

# Three Dissimilarity Measures to Contrast Dendrograms

Guillermo Restrepo

Laboratorio de Química Teórica, Universidad de Pamplona, Pamplona, Colombia, and Environmental Chemistry and Ecotoxicology, University of Bayreuth, Bayreuth, Germany

Héber Mesa

Departamento de Matemáticas, Universidad del Valle, Cali, Colombia

Eugenio J. Llanos

Corporación SCIO, Carrera 13 No. 13 - 24 oficina 721, Bogotá, Colombia

Received November 19, 2006

We discussed three dissimilarity measures between dendrograms defined over the same set, they are triples, partition, and cluster indices. All of them decompose the dendrograms into subsets. In the case of triples and partition indices, these subsets correspond to binary partitions containing some clusters, while in the cluster index, a novel dissimilarity method introduced in this paper, the subsets are exclusively clusters. In chemical applications, the dendrograms gather clusters that contain similarity information of the data set under study. Thereby, the cluster index is the most suitable dissimilarity measure between dendrograms resulting from chemical investigation. An application example of the three measures is shown to remark upon the advantages of the cluster index over the other two methods in similarity studies. Finally, the cluster index is used to measure the differences between five dendrograms obtained when applying five common hierarchical clustering algorithms on a database of 1000 molecules.

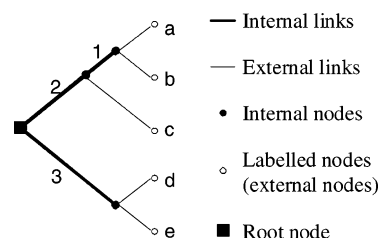
## INTRODUCTION

Hierarchical cluster analysis,<sup>1</sup> HCA, has become a standard method in searching for similarities among data sets;<sup>2,3</sup> its applications are related to the partitioning of a set into similarity classes<sup>4</sup> that are represented as clusters. HCA constitutes a method for classifying the original set with which it is possible to study the behavior of a member of determined class and finally generalize such knowledge to the other members of the class. This procedure endows the set under study with a mathematical structure,<sup>5</sup> namely, a topology.<sup>6–12</sup> In general, a HCA study begins defining the set  $Q$  of work by means of the features of its elements and then looking for the (dis)similarities among the elements using a (dis)similarity function, DF. Afterward, when a grouping methodology, GM, is used, similar elements are clustered and represented graphically in a dendrogram (rooted acyclic-connected binary graph; Figure 1), where clusters appear as branches of the dendrogram.

The total number of dendrograms,<sup>13</sup>  $|F|$ , which can be defined over  $Q$ , whose cardinality  $|Q|$  is  $n$ , grows with  $n$  according to

$$|F| = \frac{(2n - 3)!}{2^{n-2}(n - 2)!}$$

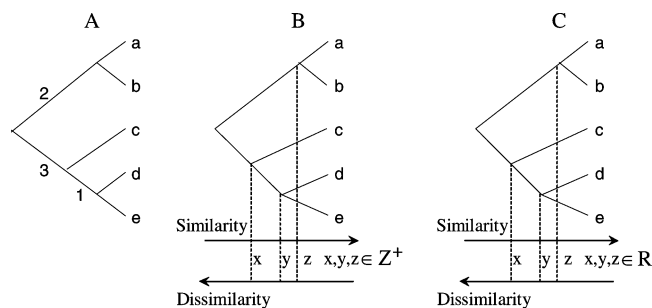
When applying different DFs and GMs (different HCA algorithms) over  $Q$ , several and different results might come up as a consequence of some bias of clustering algorithms toward particular cluster properties<sup>2,3</sup> or as the effect of the lack of “natural clusters”<sup>3</sup> in  $Q$ . Then, a question arising from this discussion is, how can we measure the dissimilarity



**Figure 1.** A dendrogram and types of nodes and links characterizing it (numbers label internal links).

between these results (it is between the corresponding dendrograms)? A contrary situation may occur if the similarity relationships among the elements are strong enough and almost invariant to different HCA algorithms; in this case, even for a set of large cardinality with a large number  $|F|$  of possible dendrograms, it is likely to find similar clusters (natural ones) in their resultant trees. Hence, a possible answer to the question on the contrast of HCA results can be addressed to the contrast of their respective dendrograms by the comparison of their clusters.

Two main mathematical methods have been proposed since the 1960s<sup>14</sup> for contrasting dendrograms. One of them defines a new tree representing a consensus or area of agreement,<sup>14</sup> and the other method defines for any pair of dendrograms a (dis)similarity measure indicating the extent of (dis)agreement.<sup>14</sup> Since this paper deals with dissimilarity measures, we describe them in detail. A detailed discussion on consensus techniques appears in refs 15 and 16. The majority of methods designed to measure dissimilarity between dendrograms have been developed for applications to biological trees (evolutionary trees in the majority of



**Figure 2.** Types of structural dendrograms: (A) bare, (B) ranked, and (C) valued dendrograms.

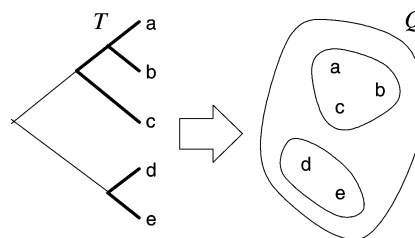
cases), and several of these techniques are related to each other<sup>17</sup> because of the resemblances among structural units of the dendrograms they assess. According to Steel and Penny,<sup>18</sup> there are different structural features of a tree because: “[...] there is no “obvious” or “natural” way to measure the distance between two trees, unlike the comparison of two numbers (where one subtracts the smaller number from the larger)”. Some examples of dissimilarity measures are partition metric<sup>19</sup> or symmetric difference,<sup>15</sup> quartets distance,<sup>14</sup> triples distance,<sup>20</sup> nearest neighborhood interchange metric,<sup>21,22</sup> and some others based on differences in the lengths of the paths between pairs of elements in  $Q$ .<sup>15,23</sup> Although none of these methods actually consider clusters as units to contrast dendrograms, it is crucial, particularly in chemistry, to assess the resemblance between dendrograms using their clusters because they are the pieces of the dendrograms containing the similarity information of  $Q$ . When one applies HCA to a data set, the interpretation of dendrograms is carried out on their clusters; therefore, it would be important to contrast dendrograms on the basis of their clusters. In this paper, we developed a new dissimilarity method dealing exclusively with clusters. This method and the other two discussed in this paper consider a dendrogram as a structure (a graph) able to be decomposed into substructures (subgraphs). A structural classification of dendrograms<sup>23</sup> is given in the following.

**Bare Dendrograms:** They only show the similarity relationships among the elements of  $Q$  without a scale of (dis)similarity (Figure 2A).

**Ranked Dendrograms:** They have internal nodes ranked on an ordinal scale of (dis)similarity (Figure 2B).

**Valued Dendrograms:** They possess internal nodes which have been assigned to a continuous (dis)similarity scale in the real numbers with at least an interval-scale interpretation (Figure 2C).

