

A FUZZY RELATION APPROACH TO SINGLE LINKAGE

MAMAN A. DJAUHARI *

Abstract. *Single linkage is equivalent to sub-dominant ultrametric. Many algorithms are available for constructing these two objects. But all of them, except one which was proposed by Gondran, are very tedious because of the lack of algebraic structure. Gondran used a special algebraic system as theoretical bases. But it seems rather artificial. In this paper, we propose a more formal approach based on fuzzy relation. The main result presented here is the equivalence between sub-dominant ultrametric and the min-max transitive closure of a symmetric and anti reflexive fuzzy relation. This property enables us to construct an easy and efficient algorithm. At the end of this paper we will find its relationship with Gondran's approach.*

Keywords. *Single linkage, sub-dominant ultrametric and Gondran's algebraic system.*

1. Introduction. Indexed hierarchical clusters (IHC) on a set I , with $\text{card}(I) = n$ and a dissimilarity d on I , is constructed based on the fact that there is a one to one correspondence between the set of all IHCs on I and the set of all ultrametrics on I (see [5], [8]). Single linkage is one of the most popular method for constructing an IHC on I (see

* Department of Mathematics, Institut Teknologi Bandung,
Bandung 40132, Indonesia

[1], [2], [3], [6], [8]). In this method, the distance δ between two clusters A and B is defined by

$$\delta(A, B) = \min_{\substack{i \text{ in } A \\ j \text{ in } B}} d(i, j)$$

We find the important role of single linkage in many applications such as in statistics (especially in cluster analysis and in detecting influential subsets), numerical taxonomy and automatic classification. We know that the algorithms of single linkage are tedious if n is sufficiently large.

Fortunately single linkage is equivalent to sub-dominant ultrametric (see [2]) where, by definition, the sub-dominant ultrametric (USD) u of d is

$$u = \text{Sup} \{ \delta \mid \delta \text{ ultrametric} \leq d \}.$$

These facts motivate us to study the USD by means of fuzzy relation approach. In Section 2 we review some important theorems cited in [4]. The main result of this work is presented in Section 3 and Section 4; here it will be shown that the USD is equivalent to the min-max transitive closure. This property enables us to construct an easy and efficient algorithm. At the end of this paper we will find its relationship with Gondran's approach.

2. Transitive closure. All symbols and properties of fuzzy relation we use here are borrowed from [4] and [5]. Suppose \tilde{R} is a max-min transitive fuzzy relation on I . In other words,

$$\mu_{\tilde{R}}(x, z) \leq \bigvee_y \{ \mu_{\tilde{R}}(x, y) \wedge \mu_{\tilde{R}}(y, z) \}$$

for all x, y, z in I .

Max-min transitivity of a fuzzy relation \tilde{R} can be verified through the following notion of max-min composition \circ of \tilde{R} .

$\tilde{R}^{\circ 2} = \tilde{R} \circ \tilde{R}$ is a fuzzy relation, where

$$\mu_{\tilde{R}}^{02}(x, z) = \bigvee_y \{ \mu_{\tilde{R}}(x, y) \wedge \mu_{\tilde{R}}(y, z) \}$$

for all x, y, z in I .

In [4] we have the following properties for checking the max-min transitivity of \tilde{R} .

Theorem 1. *If $\tilde{R}^{02} = \tilde{R}$, then \tilde{R} is max-min transitive.*

Theorem 2. *\tilde{R} is max-min transitive if and only if $\tilde{R}^{02} \subseteq \tilde{R}$.*

Since I is finite, the max-min transitive closure $\hat{\tilde{R}}$ of \tilde{R} has the following representation :

$$\hat{\tilde{R}} = \tilde{R} \cup \tilde{R}^{02} \cup \tilde{R}^{03} \cup \dots \cup \tilde{R}^{ok},$$

for an integer k ; $1 \leq k \leq n$, where $\tilde{R}^{ok} = \tilde{R} \circ \tilde{R} \circ \dots \circ \tilde{R}$,

k times max-min composition of \tilde{R} .

Those theorems give us the following theorem (see [4]) which is a standard approach to construct the max-min transitive closure.

