# An efficient algorithm for a complete link method

D. Defays

*Service de Mathématiques appliquées à la Psychologie, Université de Liège au Sart-Tilman,
Par 4000 Liège 1, Belgium*

An improved algorithm for a complete linkage clustering is discussed. The algorithm is based,
like the algorithm for the single link cluster method (Slink) presented by Sibson (1973), on a
compact representation of a dendrogram: the pointer representation. This approach offers economy
in computation. The algorithm is easily programmable.

## 1. Introduction

Two of the well-known methods of cluster analysis are the
single and complete linkage clustering. The first one was
developed by Florek *et al.* (1951), Sneath (1957), and Johnson
(1971). An optimally efficient algorithm proposed by Sibson
(1973) made its application feasible for a number of OTU's well
into the range $10^3$ to $10^4$. To avoid the extremes of this first
method—the well-known 'chaining' effect—it may be necessary
to apply on the same set of data alternative hierarchic methods.
We showed (Defays, 1975) in a fuzzy sets context that complete
linkage clustering, developed by Lance and Williams (1967),
and Johnson (1967), among others, generates one or some of
the minimal ultrametric dissimilarities superior to the initial
dis-similarity. The present paper provides an efficient algorithm
for carrying out one of the minimal superior ultrametric dis-
similarities. Like Sibson's algorithm, Slink, it enables a
complete link cluster analysis to be applied on an unprece-
dented scale. This method is dependent on the labelling of
the objects. Modifications of the labelling permit us to
obtain different minimal superior ultrametric dissimilarities.

## 2. Notation, terminology and preliminary definitions

A fuzzy relation $R$ is defined as a fuzzy collection of ordered
pairs. If $X = \{i|i = 1, \ldots, N\}$, a fuzzy relation on $X$ is
characterised by a membership function $R(.,.)$ which associ-
ates with each pair $(i, j)$ its 'grade of membership', or in this
case the dissimilarity between $i$ and $j$, $R(i, j) \in [0, \infty]$. In this
paper, we consider fuzzy relations $R$ satisfying the two
conditions:

$$\forall i \in X, R(i, i) = 0 \text{ (reflexivity) },$$

$$\forall i, j \in X, R(i, j) = R(j, i) \text{ (symmetry) }.$$

If $R$ and $Q$ are defined on $X$, the min-max composition of $R$
and $Q$ is denoted by $R \circ Q$ and is defined by

$$R \circ Q(i, j) = \wedge \{(Q(i, k) \vee R(k, j))|k \in X\}, i \in X, j \in X.$$

The $r$-fold composition $R \circ R \circ R \ldots \circ R$ is denoted by $R^r$.
$R$ is transitive if $R^2 \supset R$. We shall call an ultrametric relation,
a fuzzy relation which is reflexive, symmetric and transitive.
Note that the membership function of an ultrametric relation
is an ultrametric dissimilarity. We showed (Defays, 1975) that
if $R$ is a fuzzy reflexive symmetric relation its transitive closure
$\bar{R} = R^{N-1}$ may be obtained by a single linkage clustering and
that complete linkage clustering gives one (or some) minimal
ultrametric relation (MUR) superior to $R$.

Like Sibson, we define a pointer representation as a pair
$(\pi, \lambda)$ of functions $(\pi: 1, 2, \ldots, N \to 1, 2, \ldots, N$ and $\lambda: 1, 2,
\ldots, N \to [0. \infty])$ which satisfies the following conditions:

$$\pi(N) = N.$$
$$\lambda(N) = \infty.$$
$$\forall i < N, i < \pi(i).$$
$$\forall i < N, \lambda(i) < \lambda(\pi(i)).$$

This is the definition given by $R$. Sibson showed, in a slightly
different context, that there is a natural 1-1 correspondence
between pointer representations and ultrametric relations.
Suppose first that $L$ is an ultrametric relation on $X$. Define $\pi, \lambda$
by

$$\pi(N) = N;$$
$$\lambda(N) = \infty;$$

and for $i < N$,

$$\lambda(i) = \wedge \{L(i, j)|j > i\};$$
$$\pi(i) = \vee \{j|L(i, j) = \lambda(i)\};$$

$(\pi, \lambda)$ is called the pointer representation of $L$. Reciprocally,
suppose $(\pi, \lambda)$ is a pointer representation. Define $R$ by its
membership function:

$$R(i, j) = \lambda(i) \text{ if } j = \pi(i) > i;$$
$$= \lambda(j) \text{ if } i = \pi(j) > j;$$
$$= 0 \text{ if } i = j;$$
$$= \infty \text{ otherwise }.$$

$L = \bar{R}$ is an ultrametric relation associated with $(\pi, \lambda)$. It may
be shown that these two transformations are mutually inverse.