Bare dendrograms are topological in nature because the relationships (or links) among the elements they cluster are the only features of interest in such structures. Hence, it is irrelevant if the links joining two nodes are longer or shorter. In these kinds of dendrograms, it can only be stated either that “ $a$  is related to  $b$ ”,  $a$  and  $b$  elements being in the same cluster, or that “ $a$  is not related to  $b$ ”, otherwise. Consequently, the notion of an equivalence relation can be attached to the cluster definition, and then a cluster becomes an equivalence class where the mathematical relation is a similarity relationship.<sup>24,25</sup> On the other hand, ranked and valued dendrograms are regarded as geometrical objects where the membership of an element to a cluster is ruled by the presence or absence of links and by the (dis)similarity



**Figure 3.** Three subtrees ( $T$  and the two bold graphs) and their associated subsets ( $Q$  and  $\{a, b, c\}$  and  $\{d, e\}$ ).

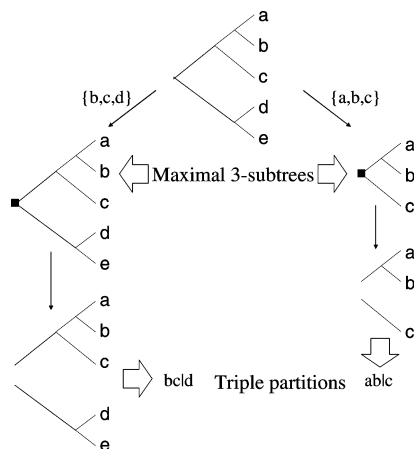
scale. Although they have their differences, these three kinds of dendrograms hold, besides their common graph-classification as complete secondary trees,<sup>26</sup> a metric relationship; all of them are ultrametric trees because any pair  $(r, s)$  of elements in  $Q$  is related by  $d(r, s)$ ,  $d$  being a function fulfilling the metric properties<sup>27</sup> and also the ultrametric property:<sup>23</sup>  $d(r, t) \leq \max\{d(r, s), d(s, t)\}$  with  $r, s, t \in Q$ . For example, for valued dendrograms,  $d(r, s)$  may be defined as the (dis)similarity value of the closest internal node connecting  $r$  and  $s$ ; for ranked trees, it may correspond to the order of the internal nodes and for bare dendrograms to the number of external nodes connected to the internal ones.<sup>23</sup>

Dissimilarity measures between dendrograms, in general, can be divided into two classes:<sup>23</sup> those transforming one dendrogram into another one and those considering a dendrogram as a simpler structure (i.e., sets, partitions, or incidence matrices). The disadvantage of the first ones is that they are often very hard to compute in contrast to the second ones, which are generally quite tractable.<sup>16,23</sup> In the following, we discuss three dissimilarity measures dealing specifically with bare and ranked dendrograms. A discussion on valued trees appears in refs 15 and 23. Two of the analyzed methods in this paper are reported in the literature, and the third one is a novel dissimilarity measure based on the contrast of clusters, which makes it highly appropriate for chemical similarity applications.

#### MEASURING DISSIMILARITY AMONG SETS OF DENDROGRAMS

For the sake of clarity, we define some relevant terms;  $Q$  is the set of objects to classify using HCA;  $|Q|$ , the cardinality of  $Q$ , is represented by  $n$ ;  $T$  is a dendrogram (tree) on  $Q$ ;  $C$  is a cluster of  $T$ ; and  $a|b|c|...$  is the representation of any partition  $\{\{a\}, \{b, c\}, \dots\}$ . When considering  $T$  as a graph (tree), then a cluster may be regarded as a subgraph or subtree of  $T$ . Normally, the extraction of a cluster from  $T$  has two viewpoints: one is to consider the dendrogram as a graph; in that case, the dendrogram becomes a rooted acyclic-connected binary graph (tree), and its clusters may be regarded as its subtrees (A1); the other viewpoint comes from the set theory, and  $T$  is considered as a subset of  $Q$ . Note that a cluster, besides being a proper set of  $Q$  ( $C \subset Q$ ), can also be the same set  $Q$  ( $C = Q$ ). For that reason, in general, we write  $C \subseteq Q$ . Then, having a subtree from  $T$ , a subset in  $Q$  can be associated to it. In this paper, when we refer to a cluster, it is because it has a subtree in  $T$  and also an associated subset in  $Q$ . We show in Figure 3 three subtrees of  $T$  and their corresponding associated subsets. Note that  $T$  is a subtree (A1), and its associated subset is  $Q$ .

**Triples Index.** This method<sup>20</sup> was designed to deal with rooted binary trees (dendrograms), in contrast to the similar



**Figure 4.** Partitioning the triples  $\{a,b,c\}$  and  $\{b,c,d\}$  according to the triples index method (in the example of application of the three dissimilarity indices appears the complete list of triple partitions).

methodology, quartets distance,<sup>14</sup> developed to treat unrooted trees. In triples index,<sup>20</sup>  $h = \{i, j, k\}$  is defined as a **triple**, where  $i, j, k \in Q$ . The total number of triples in  $Q$  is  $t = \binom{n}{3}$ . Given a triple  $h$ , the next step of the procedure is to look for the maximal 3-subtree (A2) containing the elements in  $h$ . Afterward, the root node of the maximal 3-subtree is deleted, inducing a binary partition  $ij|k$  on the triple  $h$ , which is called triple partition TP (Figure 4). Because this dissimilarity measure contrasts pairs of dendrograms, the procedure described previously is carried out over the dendrograms of interest  $T_1$  and  $T_2$ . In order to contrast the triple partition of the triple  $h$  in  $T_1[TP(T_1)]$  and  $T_2[TP(T_2)]$ , the symmetric difference (A3) given by  $SD = TP(T_1) \Delta TP(T_2)$  is calculated.  $SD$  shows the different subsets between  $T_1$  and  $T_2$  regarding the triple  $h$  whose number is given by  $|SD| = |TP(T_1)| + |TP(T_2)| - 2|TP(T_1) \cap TP(T_2)|$  (A3). In general, if two binary partitions  $TP_1$  and  $TP_2$  are built on a set of three elements, there exist only two possible values  $|SD|$  can take, namely, 0 or 4; 0 corresponds to equal partitions ( $TP_1 = TP_2$ ) and 4 to different partitions ( $TP_1 \cap TP_2 = \emptyset$ ). Now, the indicator function  $I_h$  is defined as

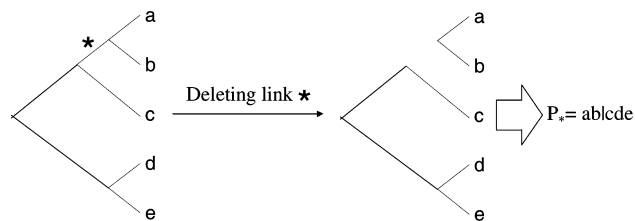
$$I_h = \begin{cases} 1 & \text{if } |SD| = 4 \\ 0 & \text{if } |SD| = 0 \end{cases}$$