Theorem 3. *Suppose \tilde{R} is a fuzzy relation on I . Let*

$$l^*(x, y) = \bigvee_{c \text{ di } \zeta} l(c), \text{ where}$$

a). $\zeta = \{c | c = (x = x_{i_1}, x_{i_2}, \dots, x_{i_r} = y)\}$ is a chain

from x to y

b). $l(c) = \mu_{\tilde{R}}(x_{i_1}, x_{i_2}) \wedge \mu_{\tilde{R}}(x_{i_2}, x_{i_3}) \wedge \dots \wedge \mu_{\tilde{R}}(x_{i_{r-1}}, x_{i_r})$.

Then $\mu_{\hat{\tilde{R}}}(x, y) = l^*(x, y)$, for all x and y in I .

In practice, this theorem is still difficult to be

implemented. An easier way, for which \tilde{R} is a dissimilarity, will be shown in the next section.

Like the max-min transitive closure $\hat{\tilde{R}}$, the min-max transitive closure $\check{\tilde{R}}$ of \tilde{R} has the following representation.

$$\check{\tilde{R}} = \tilde{R} \overset{*2}{\cap} \tilde{R} \overset{*3}{\cap} \tilde{R} \cap \dots \cap \tilde{R} \overset{*k}{}$$

for an integer k , $1 \leq k \leq n$, where $\tilde{R}^{*k} = \tilde{R} * \tilde{R} * \dots * \tilde{R}$, k times min-max composition $*$ of \tilde{R} and

$$\mu_{\check{\tilde{R}}^{*2}}(x,z) = \bigwedge_y \{ \mu_{\tilde{R}}(x,y) \vee \mu_{\tilde{R}}(y,z) \}$$

The following theorem enables us to work either with $\hat{\tilde{R}}$ or with $\check{\tilde{R}}$.

Theorem 4. $\check{\tilde{R}} = \check{\check{\tilde{R}}}$ and $\hat{\tilde{R}} = \hat{\hat{\tilde{R}}}$.

Corollary. By De Morgan's rule, we have $\hat{\tilde{R}} = \check{\check{\tilde{R}}}$ and $\check{\tilde{R}} = \hat{\hat{\tilde{R}}}$.

3. Equivalence between USD and min-max transitive closure.

Suppose D is a dissimilarity matrix associated with dissimilarity index d on I . It means that the (i,j) th element of D is $d(i,j)$, where i and j in I . Suppose also that U is an ultrametric matrix on I . Then the (i,j) th element of U is $u(i,j)$, where u is an ultrametric on I . From fuzzy relation's point of view, we know that

(i) D is a symmetric and anti-reflexive fuzzy relation,

where

$$\mu_D(x,y) = d(x,y) \text{ for all } x \text{ and } y \text{ in } I.$$

(ii) U is a symmetric, anti-reflexive and min-max transitive fuzzy relation. It means that

$$\mu_U(x,z) \leq \bigwedge_y \{ \mu_U(x,y) \vee \mu_U(y,z) \}$$

for all x, y and z in I , where $\mu_U(x,y) = U(x,y)$.

In the next paragraph, we propose two propositions that enable us to find the USD more easily. These two propositions are the main result of this work.

Proposition 1. *If \bar{R} is a dissimilarity on I , then $\bar{R}^{\vee} = \bar{R}^{*k}$ for an integer k ; $1 \leq k \leq n$.*

Proof.

From **Theorem 4** we know that $\bar{R}^{\vee} = \bar{R}$. Hence,

$$\begin{aligned} \bar{R} &= \overline{(\bar{R} \cup \bar{R}^{o2} \cup \bar{R}^{o3} \cup \dots \cup \bar{R}^{ok})} \text{ for an integer } k; 1 \leq k \leq n. \\ &= \overline{(\bar{R} \cap \bar{R}^{o2} \cap \bar{R}^{o3} \cap \dots \cap \bar{R}^{ok})} \text{ by De Morgan's rule.} \end{aligned}$$