## 3. Algorithm

The problem is to find one of the MUR $L$ superior to a reflexive
symmetric fuzzy relation $R . \bar{R}$ and $L$ will then be two extreme
clusterings of $X$. The interest of $L$ is to shade the results obtained
by the single linkage clustering. We shall note $R_n$ the restriction
of $R$ to the first $n$ OTU's $\{1, 2, \ldots, n\}$ of $X$. If $L_n$ is a MUR
superior to $R_n$, as we shall show it, a MUR $L_{n+1}$ superior to
$R_{n+1}$ may be easily obtained from $L_n$. Like in Slink, the reason
for considering a pointer representation is that it can be up-
dated on the inclusion of a new OTU in an efficient way.
Quantities defined on the first $n$ elements will be given subscript
$n$. So, we shall note $(\pi_n, \lambda_n)$ as the pointer representation of a
MUR $L_n$ superior to $R_n$. The present paper gives a method to
generate from $(\pi_n, \lambda_n)$ the pointer representation $(\pi_{n+1}, \lambda_{n+1})$ of
a MUR $L_{n+1}$ superior to $R_{n+1}$.

For given $n$ we define $\mu_n(i)$ recursively on $i$:

$$\mu_n(i) = \vee \{R(i, n + 1), \mu_n(j)|j : \pi_n(j) = i, \lambda_n(j) < \mu_n(j)\}$$

and then, $v_n(n - i)$ which, when unset, will be noted $*$,

$$v_n(n - i) = \mu_n(n - i)$$

if

$$\lambda_n(n - i) \geqslant \vee \{\mu_n(n - i), \mu_n(\pi_n(n - i))\}$$

and if $v_n(\pi_n(n - i)) \neq *$,

$$v_n(n - i) = * \text{ otherwise }.$$

If $a = \vee \{i|\forall j, v_n(j) \geqslant v_n(i) \text{ or } v_n(j) = *\}$, we may define $(\pi, \lambda)$
which we shall prove to be the pointer representation of a
MUR $L_{n+1}$ superior to $R_{n+1}$ as follows:

$$\pi(n + 1) = n + 1;$$
$$\lambda(n + 1) = \infty;$$

$$\pi(a) = n + 1 ;$$
$$\lambda(a) = v_n(a) ;$$
$$\forall k \geqslant 1, \pi(\pi_n^k(a)) = n + 1 ;$$
$$\forall k \geqslant 1, \lambda(\pi_n^k(a)) = \lambda_n(\pi_n^{k-1}(a))$$

if we note $\pi_n^0(a) = a$ and if $\pi_n^{k-1}(a) < n$ and recursively on $i$, if for all $k \geqslant 0$, $i \neq \pi_n^k(a)$:

$$\lambda(i) = \lambda_n(i) ;$$
$$\pi(i) = \pi_n(i)$$

except that if $\pi(\pi_n(i)) = n + 1$ and $\lambda_n(i) \geqslant \lambda(\pi_n(i))$, $\pi(i) = n + 1$.

### Theorem

$(\pi, \lambda)$ is the pointer representation of a MUR $L_{n+1}$ superior to $R_{n+1}$.