$I_h$  yields a value of 1 if the partition of the considered triple  $h$  in  $T_1$  is different from the partition of  $h$  in  $T_2$ . If those partitions are the same, then  $I_h = 0$ . The triples distance between dendrograms  $T_1$  and  $T_2$  is defined as

$$S(T_1, T_2) = \sum_{h=1}^t I_h$$

Hence,  $S(T_1, T_2)$  counts how many triples are different out of the total  $t$ . In order to restrict the values  $S(T_1, T_2)$  can take to the interval  $[0, 1] \in R$ , we normalized  $S(T_1, T_2)$  looking for the maximum and minimum values it can take. Hence,  $\max[S(T_1, T_2)] = t$ , which means that all  $t$  triple partitions in  $T_1$  and  $T_2$  are different. On the other hand,  $\min[S(T_1, T_2)] = 0$ , meaning that all the  $t$  triple partitions in  $T_1$  and  $T_2$  are the same. Having these maximum and minimum values for  $S(T_1, T_2)$ , we define the triples index as

$$\bar{S}(T_1, T_2) = \frac{S(T_1, T_2)}{t}$$



**Figure 5.** Deletion of a link in the partition index method.

**Table 1.** Partitions for  $T_1$  and  $T_2$  (Figures 1 and 2, Respectively) by Deleting Their Internal Links

$r$ th deleted link	$P_r(T_1)$	$P_r(T_2)$
1	$ab cde$	$abc de$
2	$abc de$	$ab cde$
3	$abc de$	$ab cde$

Thus,  $\bar{S}(T_1, T_2)$  is a dissimilarity measure between  $T_1$  and  $T_2$ . In this method, it is considered that a dendrogram is completely determined by the way in which its triples are partitioned in  $Q$ . However, it might be possible to consider other kinds of subtrees, namely, quartets, quintets, and in general  $s$ -tets. Note that  $s = 2$  is not informative, since we look for partitions of the  $s$ -tets to do contrasts; in this case, a duo always yields the same partitions, and then they cannot be used to look for differences between trees. A generalization of this procedure is to look for all the  $\binom{n}{3}$   $s$ -tets in  $Q$  and to study how they are partitioned in the two trees. We are preparing another paper regarding this family of  $s$  indices and their features. We show an example of the application of this method in the next section.

**Partition Index.** This method<sup>19</sup> was conceived to contrast unrooted as well as rooted binary trees (dendrograms). Nowadays, it has been called partition metric and is one of the most known measures of dissimilarity between trees. In fact, it is included in several software packages, for instance, COMPONENT,<sup>28,29</sup> PHYLIP,<sup>30</sup> and PAUP.<sup>31</sup> Its procedure is based on the contrast of partitions generated by removing internal links (Figure 1) of a dendrogram. It does not remove external links (Figure 1) because the generated partitions do not contribute to differentiate the contrasted dendrograms (A4). The first step in this method is the deletion of an internal link  $r$  from  $T$ ,  $n - 2$  being the total number of internal links in a dendrogram.<sup>19</sup> The deletion of this link produces two subgraphs of  $T$  whose associated subsets become a binary partition  $P_r$  of  $Q$ . Consider, for instance, the dendrogram shown in Figure 5; if the link marked \* is removed, then two disjoint subsets are produced:  $\{a,b\}$  and  $\{c,d,e\}$ , which are gathered in  $P^*$  (Figure 5). Note that the deletion of a link always produces a partition  $A|A^C$ , where either  $A$  or  $A^C$  is a cluster (subtree) (A1). The key of this method is the contrast of these binary partitions  $P_r$ 's, and for that reason, it is important to know their number. Although  $n - 2$  internal links yield  $n - 2$   $P_r$ 's, this is not the number of partitions to contrast because there is always a redundant partition. We can explain this analyzing the  $P_r$ 's of the dendrogram in Figure 1, which appear in Table 1 (second column), where  $P_2(T_1) = P_3(T_1)$ . In general, there are always two equal  $P_r$ 's produced by deleting the connected links to the root node (Figure 1). Then, the total number of partitions to contrast is  $n - 2 - 1 = n - 3$ . Now, PT is defined as the collection of partitions to contrast and  $P(T_1, T_2) = PT_1 \Delta PT_2$  as their symmetric difference (A3), where

$P(T_1, T_2)$  yields the different partitions between  $T_1$  and  $T_2$ . The partition metric is defined as the cardinality of  $P(T_1, T_2)$ ,<sup>19</sup> which is given by  $|P(T_1, T_2)| = 2(n - 3) - 2 |PT_1 \cap PT_2|$  because of A3 and  $|PT_1| = |PT_2| = n - 3$ . If  $m$  is assumed as the number of common partitions to  $T_1$  and  $T_2$ , then  $|P(T_1, T_2)| = 2(n - 3 - m)$ . In order to obtain a dissimilarity index for two dendrograms ranging between 0 and 1, we normalized  $|P(T_1, T_2)|$ , which depends on  $m$  because  $n$  is a fixed value for  $T_1$  and  $T_2$ . Hence, the normalization factor of  $|P(T_1, T_2)|$  depends on the minimum and maximum values  $m$  can take.  $\max(m)$  is reached when all partitions are equal for both trees, then  $\max(m) = n - 3$ , and  $\min(m)$  occurs when  $T_1$  and  $T_2$  have the minimum number of common partitions between them; it is  $\min(m) = 0$ . Now,  $\min(m)$  and  $\max(m)$  determine the maximum and minimum values of  $|P(T_1, T_2)|$ , respectively. Thus,  $0 \leq |P(T_1, T_2)| \leq 2(n - 3)$  is obtained, from which  $|P(T_1, T_2)|$  is normalized to the partition index  $PI(T_1, T_2)$  given by

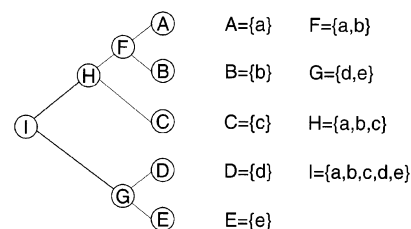
$$PI(T_1, T_2) = 1 - \frac{m}{n - 3}$$

When the partition index takes a value of zero, the contrasted dendrograms have all their partitions in common; if  $PI(T_1, T_2) = 1$ , it is because there are no common partitions to  $T_1$  and  $T_2$ . In appendix A5, we prove that partition index results are the same if the link deletion includes either all the links in  $T$  or only the internal ones. We show an example of the application of this method in the next section.

**Remarks on the Concept of Cluster.** Penny et al.<sup>19</sup> justify the use of partitions as objects to contrast in the partition metric following the results of Waterman and Smith,<sup>22</sup> to whom a tree is represented by the binary partitions produced in the partition metric method. In general, any method extracting structural units from a dendrogram can be considered as a relation  $R$  between  $T$  and the power set of  $Q$ . However, there is not a unique  $R$ ; moreover, it can be possible to have several relations between  $T$  and different aggregations of the  $2^n - 1$  nonempty subsets in  $Q$ . Then, in principle, it can be possible to find several subsets or structural units representing  $T$ . Nevertheless, if the interest in the definition of  $T$  is not concerned exclusively with its representation but also with the similarities shown by  $T$ , then the most appropriate units for reaching this goal are the clusters of  $T$ . In the following, we show how the concept of a cluster can be regarded, first, as a structural unit describing and reconstructing  $T$  and, second, as an equivalence class containing similarity information.