We know that

$$\begin{aligned} \mu_{\bar{R}^{o2}}(x,y) &= \max_{(a,b)} \{ \mu_{\bar{R}^{o2}}(a,b) \} - \mu_{\bar{R}^{o2}}(x,y) \\ &= \alpha - \max_z \{ \min \{ \mu_{\bar{R}}(x,z), \mu_{\bar{R}}(z,y) \} \} \\ &\quad \text{where } \alpha = \max_{(a,b)} \{ \mu_{\bar{R}^{o2}}(a,b) \} \\ &= \mu_{\bar{R}}(x,x) \\ &= \alpha - \max_z \{ \alpha - \max \{ \alpha - \mu_{\bar{R}}(x,z), \alpha - \mu_{\bar{R}}(z,y) \} \} \end{aligned}$$

$$\begin{aligned}
 &= \alpha - \max_z \{ \alpha - \max_{\tilde{R}} \{ \mu_{\tilde{R}}(x,z), \mu_{\tilde{R}}(z,y) \} \} \\
 &= \min_z \{ \max_{\tilde{R}} \{ \mu_{\tilde{R}}(x,z), \mu_{\tilde{R}}(z,y) \} \} \\
 &= \mu_{\tilde{R}^{*2}}(x,y)
 \end{aligned}$$

Hence $\overline{R^{o2}} = \tilde{R}^{*2}$.

In general we have $\overline{R^{om}} = \tilde{R}^{*m}$; $m = 1, 2, \dots, k$. Hence,

$$\tilde{R} = \tilde{R} \cap \tilde{R}^{*2} \cap \dots \cap \tilde{R}^{*k}$$

for an integer k ; $1 \leq k \leq n$.

Now we show that the right hand side equals \tilde{R}^{*k} .

By definition,

$$\mu_{\tilde{R}^{*2}}(x,z) = \bigwedge_y \{ \mu_{\tilde{R}}(x,y) \vee \mu_{\tilde{R}}(y,z) \},$$

for all x, y and z in I . Especially if $y = z$, then

$$\mu_{\tilde{R}^{*2}}(x,z) \leq \mu_{\tilde{R}}(x,z) \vee \mu_{\tilde{R}}(z,z)$$

But $\mu_{\tilde{R}}(z,z) = 0$, because \tilde{R} is a dissimilarity. Hence,

$$\mu_{\tilde{R}^{*2}}(x,z) \leq \mu_{\tilde{R}}(x,z)$$

for all x and z in I or $\tilde{R}^{*2} \subseteq \tilde{R}$.

In general we have $\tilde{R}^{*k} \subseteq \dots \subseteq \tilde{R}^{*3} \subseteq \tilde{R}^{*2} \subseteq \tilde{R}$.

It implies that $\overset{\vee}{\underset{\sim}{R}} = \tilde{R}^{*k}$ and the proof is complete.

Corollary. $\overset{\vee}{\underset{\sim}{R}} = \tilde{R}^{*k} \subseteq \tilde{R}$ or $\mu_{\overset{\vee}{\underset{\sim}{R}}}(x,y) \leq \mu_{\tilde{R}}(x,y)$

for all x and y in I .

The second proposition will show us how to construct easily the USD or to realize easily the single linkage.

Proposition 2. If \tilde{R} is a dissimilarity on I , then $\overset{\vee}{\underset{\sim}{R}}$ is the USD of \tilde{R} .

Proof.

Theorem 3 tells us that

$$\mu_{\overset{\vee}{\underset{\sim}{R}}}(x,y) = \mu_{\tilde{R}}^{\wedge}(x,y) = \max_{c \text{ di } \mathcal{C}} l(c)$$

where $c = (x = x_{i_1}, x_{i_2}, \dots, x_{i_r} = y)$ is a chain from x to y .

If $\alpha = \mu_{\tilde{R}}(x,x)$ for all x in I , then

$$\begin{aligned} \mu_{\overset{\vee}{\underset{\sim}{R}}}(x,y) &= \max_c \{ \min_k \{ \mu_{\tilde{R}}(x_{i_k}, x_{i_{k+1}}) \} \} \\ &= \max_c \{ \min \{ \mu_{\tilde{R}}(x_{i_1}, x_{i_2}), \dots, \mu_{\tilde{R}}(x_{i_{r-1}}, x_{i_r}) \} \} \\ &= \max_c \{ \alpha - \max \{ \alpha - \mu_{\tilde{R}}(x_{i_1}, x_{i_2}), \dots, \alpha - \mu_{\tilde{R}}(x_{i_{r-1}}, x_{i_r}) \} \} \end{aligned}$$