### Proof

Let us show first that $(\pi, \lambda)$ is a pointer representation. We have defined $\pi(n + 1) = n + 1$ and $\lambda(n + 1) = 0$. Since $\pi(i) \geqslant \pi_n(i)$, for $i < n$ we have $i < \pi_n(i) \leqslant \pi(i)$. If $i = n$, since $n = \pi_n^k(a)$ for some $k \geqslant 0$, we have $n < \pi(n) = n + 1$. In all cases, if $i < n + 1$, we have $i < \pi(i)$. If $\pi(i) \neq n + 1$, we have $\lambda(i) = \lambda_n(i)$. Then, if $\pi_n(i) < n$, we have $\lambda(i) = \lambda_n(i) < \lambda_n(\pi_n(i))$. If $\lambda(\pi(i)) = \lambda_n(\pi_n(i))$, then $\lambda(i) < \lambda(\pi(i))$. If $\lambda(\pi(i)) = \lambda_n(\pi_n(i)) \neq \lambda_n(\pi_n(i))$, then $\pi(\pi_n(i)) = n + 1$. But, since $\pi(i) \neq n + 1$, we have $\lambda_n(i) < \lambda(\pi_n(i)) = \lambda(\pi(i))$. If $\pi_n(i) = n$, then $\pi(\pi_n(i)) = n + 1$. But, since $\pi(i) \neq n + 1$, we have $\lambda(i) = \lambda_n(i) < \lambda(\pi_n(i)) = \lambda(\pi(i))$. In all cases, $\pi(i) \neq n + 1$ implies $\lambda(i) < \lambda(\pi(i))$ and $(\pi, \lambda)$ is a pointer representation.

Let us show now that $(\pi, \lambda)$ is the pointer representation of a relation superior to $R_{n+1}$. We shall note $L$ the ultrametric relation associated with $(\pi, \lambda)$. First of all, it is easy to prove that $L(i,j) \leqslant L_n(i,j)$ for $i, j \leqslant n$. For, if $\pi(i) \neq \pi_n(i)$, and if $i = a$ or $i = \pi_n^k(a)$ for some $k \geqslant 1$, we have $L(i, n + 1) \leqslant \lambda(i) \leqslant \lambda_n(i)$ and $L(\pi_n(i), n + 1) \leqslant \lambda(\pi_n(i)) = \lambda_n(i)$. In virtue of transitivity, we must have $L(i, \pi_n(i)) \leqslant \lambda_n(i)$; if $\pi(i) \neq \pi_n(i)$, and $i \neq a$ and $i \neq \pi_n^k(a)$ for all $k \geqslant 1$, we have $\pi_n(i) = a$ or $\pi_n(i) = \pi_n^k(a)$ for some $k \geqslant 1$ and $\lambda_n(i) \geqslant \lambda(\pi_n(i))$. Since $L(i, n + 1) \leqslant \lambda_n(i)$ and $L(\pi_n(i), n + 1) \leqslant \lambda(\pi_n(i)) \leqslant \lambda_n(i)$, in virtue of transitivity, $L(i, \pi_n(i)) \leqslant \lambda_n(i)$. If $\pi(i) = \pi_n(i)$, we have $L(i, \pi_n(i)) \leqslant \lambda_n(i)$. Therefore, we shall have $L(i,j) \leqslant L_n(i,j)$ for all $i, j \leqslant n$. To prove that $L \supset R_{n+1}$, it is sufficient to prove that for all $h \in [0, \infty]$, $L^h \supset R_{n+1}^h$ if we note

$$L^h = \{(i,j) | L(i,j) \leqslant h\}, \quad R_{n+1}^{h} = \{(i,j) | R_{n+1}(i,j) \leqslant h\} .$$

If $L$ is an ultrametric relation, Zadeh (1971) has shown that for all $h \in [0, \infty]$, $L^h$ is an equivalence. Let $C$ be a class of the equivalence $L^h$ and let us prove that for all $i, j \in C$, $R(i, j) \leqslant h$. If $C$ is also a class of the equivalence $L_n^h = \{(i,j) | L_n(i,j) \leqslant h\}$, since $L_n \supset R_n$ the assertion is established. If $C$ is not a class of $L_n^h$, since $L_n$ is a MUR superior to $R_n$, we have $n + 1 \in C$ and $C - \{n + 1\}$ is a class of $L_n^h$. The assertion will be established if we prove that for all $i \in C$, $R(i, n + 1) \leqslant h$. For all $i < n + 1$, we define $\sigma(i)$ to be $\lambda(\pi^{k-1}(i))$ if $\pi^{k-1}(i) \neq \pi^k(i) = n + 1$ and we define $\sigma(n + 1) = 0$. By construction of $L$, $C = \{i | \sigma(i) \leqslant h\}$. This assertion may be easily established. The proof may be found in (Sibson, 1973) and is not given. We suppose $C - \{n + 1\} \neq \phi$. In this case $a \in C$ for, if $i \in C$ and $i \neq n + 1$, we have $\sigma(i) \geqslant \lambda(a)$. Let us show first that if $\pi(j) = n + 1$ and $j \in C$, $R(j, n + 1) \leqslant \mu_n(j) \leqslant h$. If $j = a$, it is true. If $\pi_n(j) = a$ and $\pi(j) = n + 1$, we have $\lambda_n(j) \geqslant \lambda(a)$. If $\mu_n(j) > \lambda_n(j)$, we have $h \geqslant \mu_n(a) = \lambda(a) \geqslant \mu_n(j)$. If $\mu_n(j) \leqslant \lambda_n(j)$, since $\lambda_n(j) = \lambda(j)$ and $\sigma(j) = \lambda(j) \leqslant h$, we have $\mu_n(j) \leqslant h$. In all cases, if $\pi_n(j) = a$ and $\pi(j) = n + 1$, we have $\mu_n(j) \leqslant h$ if $j \in C$. If $j = \pi_n^k(a)$ for some $k \geqslant 1$ and $j \in C$, we have $h \geqslant \sigma(j) = \lambda(j) = \lambda(\pi_n^k(a)) = \lambda_n(\pi_n^{k-1}(a)) \geqslant \mu_n(\pi_n^k(a)) = \mu_n(j)$ for $v_n(\pi^k(a)) \neq *$. If $\pi_n(j) = \pi_n^k(a)$ for some $k \geqslant 1$ and $\lambda_n(j) \geqslant$