Gusfield<sup>32</sup> has noted that a dendrogram is represented by its clusters, a statement that can be formalized through the concept of intersection graphs<sup>33</sup> as follows: let  $J$  be a collection of sets; the intersection graph of  $J$  is the graph obtained by assigning to each set in  $J$  a distinct vertex. A line is drawn between two vertices if the intersection between the two sets associated with each vertex is nonempty. Hence, the dendrogram shown in Figure 1 can be considered as the intersection graph (Figure 6) of the subsets  $A$  to  $I$  shown in Figure 6. Additionally, according to hypergraph theory,<sup>34</sup> the clusters can be regarded as the vertices of the hypergraph (dendrogram) and the hyperedges as the intersection between any two different clusters.

Although the reconstruction of a graph from the collection of its one-vertex-deleted subgraphs is still an open question



**Figure 6.** A dendrogram as an intersection graph of the subsets  $A$  to  $I$  and as a hypergraph of the subsets  $A$  to  $I$  and their intersections.

in graph theory (reconstruction<sup>35,36</sup> or Ulam's conjecture),<sup>37</sup> its particularization to the case of trees<sup>38</sup> and some other special graphs has been proved.<sup>37</sup> Since a dendrogram is a tree (complete secondary tree)<sup>26</sup> and it can be defined as an intersection graph or as a hypergraph of its clusters, it can therefore be reconstructed from its clusters, which are also obtained from the one-vertex-deleted subgraphs.

The second, and most relevant, viewpoint of a cluster is its ability to show similarities between the elements in  $Q$ , which is the main reason why cluster analysis is broadly used in drug discovery processes and molecular diversity studies.<sup>2</sup> Restrepo and Brüggemann<sup>39</sup> recently showed that, if the similarity between the elements in  $Q$  is regarded as an equivalence relation  $R$ , then the elements in a cluster constitute an equivalence class,  $R$  being a similarity relation. In that case, the set of all the clusters in  $Q$  becomes the quotient of  $Q$  by  $R$  ( $Q/R$ ). In summary, clusters determine similarity classes or similarity neighborhoods in  $Q$ ,<sup>6,10</sup> and also they contain structural information of the dendrogram; for these reasons, they can be used to characterize a dendrogram and also to characterize the similarity relationships in the set  $Q$ , which is one of the targets of HCA in chemistry. In the following, we propose a new method for measuring the dissimilarity between two dendrograms on the basis of the contrast of their clusters.

**Cluster Index.** We understand the question about the dissimilarity between two dendrograms as "how dissimilar are their clusters", because, as we have remarked, the applications of HCA in chemistry are related to the notion of similarity, expressed in the clusters. The cluster index, here described, follows the counting ideas shown in the previous methods using symmetric difference, but it considers as raw material the clusters of the contrasted dendrograms. The total number of clusters in a dendrogram is  $2n - 1$ , as we show in the following proposition.

**Proposition.** Let  $Q$  be a finite set,  $T$  a dendrogram over  $Q$ , and  $RT$  the collection of clusters of  $T$ ; then, the number of clusters in  $T$  is given by  $|RT| = 2n - 1$ .

**Proof.** Let us use induction over  $|Q|$ . If  $Q = \{x_1\}$ , then there exists only one dendrogram with a unique cluster  $\{x_1\}$ . In this case,  $|RT| = 1 = 2(1) - 1 = 2|Q| - 1$ . Suppose  $|Q| = k$  with  $k \leq n$ , then for any dendrogram  $T$ , the number of its clusters is given by  $|RT| = 2k - 1 = 2|Q| - 1$ . Let  $Q$  be a set such that  $|Q| = n + 1$  and  $T$  be a dendrogram. Every  $T$  has a root node (Figure 1) splitting  $T$  into subtrees (clusters)  $T_1$  and  $T_2$ , where  $T_1$  is a dendrogram over a set  $A \subset Q$  and  $T_2$  a dendrogram over  $A^c$ . Since  $|RT|$  is the number of clusters in  $T$ , then  $|RT| = |RT_1| + |RT_2| + 1$ .  $|RT|$  counts the clusters in  $T_1$  and  $T_2$  given by  $|RT_1|$  and  $|RT_2|$  respectively and also the largest cluster corresponding to the complete dendrogram  $T$  that is counted by the addition of 1 in  $|RT|$ .



On the other hand,  $|A| < |Q| = n + 1$  and  $|A^C| < |Q| = n + 1$ ; then,  $|A| \leq n$  and  $|A^C| \leq n$ . Using the hypothesis of induction, we have that  $|RT_1| = 2|A| - 1$  and  $|RT_2| = 2|A^C| - 1$ . Note that  $|A| + |A^C| = |Q|$  because  $A$  and  $A^C$  are disjoint; then,  $|RT| = 2|A| - 1 + 2|A^C| - 1 + 1 = 2|Q| - 1$ . Now, since  $Q$  is made from  $n + 1$  elements, then  $|Q| = n + 1$  and  $|RT| = 2(n + 1) - 1$ .

When contrasting the set of clusters of  $T_1$  with the ones of  $T_2$ , there are always  $n + 1$  common trivial clusters, which are the  $n$  single clusters  $\{x\}$ ,  $x \in Q$ , and the complete set  $Q$ . For this reason, if the goal is to measure the difference between two dendrograms, then these trivial clusters must not be considered, and the total number of clusters to contrast becomes  $n - 2$ . We call  $CT_1$  and  $CT_2$  the clusters to contrast for  $T_1$  and  $T_2$ , respectively, and their contrast is carried out by the symmetric difference  $C(T_1, T_2) = CT_1 \Delta CT_2$ , which yields the different clusters between  $T_1$  and  $T_2$ . Now, we call the cluster metric the cardinality of  $C(T_1, T_2)$  (A3), given by  $|C(T_1, T_2)| = 2(n - 2) - 2|CT_1 \cap CT_2|$ . If we assume  $c$  as the number of clusters common to  $CT_1$  and  $CT_2$ , then  $|C(T_1, T_2)| = 2(n - 2 - c)$ , which depends on  $c$  because  $n$  is a fixed value for  $T_1$  and  $T_2$ . For that reason, we studied the maximum and minimum values  $c$  can reach.  $\max(c) = n - 2$ , meaning that all the possible clusters in  $T_1$  are present in  $T_2$  and  $\min(c) = 0$ . These maximum and minimum values of  $c$  determine the minimum and maximum values of  $|C(T_1, T_2)|$ , respectively, yielding  $0 \leq |C(T_1, T_2)| \leq 2(n - 2)$ . Now, we call the cluster index,  $CI(T_1, T_2)$ , the rank normalization of  $|C(T_1, T_2)|$ , which is given by

$$CI(T_1, T_2) = 1 - \frac{c}{n - 2}$$

When the cluster index reaches a value of zero, it is because the contrasted dendrograms have all their clusters in common; if that value is 1, all of their contrasted clusters are different. One question arising from the  $CI(T_1, T_2)$  expression is whether or not it changes when considering one, two, or in general  $k$  trivial clusters. We show in A6 that  $CI(T_1, T_2)$  yields the same result, when adding  $k$  trivial clusters to the clusters to contrast.