$$\begin{aligned}
&= \max_c \{ \alpha - \max \{ \mu_{\bar{R}}(x_{i_1}, x_{i_2}), \dots, \mu_{\bar{R}}(x_{i_{r-1}}, x_{i_r}) \} \} \\
&= \alpha - \min_c \{ \alpha - \{ \alpha - \max \{ \mu_{\bar{R}}(x_{i_1}, x_{i_2}), \dots, \mu_{\bar{R}}(x_{i_{r-1}}, x_{i_r}) \} \} \} \\
&= \alpha - \min_c \{ \max_k \{ \mu_{\bar{R}}(x_{i_k}, x_{i_{k+1}}) \} \}
\end{aligned}$$

This equality shows that

$$\mu_{\bar{R}}^{\vee}(x, y) = \min_c \{ \max_k \{ \mu_{\bar{R}}(x_{i_k}, x_{i_{k+1}}) \} \}$$

Now we will show that \bar{R}^{\vee} is the USD of \bar{R} .

(i). It is clear that $\mu_{\bar{R}}^{\vee}(x, y) \leq \mu_{\bar{R}}(x, y)$ for all x and y

in I , since $\bar{R}^{\vee} = \bar{R}^{*k} \subseteq \bar{R}$ (see **Corollary to Proposition 1**).

(ii). If $c = (x = x_{i_1}, x_{i_2}, \dots, x_{i_r} = y)$ is a chain from x

to y , we note that $L(c) = \max_k \{ \mu_{\bar{R}}(x_{i_k}, x_{i_{k+1}}) \}$.

Suppose c_1 is a chain from x to y and c_2 is a chain from y to

z , such that $\mu_{\bar{R}}^{\vee}(x, y) = L(c_1)$ and $\mu_{\bar{R}}^{\vee}(y, z) = L(c_2)$.

Suppose also that c_3 is a chain from x to z , constructed from c_1 and c_2 such that

$$L(c_3) = \max \{ L(c_1), L(c_2) \}$$

In this case,

$$L(c_3) = \max \left\{ \underset{\sim}{\mu}_{\bar{R}}^{\vee}(x,y), \underset{\sim}{\mu}_{\bar{R}}^{\vee}(y,z) \right\}$$

Then we have

$$\begin{aligned} \underset{\sim}{\mu}_{\bar{R}}^{\vee}(x,z) &= \min_{c \text{ di } \zeta} L(c) \leq L(c_3) \\ &\leq \max \left\{ \underset{\sim}{\mu}_{\bar{R}}^{\vee}(x,y), \underset{\sim}{\mu}_{\bar{R}}^{\vee}(y,z) \right\} \end{aligned}$$

It implies that $\underset{\sim}{\mu}_{\bar{R}}^{\vee}$ is an ultrametric on I.

(iii). Suppose U is the USD of $\underset{\sim}{\mu}_{\bar{R}}^{\vee}$. Now we show that $U = \underset{\sim}{\mu}_{\bar{R}}^{\vee}$. Consider a chain $c_1 = (x = x_{i_1}, x_{i_2}, \dots, x_{i_r} = y)$ from x to y where $\underset{\sim}{\mu}_{\bar{R}}^{\vee}(x,y) = L(c_1)$. Then,

a). $\mu_U(x,y) \leq \max \{ \mu_U(x,z), \mu_U(y,z) \}$ for all x, y and z in I, because U is an ultrametric. Especially,

$$\mu_U(x,y) \leq \max \{ \mu_U(x, x_{i_k}), \mu_U(x_{i_k}; y) \}$$

for all $k = 1, 2, \dots, r$. Hence,

$$\begin{aligned} \mu_U(x,y) &\leq \max \{ \mu_U(x, x_{i_2}), \mu_U(x_{i_2}; y) \} \\ &\leq \max \{ \mu_U(x, x_{i_2}), \max \{ \mu_U(x_{i_2}, x_{i_3}), \mu_U(x_{i_3}, y) \} \} \\ &\leq \max \{ \mu_U(x, x_{i_2}), \mu_U(x_{i_2}, x_{i_3}), \mu_U(x_{i_3}, y) \} \end{aligned}$$