$\lambda(\pi_n^k(a))$, we have $\mu_n(j) \leqslant h$ too for if $\mu_n(j) \leqslant \lambda_n(j)$, we have $\mu_n(j) \leqslant \lambda_n(j) = \lambda(j) = \sigma(j) \leqslant h$ and if $\mu_n(j) > \lambda_n(j)$, since $h \geqslant \sigma(j) = \lambda(j) = \lambda_n(j) \geqslant \lambda(\pi_n^k(a))$ we have $\pi_n^k(a) \in C$ and $h \geqslant \mu_n(\pi_n^k(a)) \geqslant \mu_n(j)$. Let us show now that if $\pi(j) \neq \pi^2(j) = n + 1$ and $j \in C$, we have $R(j, n + 1) \leqslant \mu_n(j) \leqslant h$. If $\lambda_n(j) < \mu_n(j)$, recursively we have $h \geqslant \mu_n(\pi(j)) \geqslant \mu_n(j)$; if $\lambda_n(j) \geqslant \mu_n(j)$, since $\pi(j) \neq n + 1$ we have $\pi(j) = \pi_n(j)$ and $\mu_n(j) \leqslant \lambda_n(j) = \lambda(j) < \lambda(\pi(j)) = \sigma(j) \leqslant h$. If $\pi^2(j) \neq \pi^3(j) = n + 1$ and $j \in C$, it can be established as precedently that $h \geqslant \mu_n(j) \geqslant R(j, n + 1)$. Recursively, it can be shown that for all $i \in C$, $h \geqslant \mu_n(i) \geqslant R(i, n + 1)$. Thus we have $L \supset R_{n+1}$.