#### EXAMPLE OF APPLICATION OF THE THREE DISSIMILARITY INDICES

In order to show the way in which the three dissimilarity measures here discussed work, we propose the following example. Suppose the dendrograms shown in Figures 1 and 2 defined over  $Q = \{a, b, c, d, e\}$ ; since methods here-described treat bared and ranked trees, then whatever dendrogram, either Figure 2A or B, can be considered in this example.

**Triples Index.** In this case  $t = 10$ ; then,  $Q$  has 10 possible triples, which are listed in Table 2. Partitions  $TP_h$  for  $T_1$  and  $T_2$  are shown in the third and fourth columns in Table 2, respectively. In addition,  $TP(T_1)$  and  $TP(T_2)$  are the sets gathering the cells in the third and fourth columns of Table 2, respectively.

The column labeled  $I_h$  shows the indicator function values for the symmetric difference between each  $h$ th couple of partitions. Thus, it is clear that there are four triples out of 10 for which the partitions are different. Hence,  $\bar{S}(T_1, T_2) = 4/10 = 0.4$ , which means that the dendrograms in Figures 1 and 2 are 40% dissimilar.

**Table 2.** Triples, Their Partitions and Indicator Function Values for Dendrograms of Figures 1 and 2

$h$	triples	$TP_h(T_1)$	$TP_h(T_2)$	$I_h$
1	<i>abc</i>	<i>ab c</i>	<i>ab c</i>	0
2	<i>abd</i>	<i>ab d</i>	<i>ab d</i>	0
3	<i>abe</i>	<i>ab e</i>	<i>ab e</i>	0
4	<i>bcd</i>	<i>bc d</i>	<i>cd b</i>	1
5	<i>bce</i>	<i>bc e</i>	<i>ce b</i>	1
6	<i>cde</i>	<i>de c</i>	<i>de c</i>	0
7	<i>cda</i>	<i>ac d</i>	<i>cd a</i>	1
8	<i>edb</i>	<i>ed b</i>	<i>ed b</i>	0
9	<i>ace</i>	<i>ac e</i>	<i>ce a</i>	1
10	<i>ade</i>	<i>de a</i>	<i>de a</i>	0

**Partition Index.** In this case,  $n - 3 = 2$ ; then, each tree  $T_1$  and  $T_2$  has two possible nonredundant partitions of  $Q$  by removing their internal links. These partitions  $P_r$  are shown in Table 1.

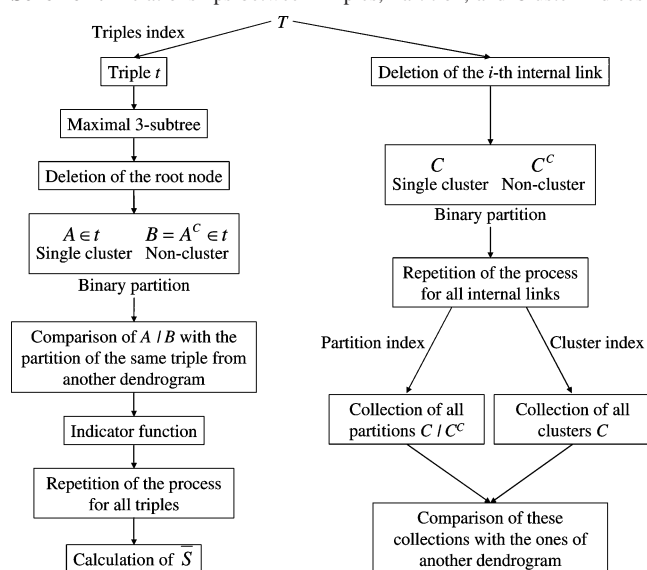
Thus,  $PT_1 = \{ab|cde, abc|de\}$  and  $PT_2 = \{abc|de, ab|cde\}$ , and according to Table 1, we conclude that  $m = 2$ , and then  $PI(T_1, T_2) = 0$ . It means that the two partitions are common to the dendrogram in Figures 1 and 2. In other words, they are 0% dissimilar (100% similar). We give further details about the interpretation of this result in the next section.

**Cluster Index.** The collections of nontrivial clusters for  $T_1$  and  $T_2$  are  $CT_1 = \{\{ab\}, \{de\}, \{abc\}\}$  and  $CT_2 = \{\{ab\}, \{de\}, \{cde\}\}$ . From  $CT_1$  and  $CT_2$ , we conclude that  $c = 2$  since  $\{a, b\}$  and  $\{d, e\}$  are common clusters to  $CT_1$  and  $CT_2$ . Finally,  $CI(T_1, T_2) = 1/3 = 0.\bar{3}$ . In other words, dendrograms in Figures 1 and 2 are about 33% dissimilar.

In summary, the results of the application of the three methods here described are  $\bar{S}(T_1, T_2) = 0.4$ ,  $PI(T_1, T_2) = 0$ , and  $CI(T_1, T_2) = 0.\bar{3}$ , which show that these methods assess different aspects of the structure of a dendrogram, as we have remarked above.

It is surprising to have a value of  $PI(T_1, T_2) = 0$  for the above example because  $T_1$  and  $T_2$  are topologically different and one expects  $PI(T_1, T_2) > 0$ . That is, nodes in  $T_1$  and  $T_2$  have different connectivities. Nevertheless, it is important to state that all the methods here-described assess the dissimilarity between dendrograms, decomposing each tree into subsets and then contrasting those subsets. For that reason  $PI(T_1, T_2) = 0$  does not necessarily mean  $T_1 = T_2$ ; what it really means is  $PT_1 = PT_2$ . That is, the binary partitions of  $T_1$  are the same as the ones of  $T_2$ . In general, each dissimilarity measure here-discussed can be regarded as a function  $d(T_1, T_2)$ , assessing a particular structural aspect of  $T_1$  and  $T_2$ ; then,  $d(T_1, T_2) = 0$  is reached if the “contrasted units” (not the dendrograms as they are) are exactly the same for  $T_1$  and  $T_2$ . Dendrograms are mathematically represented by subsets, and a complete match of these subsets ought to be interpreted only in terms of the subsets. It is not correct to state that two dendrograms are the “same” because one particular dissimilarity measure yields  $d(T_1, T_2) = 0$ . Perhaps, if the idea of a holistic dissimilarity measure is searched, that is, if we want to consider several structural aspects of the dendrogram, then a first attempt for having such a measure is the combination of  $f$  specific dissimilarity measures in this way:

$$D(T_1, T_2) = \frac{\sum_{i=1}^f w_i d_i(T_1, T_2)}{\sum_{i=1}^f w_i}$$

**Scheme 1.** Relationships between Triples, Partition, and Cluster Indices

where  $d_i(T_1, T_2)$  is a dissimilarity measure and  $w_i$  is a weight factor sizing the priority or the importance of the  $i$ th dissimilarity measure (note that  $D(T_1, T_2)$  is based on the general similarity coefficient suggested by Gower).<sup>40</sup>

