal{C}$  such that

$$G(\{y_{\alpha,1,3}, y_{\alpha,2,3}, y_{\alpha,3,3}\}) \in E.$$



It is left to the reader to show that  $\mathcal{C}$  is an exact cover in the sense of problem EC 3 with respect to the instance  $(X, \mathcal{C})$ . (Observe that if  $E$  contains a “central” triangle

$$G(\{y_{\alpha,1,3}, y_{\alpha,2,3}, y_{\alpha,3,3}\})$$

then it contains also 3 triangles

$$G(\{x_{\alpha,\beta}, y_{\alpha,\beta,1}, y_{\alpha,\beta,2}\}), \text{ where } \beta \in \{1, 2, 3\}.$$

This completes the proof.  $\square$

For  $i \in \{1, 2, \dots, \text{card}(\Omega)\}$  let  $m_j(i)$  denote the maximum number of edges in a subgraph  $G(I)$  with  $\text{card}(I) = i$ .

**Lemma 7.** *It holds that*

$$m_j(i) \leq 2i - 3, \quad i \geq 1.$$

*Proof.* Let us consider an arbitrary  $I \subset \Omega$  with  $\text{card}(I) = i$ . Let  $i_\alpha$  denote the number of all edges in the subgraph  $G(V_\alpha \cap I)$ , where  $\alpha \in \{1, 2, \dots, a\}$  and let  $A$  denote the set of all  $\alpha$  with the property  $(V_\alpha \setminus X) \cap I \neq \emptyset$ . Since each graph  $G_\alpha$  is a block of graph  $\mathbf{G}$  (cf. e.g. [14], p. 529) we have

$$m_j(i) = \max \left\{ \sum_{\alpha \in A} i_\alpha \mid I \subset \Omega, \text{card}(I) = i \right\}. \tag{13}$$

Examining the graph  $G_\alpha$  (cf. Fig. 2) one can easily verify the following values of  $m_j(i)$  in  $G_\alpha$ :

$$\begin{aligned} m_j(1) &= 0, & m_j(2) &= 1, & m_j(3) &= 3, & m_j(4) &= 4, & m_j(5) &= 6, \\ m_j(6) &= 8, & m_j(7) &= 10, & m_j(8) &= 12, & m_j(9) &= 15, & m_j(10) &= 17, \\ m_j(11) &= 19, & m_j(12) &= 21. \end{aligned}$$

Hence

$$i_\alpha \leq m_j(\text{card}(V_\alpha \cap I)) \leq 2 \text{card}(V_\alpha \cap I) - 3 / \text{card}(A), \quad \alpha \in A. \tag{14}$$

Altogether (13) and (14) yield

$$m_j(i) \leq 2 \cdot \text{card}(I) - 3 = 2i - 3, \quad i \geq 1.$$

The proof is concluded.  $\square$

**Lemma 8.** *Let  $\emptyset \neq I \subseteq \Omega$ ,  $i = \text{card}(I)$ , and let  $j$  be the number of all edges in  $G(I)$ . Then*

$$\binom{i}{2} - 2j + i \geq 0. \tag{15}$$

Moreover

$$\binom{i}{2} - 2j + i = 0$$

if and only if  $i = 3$  and  $G(I)$  is a triangle.

*Proof.* Obviously it holds that

$$2i - 3 \leq \left( \binom{i}{2} + i \right) / 2, \quad i \geq 1.$$

Thus by applying Lemma 7 the inequality (15) follows. Moreover if  $G(I)$  is a triangle then

$$\binom{i}{2} - 2j + i = 3 - 6 + 3 = 0.$$

Conversely, if

$$\binom{i}{2} - 2j + i = 0$$

then it follows from the definition of  $\mathbf{G}$  and from Lemma 7 that  $i = j = 3$  and  $G(I)$  is a triangle.  $\square$

Now, let the instance  $(\Omega, D)$  for  ${}^b\text{HIC}_3$  be defined by (8)-(12). For each partition  $\{I_1, I_2, \dots, I_r\}$  of  $\Omega$  let us set

$$\Psi(\{I_1, I_2, \dots, I_r\}) \stackrel{\text{df}}{=} \sum_{\rho=1}^r \left( \binom{i_\rho}{2} - 2j_\rho \right),$$

where  $i_\rho, j_\rho$  ( $\rho = 1, 2, \dots, r$ ) are defined as in Lemma 4. It is easily observed that  $j_\rho$  equals the number of all edges in the graph  $G(I_\rho)$ .