In general we have

$$\mu_U(x, y) \leq \max_k \{ \mu_U(x_{i_k}, x_{i_{k+1}}) \}, \quad 1 \leq k \leq r - 1.$$

b). U is the USD of \bar{R} . Then by definition, $U \subseteq \bar{R}$ or

$$\mu_U(x, y) \leq \mu_{\bar{R}}(x, y), \quad \text{for all } x \text{ and } y \text{ in } I.$$

From a) and b), we have

$$\begin{aligned} \mu_U(x, y) &\leq \max_k \{ \mu_{\bar{R}}(x_{i_k}, x_{i_{k+1}}) \}, \quad 1 \leq k \leq r - 1 \\ &\leq L(c_1) = \mu_{\bar{R}}(x, y) \text{ or} \end{aligned}$$

$$\mu_U(x, y) \leq \mu_{\bar{R}}(x, y).$$

It has been shown that \bar{R} is an ultrametric and U is the USD of \bar{R} . Hence the inequality $\mu_U(x, y) \leq \mu_{\bar{R}}(x, y)$ gives us

$$\mu_U(x, y) = \mu_{\bar{R}}(x, y) \text{ or } U = \bar{R}. \quad \text{It implies that } \bar{R} \text{ is the USD}$$

of \bar{R} and the proof is complete.

4. Conclusion. Proposition 2 shows us how to construct easily the single linkage. Its simple algorithm is given in Proposition 1. In practice, if D is a dissimilarity matrix, then the single linkage on D can be easily found by constructing the sequence D^2, D^4, D^8, \dots etc, where matrix multiplication is defined as usual but the sum and the product of two real numbers a and b are defined by $a + b = \min \{a, b\}$ and $a \cdot b = \max \{a, b\}$.

If $D^{2^{(k+1)}} = D^{2^k}$ for an integer k , then D^{2^k} is the result for single linkage on D .

5. Relationship with Gondran's approach.

Consider the algebraic system $(R^+, +, \cdot)$ where

$$a + b = \min \{a, b\} \text{ and } a \cdot b = \max \{a, b\}$$

for all a and b in R^+ . If \mathcal{G} is the set of all $n \times n$ matrices over $(R^+, +, \cdot)$, the sum \oplus and the product $*$ of two matrices are defined as usual, we know that $(\mathcal{G}, \oplus, *)$ is the Gondran's algebraic system (see [3], [8], and [9]). By means of the two propositions we find that this algebraic system is the system of fuzzy relation with the two operations V and \wedge . Furthermore the sequence D^2, D^4, D^8, \dots etc., and thus the solution D^{2^k} where $D^{2^{(k+1)}} = D^{2^k}$ for an integer k , is equal to the solution for single linkage on D given by Gondran (see again [3], [8], and [9]).

REFERENCES

1. Benzecri J.P. *L'analyse des données; la taxinomie*. Dunod - Paris, 1980.
2. Caillez F. and Pages J.P. *Introduction à l'analyse des données*. Smash - Paris, 1976.
3. Jambu M. *Classification automatique pour l'analyse des données*. Dunod - Paris 1978.
4. Kaufmann A. *Introduction à la théorie des sousensemble flous; element théoriques de base*. 2eme edition. Masson Paris 1977.
5. Kaufmann A. *Introduction à la théorie des sousensemble flous; application à la classification et à la reconnaissance des formes aux automates et aux system; aux choix des critères*. Masson - Paris 1975.
6. Jardine N. and Sibson P. *Mathematical taxonomy*. John Wiley and Sons 1977.
7. Pal S.K. and Majumder D.K.D. *Fuzzy Mathematical approach to pattern recognition*. Wiley Eastern Ltd 1986.
8. Roux. *Classification automatique*. Ecole d' Ete d'Analyse Numérique, Paris 1975.

9. Van Cutsem B. *Ultrametric, distance, \emptyset - distances maximum dominées par une dissimilarité donnée*. Statistique et Analyse des Données, Vol 8 No. 2, 1983.