#### SIMILARITIES AND DIFFERENCES AMONG THE THREE METHODS

It is possible to study the relationships between the three dissimilarity methods, analyzing their procedures. We summarize in Scheme 1 the three methodologies, showing how each one works on a dendrogram  $T$ . The triples index contrasts dendrograms, studying their triples binary partitions  $A|B$ , where  $A$  is always a single cluster (Scheme 1) and the other part,  $B$ , is not a cluster. The partition index also contrasts binary partitions  $C|C^c$  of  $Q$ , and  $C$  is always a cluster. Thus, both, partition and triples indices, always contrast partitions including at least one cluster. Nevertheless, the whole contrast is not based exclusively on clusters but on a mixture of 50% clusters and 50% nonclusters. In other words, 50% of the raw material of triples and partition indices contains similarity information, and the other 50% does not. For this reason, their results cannot be interpreted in terms of similarities between the elements in  $Q$ . On the other hand, according to Scheme 1, partition and cluster indices can be regarded as similar in the initial step. If we consider the process of looking for clusters as a link deletion, then both methods remove links on the dendrogram and produce partitions  $C|C^c$  where  $C$  is a cluster and  $C^c$  is not. The difference between these two methods begins when the partition index considers, as units to contrast two dendrograms, the complete partition  $C|C^c$ ; thereby, it is attached to the similarity information, expressed as  $C$ , nonsimilarity information, given by  $C^c$ . On the contrary, the cluster index only takes clusters  $C$  as units to contrast, discarding  $C^c$ , that is, discarding the nonsimilarity information of  $Q$ . We mentioned above some remarks on the concept of cluster and its importance in expressing neighborhood and similarity relationships among its elements. Accordingly, if we want to build a method looking for dissimilarity between dendrograms on the basis of their clusters, then the most appropriate of the three methods here discussed is the cluster index

because it is the only one completely based on the contrast of clusters, without adding other kinds of mathematical objects which "contaminate" the similarity expressed in the clusters. In few words, the cluster index is created to interpret a dendrogram as a collection of clusters, and it always contrasts exclusively clusters.

By the analysis of dendrograms  $T_1$  and  $T_2$  shown in Figures 1 and 2, respectively, we can see that  $CT_1 = \{ab, de, abc\}$  and  $CT_2 = \{ab, de, cde\}$ . Hence, these dendrograms are differentiated by the clusters  $abc$  and  $cde$ . However, if we attach to the clusters in  $CT_1$  and  $CT_2$  their complements, then we obtain  $PT_1$  and  $PT_2$  in the partition metric, which are the sets including the binary partitions of the first and second columns of Table 1, respectively. Then, the discriminatory power of  $abc$  and  $cde$  is lost since  $abc$  (in the cluster metric) becomes  $abc|de$  (in the partition metric) and  $cde$  becomes  $cde|ab$ . The ability to differentiate disappears because  $abc|de = P_2(T_1)$  (Table 1) is also present in  $T_2[P_1(T_2)]$  (Table 1). Similarly,  $P_1(T_1) = ab|cde$  (Table 1) becomes  $P_2(T_2)$  in  $T_2$  (Table 1). Thus, the discriminatory power embedded in the clusters is lost when attaching their complements to them. In contrast to the partition index, the cluster index keeps the discriminatory power of the clusters.

In summary, all three methods work with clusters, the triples index with single clusters and the partition and cluster indices with the same clusters; however, triples and partition indices join to their clusters some other subsets not corresponding to clusters. On the other hand, although cluster and partition indices use the same clusters, the attachment of nonclusters to the clusters, in the case of the partition index, makes the results and their meaning change. Thus, the advantage of contrasting trees using the *similarity sense of clusters is lost when applying triples and partition indices, and it is kept in cluster index*.

#### CHEMICAL APPLICATION OF THE CLUSTER INDEX

In the previous discussion the argument arose that, in HCA chemical applications, any contrast of dendrograms must deal with the contrast of their clusters because the clusters are the entities gathering the similarity chemical information. In this section, we calculate the dissimilarity of different HCA algorithms over a chemical database; since triples and partition indices do not operate entirely on the clusters of the HCA results, then these methods are not considered in this example.

A set of 1000 molecules was randomly selected from the National Cancer Institute, NCI, database,<sup>41</sup> and they were represented by 1024-bit Barnard Chemical Information, BCI, fingerprints;<sup>42–44</sup> their similarities were calculated using the Tanimoto coefficient.<sup>4</sup> Five different GMs were applied to the Tanimoto similarity matrix yielding five dendrograms, which, because of their large size, cannot be displayed in this manuscript; however, their electronic ASCII files can be requested from G. Restrepo. The GMs employed were<sup>9</sup> single (sing), complete (comp), centroid (cent), and unweighted average (unav) linkages and Ward's method, all of them members of the sequential agglomerative hierarchical nonoverlapping methods.<sup>45</sup> The cluster index values for the contrast of the 10 pairs of dendrograms appear in Table 3.

There are 1999 clusters in each dendrogram, 1001 of which are trivial ones, and 998 are considered in the calculations

**Table 3.** Cluster Index Results for the Contrast of Five Dendrograms Obtained from the Combination of the Tanimoto Coefficient and Five Grouping Methodologies

	centroid linkage	unweighted average linkage	complete linkage	single linkage	Ward's method
centroid linkage	0				
unweighted avg linkage	0.679	0			
complete linkage	0.618	0.545	0		
single linkage	0.769	0.791	0.765	0	
Ward's method	0.669	0.511	0.537	0.787	0

of the cluster dissimilarity index. From Table 3, it can be seen that the five HCA results yield rather different outcomes since their dissimilarities are greater than 0.5, which means that more than 50% of the clusters in each dendrogram are different from those in the other four trees. In other words, more than 499 clusters are different in any contrast of dendrograms. Keeping in mind these differences, the lowest difference occurs for the couple unav–Ward (51% dissimilarity), with 488 common clusters, while the largest difference results for the couple unav–sing (79% dissimilarity), with 209 common clusters. According to Table 3, the ranges of dissimilarities are Ward, [0.511, 0.787]; sing, [0.765, 0.791]; comp, [0.537, 0.765]; unav, [0.511, 0.791]; and cent, [0.618, 0.769]. Their standard deviations are sing, 0.011; cent, 0.054; comp, 0.091; Ward, 0.110; and unav, 0.111. Hence, the most spread dissimilarities are those of unweighted average linkage, and the least ones are those corresponding to a single linkage. The following order of the dissimilarity spread can be set up: sing < cent < comp < Ward < unav. Although sing is the GM with the least spread dissimilarities, it is, in general, the most dissimilar GM when combined with the Tanimoto coefficient because of its high cluster index dissimilarity values (Table 3).