**Lemma 9.** *Let  $\{I_1, I_2, \dots, I_r\}$  be a partition of  $\Omega$ . Then*

$$\Psi(\{I_1, I_2, \dots, I_r\}) \geq -3(m + 3 \text{card}(\mathcal{C})). \tag{16}$$

Moreover,

$$\Psi(\{I_1, I_2, \dots, I_r\}) = -3(m + 3 \text{card}(\mathcal{C})) \tag{17}$$

if and only if  $G(I_\rho)$  is a triangle for each  $\rho \in \{1, 2, \dots, r\}$ , i.e.

$$\{G(I_\rho) \mid \rho = 1, 2, \dots, r\}$$

is a vertex-partition of  $G$  into triangles. (Recall that  $m = 1/3 \text{card}(X)$ .)

*Proof.* Let  $\eta_s$  denote the number of all  $\rho \in \{1, 2, \dots, r\}$  with  $\text{card}(I_\rho) = s$ . Then

$$\Psi(\{I_1, I_2, \dots, I_r\}) = \sum_{s=1}^{\infty} \left( \binom{s}{2} - 2 \cdot {}^m j(s) \right) \cdot \eta_s$$

On the other hand

$$\sum_{s=1}^{\infty} s \cdot \eta_s = \text{card}(\Omega) = 3(m + 3 \text{card}(\mathcal{C})).$$

Thus we have

$$\begin{aligned} \Psi(\{I_1, \dots, I_r\}) &= \Psi(\{I_1, \dots, I_r\}) + \left( \sum_{s=1}^{\infty} s \cdot \eta_s - 3(m + a) \right) \\ &= \sum_{s=1}^{\infty} \left( \binom{s}{2} - 2 \cdot {}^m j(s) + s \right) \cdot \eta_s - 3(m + 3 \text{card}(\mathcal{C})). \end{aligned}$$

The inequality (16) now follows immediately from Lemma 7. Moreover, if  $\{G(I_\rho) \mid \rho=1, 2, \dots, r\}$  is a vertex-partition of  $G$  into triangles then (17) clearly holds.

Conversely, let us assume that (17) holds. Then, using Lemma 7 we see successively that

- (j)  $\eta_s=0$  for  $s \neq 3$ ,
  - (jj)  $\eta_3 = m + 3 \text{card}(\mathcal{C})$ ,
  - (jjj) for each  $\rho \in \{1, 2, \dots, r\}$  the subgraph  $G(I_\rho)$  is a triangle.
- The proof is completed.  $\square$

**Lemma 10.** *It holds  $\text{EC3} \propto {}^b\text{HIC}_3$ .*

*Proof.* It is sufficient to see that given a solution of  ${}^b\text{HIC}_3$  for the instance  $(\Omega, D)$ , defined by (8)–(12), we obtain the answer to EC3 using a polynomially bounded algorithm. Indeed, solving  ${}^b\text{HIC}_3$  we obtain a partition  $\{I_1, I_2, \dots, I_r\}$  of  $\Omega$  such that  $\Psi(\{I_1, I_2, \dots, I_r\})$  is minimum. Now, by virtue of Lemma 9 and Lemma 6 we have

$$\Psi(\{I_1, I_2, \dots, I_r\}) = -3(m + 3 \text{card}(\mathcal{C}))$$

if and only if there exists a vertex-partition of  $G$  into triangles i.e. if and only if there exists an exact cover by ordered 3-tuples with respect to EC3 for the instance  $(X, \mathcal{C})$ . This completes the proof.  $\square$

The proof of the announced result (Theorem): The NP-hardness of  $\text{HIC}_q$  for  $q \geq 3$  follows from (1), Lemma 1, Lemma 5 and Lemma 10. The NP-hardness of HIC follows from (1), Lemma 2, Lemma 5 and Lemma 10.

### III. Best Approximation of Symmetric Relation by an Equivalence

In [16] the following optimization problem is proposed: Given a finite non-empty set  $Z$  and a symmetric relation  $s \subseteq Z \times Z$  we are asked to determine an equivalence relation  $e \subseteq Z \times Z$  minimizing the objective function

$$e \mapsto \text{card}(s \Delta e),$$

where  $s \Delta e = (s \setminus e) \cup (e \setminus s)$  is the symmetric difference of the relations  $s, e$  (considered as subsets of  $Z \times Z$ ). I.C. Lerman observed in [13] the importance of this problem in the hierarchical clustering.

By an immediate application of Lemma 10 we prove that the underlying decision computational problem is NP-complete: Problem  $\text{S}\Delta\text{E}$  (best approximation of a symmetric relation by an equivalence relation):

**INSTANCE:** A finite set  $Z = \{z_1, z_2, \dots, z_m\}$ , symmetric relation  $s \subseteq Z \times Z$  and a positive integer  $k$ ;

**QUERY:** Decide whether there exists an equivalence relation  $e \subseteq Z \times Z$  such that

$$\text{card}(s \Delta e) \leq k.$$

To prove this assertion we observe that problem  $S \Delta E$  is evidently in  $NP$ , and exhibit the polynomial transformation

$${}^b\text{HIC}_3 \propto S \Delta E, \quad (18)$$

where  ${}^b\text{HIC}_3$  is the following decision version of  ${}^b\text{HIC}_3$ :

INSTANCE:  $(\Omega, D, k')$ , where  $D$  is binary and  $k'$  is a positive integer;

QUERY: Decide whether there exists  $T \in \mathfrak{A}_3(\Omega)$  such that

$$F(T) = \sum_{1 \leq i < j \leq n} |d_{i,j} - u \langle T \rangle(\omega_i, \omega_j)| \leq k'. \quad //$$

(It follows immediately from Lemma 10 that  ${}^b\text{HIC}_3$  is  $NP$ -hard.)