To complete the proof, we must establish that for every ultrametric relation $L'$ such that $L \supset L' \supset R_{n+1}$ we have $L = L'$. Let us suppose $L'$ to be such an ultrametric relation $(L \supset L' \supset R_{n+1})$ and let us show first that for all $i \leqslant n$, $L'(i, n + 1) \geqslant \mu_n(i)$. If $i = 1$, it is obvious. Suppose it is true for all $i < k \leqslant n$. We shall establish that $L'(k, n + 1) \geqslant \mu_n(k)$. Since $L' \supset R_{n+1}$, we have $L'(k, n + 1) \geqslant R(k, n + 1)$. If $\pi_n(j) = k$ and $\mu_n(j) > \lambda_n(j)$, since $L'(j, k) \leqslant L(j, k) \leqslant L_n(j, k) \leqslant \lambda_n(j) < \mu_n(j)$ and since $L'(j, n + 1) \geqslant \mu_n(j)$ for $j < k$, we have in virtue of transitivity of $L'$, $L'(k, n + 1) \geqslant \mu_n(j)$ and therefore, by construction of $\mu_n(k)$, $\mu_n(k) \leqslant L'(k, n + 1)$. We note $(\pi', \lambda')$ the pointer representation of $L'$. The theorem will be established if we suppose that $L' \neq L$ and we show that it induces an absurdity. Since $L_n$ is a MUT superior to $R_n$, if $L' \neq L$, there exists $i$ with $\lambda'(i) \leqslant \lambda(i)$ and $\pi'(i) = n + 1 \neq \pi(i)$. Since $\mu_n(i) \leqslant L'(i, n + 1) \leqslant \lambda'(i)$, we have $\mu_n(i) \leqslant \lambda(i)$ and for all $k \geqslant 1$, $L'(i, n + 1) \leqslant L'(i, \pi_n^k(i)) \leqslant L_n(i, \pi_n^k(i)) \leqslant \lambda_n(\pi_n^{k-1}(i))$. Thus, in virtue of transitivity of $L'$, we have $\lambda_n(\pi_n^{k-1}(i)) \geqslant L'(\pi_n^k(i), n + 1) \geqslant \mu_n(\pi_n^k(i))$ and by construction of $v_n$, $v_n(i) = \mu_n(i)$. An immediate consequence is $v_n(a) < v_n(i) = \mu_n(i) \leqslant \lambda'(i)$ or $v_n(a) = v_n(i) = \mu_n(i) \leqslant \lambda'(i)$ and $a > i$. In $L'^{\lambda'(i)} = \{(r, s) | L'(r, s) \leqslant \lambda'(i)\}$, $i$ and $a$ must be in the same class. Since $L_n$ is a MUT superior to $R_n$ and since for all $i, j \leqslant n$, $L'_n(i, j) \geqslant L(i, j) \geqslant L'(i, j)$, $i$ and $a$ must be in the same class of $L'^{\lambda'(i)}$. Since $n + 1$ belongs to that class and $\pi(i) \neq n + 1$, it implies $\lambda(i) < \lambda'(i)$ which is absurd. This completes the proof of our theorem.
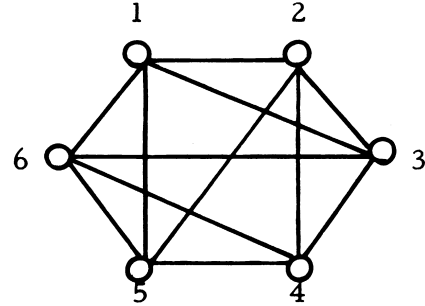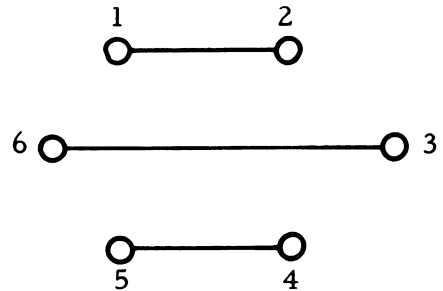


**Fig. 1 Relation R**



**Fig. 2 MUR L superior to R which cannot be obtained with our algorithm**

Notice that the choice of $a = \vee \{i | \forall j, v_n(j) \geqslant v_n(i)$ or $v_n(j) = *\}$ is arbitrary. Other choices are possible.

If we start with $\pi_1(1) = 1$ and $\lambda_1(1) = \infty$, then after $N - 1$ steps of the above recursive process, we shall obtain the pointer representation of a MUR superior to **R**. We have noticed at the beginning of this paper that the result will vary with a relabelling of the elements of $X$.

An interesting question is: is it possible by modifications of the labellings of the OTU's to obtain with our algorithm all the MUR's superior to **R**? Unhappily, this hope is not founded. The very simple example given below proves it. In **Fig. 1**, we link OTU's with dissimilarity 0. The dissimilarity between two not linked OTU's is supposed to be 1. In **Fig. 2**, the relation **L** represented with the same conventions as in Fig. 1, is a MUR superior to the relation **R** of Fig. 1, which obviously cannot be obtained with our algorithm.

If one hopes to obtain a result **L** not too far from the initial relation **R**, a choice of the labels such as

$$\sum_{j=1}^{N} R(1,j) \leqslant \sum_{j=1}^{N} R(2,j) \leqslant \ldots \leqslant \sum_{j=1}^{N} R(N,j)$$

may be judicious.

The results may be presented as in Slink; a conversion of the pointer representation into a packed representation provides readable output.