#### SUMMARY, CONCLUSIONS, AND OUTLOOK

We described three different methods (triples, partition, and cluster indices) for calculating dissimilarities between trees; each one assesses the dissimilarity between two dendrograms by the analysis of different structural aspects of a tree. In general, the triples index looks for all the possible sets of three members (triples) contained in  $Q$ , and it contrasts the structural connections of them (maximal 3-subtrees) through the analysis of their binary partitions. This contrast is carried out counting the different partitions between the dendrograms considered. The mathematical background of this method opens the possibility of considering other kinds of  $h$ -tets (quartets, quintets, and so on) to scan the structural dissimilarity between two trees. Nevertheless, these dissimilarities cannot be fully interpreted in terms of the resemblances among the elements in  $Q$  because they are not completely based on the contrast of clusters because some other subsets, which are not clusters, are considered.

The second dissimilarity method was the partition index, which builds all the possible binary partitions of a dendrogram by deleting its internal links, and then it contrasts those partitions counting the different numbers of them. Hence, its result is a measure of how many partitions are different between the two given dendrograms.

The cluster index, a novel dissimilarity measure introduced in this paper, was developed taking into consideration the

lack of methods based on the contrast of the clusters present in two dendrograms. We considered it important to have such kinds of dissimilarity measures because of the two important aspects of the concept of the cluster, namely, the possibility of reconstructing a dendrogram from its clusters and the similarity information gathered in each cluster regarding the elements contained in it. In fact, a dendrogram is mainly used in chemistry for searching similarity classes (clusters) in a given set  $Q$ . Hence, the cluster index is based on the consideration that a dissimilarity measure between dendrograms must be addressed to the assessment of the dissimilarity between their clusters.

On the other hand, through the application of a method contrasting clusters, it is possible to evaluate the following: (1) the effect of the chemical representations, that is, assessing to what extent the HCA results change if the molecules are described by dataprints,<sup>2</sup> that is, real number descriptors, such as topological indices and physicochemical properties, or if they are described by fingerprints,<sup>2</sup> that is, binary strings representing the presence or absence of 2D structural fragments or 3D pharmacophores; (2) the effect of clustering a data set  $Q$  using different dissimilarity functions and similarity coefficients,<sup>4,46</sup> such as Hamming and Soergel distances and Tanimoto and Dice coefficients; (3) the effect of applying different grouping methodologies, such as those mentioned in this paper.

Furthermore, it is possible to assess the combined effect of points 1–3; these contrasts can be done directly measuring the behavior of the similarity neighborhoods in a dendrogram, that is, in their clusters.

We introduced the cluster index as a dissimilarity measure evaluating the behavior of the clusters in two dendrograms. In its mathematical development, it was proved that the number of clusters to contrast between two dendrograms is always  $n - 2$  because the total number of clusters in a dendrogram is  $2n - 1$ , where  $n + 1$  out of  $2n - 1$  are always trivial clusters present in all couples of dendrograms.

By the comparison of the three methods described in this paper, we found that, although the three methods initially consider clusters as units to contrast dendrograms, triples and partition indices mix them with some other mathematical objects different than clusters, which causes them to not be recommended for chemical applications where the final aim is the searching of similarities in a data set, similarities that are contained in the clusters. It was shown that the only method dealing exclusively with clusters is the cluster index; thereby, it is the recommended dissimilarity method to be applied in chemical studies.

A common characteristic of the procedure followed by the three dissimilarity measures here discussed is the use of the operation of symmetric difference. This resemblance underlies the fact of describing each dendrogram as a collection of subsets, which is a constant feature of the methods discussed here. In that case, a mathematical tool for looking for particular differences in sets is the symmetric difference. In fact, this operation has been used in several structural dissimilarity measures, and it is not only restricted to the case of dendrograms; for example, Brüggemann and co-workers<sup>47</sup> have defined a dissimilarity measure between posets (partially ordered sets), describing a poset as a collection of subsets and using the symmetric difference as the tool for looking for differences.



We mentioned that each dissimilarity measure analyzes different structural aspects of a dendrogram and also discussed the possibility of having a general dissimilarity index through the weighted linear combination of different dissimilarity measures.

A chemical application of the cluster index was carried out when analyzing five HCA algorithms which combine the Tanimoto coefficient and five common grouping methodologies; the database selected was a collection of 1000 molecules from the National Cancer Institute. The dissimilarity values showed a high sensitivity against the change of grouping methodology (more than 50% of the clusters in the contrasted dendrograms were different); similar results were obtained by Adamson and Badwen when analyzing a data set of 36 chemicals.<sup>48,49</sup> The high variability found for the Tanimoto–unweighted average linkage algorithm indicates that, in the space of 10 dendrograms here considered, this dendrogram has an intermediate dissimilarity in respect to the other nine dendrograms. The most dissimilar HCA algorithm was the Tanimoto–single linkage, which means it was the dendrogram with more changes in the similarity relationships shown by its clusters.

The potential application of the cluster index was mentioned to assess the effect of varying the parameters defining a HCA study, such as the chemical representation, the dissimilarity measure, and the grouping methodology. Some other applications of this dissimilarity index are, for instance, the quantification of the effect of the ties in proximity in a given chemical data set, where the high probability of having ambiguities in the HCA results is well-known when the size of the data set increases and the (dis)similarity function employed has statistical bias toward particular proximity values.<sup>50</sup> In this particular case, the effect of a decision to overcome a tie can be assessed directly on the effect on the clusters, that is, how many clusters are affected for a particular tie.

An aspect to explore, after defining the cluster index, is the study of the distribution of its values for random trees defined over a particular set  $Q$ . Additionally, it must be proved that  $d(sT_i, sT_j) = 0$  if  $sT_i = sT_j$ , where  $sT_i$  and  $sT_j$  are structural units of  $T_i$  and  $T_j$ . This particular proof must be developed for all the methods here discussed and also for the suggested  $h$ -tets indices.

Finally, triples, partition, and cluster indices deal with bare or ranked trees, that is, with dendrograms where the branch lengths are not considered. In this case, all three methods work on the topology of the trees to be contrasted. However, several applications of HCA and the determination of the number and sort of clusters extracted from a dendrogram are based on the geometrical viewpoint of the dendrogram (valued dendrograms). In those cases, the branch lengths in the dendrogram are important, and it is interesting to develop a dissimilarity measure including these structural aspects of the contrasted dendrograms. There are only two methods<sup>15</sup> measuring dissimilarity between trees that consider branch lengths, but none of them deal exclusively with clusters. A first attempt for reaching this goal is to attach to each cluster in a dendrogram the corresponding length of its associated subtree (the ultrametric height of the cluster). In such a case, the contrast between two dendrograms keeps being based on the clusters and its similarity information but now also includes the geometrical structure of each cluster.