To prove (18) we assign to an instance  $(\Omega, D, k')$  of  ${}^b\text{HIC}_3$  the instance  $(Z, s, k)$  of  $S \Delta E$ , where

$$Z \stackrel{\text{df}}{=} \Omega \quad \text{and} \quad z_j = \omega_j \quad (j = 1, 2, \dots, m);$$

$$s \stackrel{\text{df}}{=} \{(z_i, z_j) \in Z \times Z \mid d_{i,j} \leq 1\};$$

$$k \stackrel{\text{df}}{=} 2k'.$$

Now observe that mapping

$$T \mapsto e(T \in \mathfrak{A}_3(\Omega))$$

defined by

$$e = \{(z_i, z_j) \in Z \times Z \mid u \langle T \rangle(\omega_i, \omega_j) \leq 1\},$$

maps bijectively  $\mathfrak{A}_3(\Omega)$  onto the set of all equivalences on  $Z$  (see the proof of Lemma 4), and preserves the equality

$$\text{card}(e \Delta s) = \text{card}(e \setminus s) + \text{card}(s \setminus e) = 2F(T).$$

Thus

$$\text{card}(e \Delta s) \leq k \quad \text{if and only if} \quad F(T) \leq k',$$

which completes the proof.

*Acknowledgements.* We wish to express our thanks to Dr. I. Havel for his attention to this work and to an anonymous referee for pointing out a considerable simplification of our reduction of  $EC_3$  to  ${}^b\text{HIC}_3$ .

## References

1. Anderberg, M.: Cluster Analysis for Applications. New York: Academic Press 1973
2. Brucker, P.: On the Complexity of Clustering Problems. In: Optimization and Operations Research (R. Henn, B. Korte, W. Oletti eds.), pp. 45-54. Berlin, Heidelberg, New York: Springer 1977
3. Diday, E., Bochi, S., Brossier, G., Celeux, G., Charles, C., Chifflet, R., Darcos, J., Diday, E., Diebolt, J., Fevre, P., Govaert, G., Hanani, C., Jacquet, D., Lechevallier, Y., Lemaire, J., Lemoine, Y., Molliere, J.L., Morisset, G., Ok-Sakun, Y., Rousseau, P., Sankoff, D., Schroeder, A., Sidi, J., Taleng, F.: Optimisation en classification automatique. INRIA, Rocquencourt, 1979

4. Garey, M.R., Johnson, D.S.: *Computers and Intractability: A Guide to the Theory of NP-Completeness*. San Francisco: Freeman 1979
5. Gonzales, T.: On the Computational Complexity of Clustering and Related Problems. In: *System Modelling and Optimization* (R. Drenick, F. Kozin eds.), pp. 174–182. Berlin, Heidelberg, New York: Springer 1982
6. Hartigan, J.A.: *Clustering Algorithms*. New York: John Wiley 1975
7. Hartigan, J.A.: Representation of Similarity Matrices by Trees. *JASA* **62**, 1140–1158 (1967)
8. Jambu, M., Lebeaux, M.-O.: *Cluster Analysis and Data Analysis*. Amsterdam: North-Holland 1983
9. Jardine, N., Sibson, R.: *Mathematical Taxonomy*. New York: John Wiley 1971
10. Johnson, S.C.: Hierarchical Clustering Schemes. *Psychometrika* **32**, 241–254 (1967)
11. Karp, R.M.: Reducibility among Combinatorial Problems. In: *Complexity of Computer Computations* (E.W. Miller, J.W. Thatcher, eds.), pp. 85–104. New York: Plenum Press 1972
12. Křivánek, M., Morávek, J.: On NP-Hardness in Hierarchical Clustering. In: *Proceedings COMPSTAT '84*, pp. 189–194. Vienna: Physica 1984
13. Lerman, I.C.: *Classification et analyse ordinaire des données*. Paris: Dunod 1981
14. Lovász, L.: *Combinatorial Problems and Exercises*. Budapest: Akadémiai Kiadó 1979
15. Spáth, H.: *Cluster Analysis Algorithms for Data Reduction and Classification of Objects*. London: Ellis Horwood 1980
16. Zahn, C.T.: Approximating Symmetric Relations by Equivalence Relations. *SIAM J. Appl. Math.* **12**, 840–847 (1964)

Received April 25, 1985/Október 15, 1985