## 4. The CLINK algorithm

We are going now to give a statement of the algorithm from the computational point of view. We shall try to give it as similarly as possible to the statement of the Slink algorithm. In fact, to perform our algorithm, the subroutine Slink1 of Slink has only to be modified. As in Slink, three arrays of dimension $N$ are used; we shall denote them by $\pi$, $\wedge$, $M$. Suppose that $\pi$, $\wedge$ contain $\pi_n$, $\lambda_n$. The Clink algorithm will change them into $\pi_{n+1}$, $\lambda_{n+1}$ as follows:

1. Set $\pi(n + 1)$ to $n + 1$, $\wedge(n + 1)$ to $\infty$.

2. Set $M(i)$ to $R(i, n + 1)$ for $i = 1, \ldots, n$.

3. For $i$ increasing from 1 to $n$
   if $\wedge(i) < M(i)$
   set $M(\pi(i))$ to max $\{M(\pi(i)), M(i)\}$.
   set $M(i)$ to $\infty$.

4. Set $a$ to $n$.

5. For $i$ increasing from 1 to $n$
   if $\wedge(n - i + 1) \geqslant M(\pi(n - i + 1))$
   set $a$ to $n - i + 1$ if $M(n - i + 1) < M(a)$.
   if $\wedge(n - i + 1) < M(\pi(n - i + 1))$
   set $M(n - i + 1)$ to $\infty$.

6. Set $b$ to $\pi(a)$, $c$ to $\wedge(a)$, $\pi(a)$ to $n + 1$ and $\wedge(a)$ to $M(a)$.

7. If $a < n$
   if $b < n$

set $d$ to $\pi(b)$, $e$ to $\wedge(b)$.
set $\pi(b)$ to $n + 1$, $\wedge(b)$ to $c$.
set $b$ to $d$, $c$ to $e$.
go to 7.
if $b = n$
set $\pi(b)$ to $n + 1$, $\wedge(b)$ to $c$.

8. For $i$ increasing from 1 to $n$
   if $\pi(\pi(i)) = n + 1$
   set $\pi(i)$ to $n + 1$ if $\wedge(i) \geqslant \wedge(\pi(i))$.

## Appendix   A FORTRAN CLINK program

We only give here a subroutine CLINK which must be inserted in the Slink program in the place of the subroutine Slink1. Note that the calling program for the subroutine CLINK must declare TOP which is not declared for Slink1. The measure $\Delta_1$ of classifiability will have a negative value as the result is an ultrametric relation superior to the initial relation.

```
      SUBROUTINE CLINK(NA,HA,HB,I1,NMXOBJ,TOP)
      DIMENSION NA(NMXOBJ),HA(NMXOBJ),HB(NMXOBJ)
      DO 1 J=1,I1
      NEXT=NA(J)
      IF (HA(J)-HB(J))2,1,1
2     H=HB(J)
      HB(J)=TOP
      IF (HB(NEXT)-H)3,1,1
3     HB(NEXT)=H
1     CONTINUE
      IMAX=I1
      DO 4 JI=1,I1
      J=I1-JI+1
      NEXT=NA(J)
      IF (HA(J)-HB(NEXT))6,5,5
5     IF (HB(J)-HB(IMAX))7,4,4
7     IMAX=J
      GO TO 4
6     HB(J)=TOP
4     CONTINUE
      I1S=NA(IMAX)
      H1S=HA(IMAX)
      NA(IMAX)=I1+1
      HA(IMAX)=HB(IMAX)
      IF (IMAX-I1)10,11,11
10    IF (I1S-I1)8,9,9
8     K=NA(I1S)
      HK=HA(I1S)
      NA(I1S)=I1+1
      HA(I1S)=H1S
      I1S=K
      H1S=HK
      GO TO 10
9     NA(I1S)=I1+1
      HA(I1S)=H1S
11    DO 12 J=1,I1
      NEXT=NA(J)
      IF (NA(NEXT)-I1)12,12,13
13    IF (HA(J)-HA(NEXT))12,14,14
14    NA(J)=I1+1
12    CONTINUE
      RETURN
      END
```

## References
DEFAYS, D. (1975).   Ultrametriques et relations floues, *Bull. Soc. Roy. Sc. de Liège*, No. 1-2, pp. 104-118
SIBSON, R. (1973).   An optimally efficient algorithm for the single link cluster method, *The Computer Journal*, Vol. 16, pp. 30-45
ZADEH, L. A. (1971).   Similarity relations and fuzzy orderings, *Inf. Sciences*, No. 3, pp. 177-200