It was mentioned that the representation of a dendrogram as another mathematical object, different from a graph, is computationally more tractable. In this paper, the three dissimilarity measures consider a dendrogram as a collection of subsets, but it is also possible to describe it using adjacency matrices.<sup>23</sup> These kinds of graph representations are well-known in mathematical chemistry and have been widely used in molecular dissimilarity calculations.<sup>51–53</sup> Hence, the description of a dendrogram as an adjacency matrix or as a collection of topological invariants is an alternative possibility for calculating dendrogram dissimilarities, and it would be important to contrast its advantages and disadvantages when compared to the methods here described.

## ACKNOWLEDGMENT

G.R. specially thanks P. Willett and Y. Patel at the Department of Information Studies, University of Sheffield (U. K.), and Barnard Chemical Information Limited (nowadays, Digital Chemistry) for the access they permit to the Tanimoto data matrix used in the chemical example of application of the cluster index. The authors thank the valuable comments of the reviewers of this paper. G.R. thanks COLCIENCIAS and the Universidad de Pamplona in Colombia for the grant offered during the development of this research, and H.M. thanks the Universidad del Valle for the financial support.

## APPENDIX

**A1.** Let  $C$  be a subgraph of the dendrogram  $T$ . It is said that  $C$  is a subtree iff either  $C = T$  or (1)  $C$  does not contain the root node and (2) there is a node  $p$  in  $T$  whose degree is different than 1 such that  $C$  corresponds to one of the connected subgraphs obtained by subtracting  $p$  from  $T$ .

The idea of defining  $T$  as a subtree can be better understood if we consider a tree defined over a subset of  $Q$ ; then, when that subset is  $Q$ , the subtree becomes  $T$ .  $T$  can be considered as a subtree also under the following argument. If a subtree (cluster) is the structural unit grouping elements of  $Q$ , then the subtree grouping all the elements in  $Q$  is  $T$ , which is a subtree.

**A2.** In order to define maximal 3-subtree, we first define 3-subtree. A 3-subtree is a subtree (A1) whose associated set has a cardinality less than or equal to 3. A maximal 3-subtree is any 3-subtree such that it is not possible to find another 3-subtree containing it. Two examples of maximal 3-subtrees for the dendrogram in Figure 1 appear in Figure 3 (bold graphs in the corresponding dendrogram).

**A3.** We define the operation of symmetric difference between two nonempty sets  $A$  and  $B$  by the identity  $A \Delta B = (A \cup B) - (A \cap B)$ . The cardinality of the symmetric difference is given by  $|A \Delta B| = |A \cup B| - |A \cap B| = |A| + |B| - 2|A \cap B|$ .

**A4.** The deletion of the external link joining the element  $x$  [external node (Figure 1)] to the dendrogram produces a partition  $x|(Q - x)$ . Because there are  $n$  elements in a dendrogram, then there are  $n$  partitions of the form  $x|(Q - x)$ . These partitions are always common to all dendrograms defined over  $Q$  and do not contribute to differentiate them.

**A5. Proposition.** Partition index  $PI(T_1, T_2)$  yields the same result if either all the links in  $T$  or only the internal ones are deleted.



*Proof.* The number of external and internal links in  $T$  is  $n$  and  $n - 2$ , respectively. Then, the total number of links in  $T$  is  $2n - 2$ . These links yield  $2n - 2$  binary partitions. However, if we avoid the repeated binary partition produced by deleting the internal link connected to the root node, then the number of nonredundant binary partitions obtained by deleting links is  $2n - 3$ . Hence,  $|P(T_1, T_2)| = 2(2n - 3 - m')$ ,  $m'$  being the number of common partitions between  $T_1$  and  $T_2$ . Now,  $\max(m') = 2n - 3$  and  $\min(m') = n$ , which means that  $\max[|P(T_1, T_2)|] = 2(n - 3)$  and  $\min[|P(T_1, T_2)|] = 0$ . If we call  $PI'(T_1, T_2)$  the value of the partition index when deleting all the links, then  $PI'(T_1, T_2) = (2n - 3 - m')/(n - 3)$ . But  $m' = m + n$ , where  $m$  represents the number of common binary partitions by deleting internal links from  $T$ . Then, we found that  $PI'(T_1, T_2) = 1 - [m/(n - 3)]$ , which is the same result found when deleting only internal links in  $T$  [ $PI(T_1, T_2)$ ].

**A6.**  $T$  has  $2n - 1$  total clusters (including the  $n + 1$  trivial ones), but if the trivial clusters are not considered as objects to contrast, then  $CT_i = n - 2$  for a given  $T_i$ . We called  $c$  the number of clusters common to  $T_1$  and  $T_2$ , and we found  $0 \leq c \leq n - 2$ , which yields  $0 \leq |C(T_1, T_2)| \leq 2(n - 2)$ . Now, if  $k$  trivial clusters are added to the clusters of  $T_i$  to contrast, then  $CT'_i = n - 2 + k$ ,  $CT'_i$  being the set of clusters of  $T_i$  to contrast after adding  $k$  trivial clusters. Now, if we call  $c'$  the number of common clusters between  $T_1$  and  $T_2$ , then  $|C'(T_1, T_2)| = 2(n - 2 + k - c')$ , where  $|C'(T_1, T_2)|$  is the cluster metric for this case. In this situation, we have  $k \leq c' \leq n - 2 + k$ , and for that reason,  $0 \leq |C'(T_1, T_2)| \leq 2(n - 2)$ . Now, if we call  $CI'(T_1, T_2)$  the cluster index when adding  $k$  trivial clusters to the ones to contrast, then  $CI'(T_1, T_2) = (n - 2 + k - c')/(n - 2)$ . But  $c' = c + k$ ; then,  $CI'(T_1, T_2) = 1 - [c/(n - 2)] = CI(T_1, T_2)$ .

## REFERENCES AND NOTES

- (1) Everitt, B. S. *Cluster Analysis*; Edward Arnold: Bristol, U. K., 1993; Chapter 1, pp 1–10.
- (2) Downs, G. M.; Barnard, J. M. *Clustering Methods and Their Uses in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; Wiley-VCH: Weinheim, Germany, 2002; Vol. 18, pp 1–40.
- (3) Handl, J.; Knowles, J.; Kell, D. B. Computational Cluster Validation in Post-Genomic Data Analysis. *Bioinformatics* **2005**, *21*, 3201–3212.
- (4) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 983–996.
- (5) Potter, M. *Set Theory and Its Philosophy*; Oxford University press: Oxford, U. K., 2004; p 72.
- (6) Restrepo, G.; Mesa, H.; Llanos, E. J.; Villaveces, J. L. Topological Study of the Periodic System. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 68–75.
- (7) Restrepo, G.; Villaveces, J. L. From Trees (Dendrograms and Consensus Trees) to Topology. *Croat. Chem. Acta* **2005**, *78*, 275–281.
- (8) Restrepo, G.; Mesa, H.; Villaveces, J. L. On the Topological Sense of Chemical Sets. *J. Math. Chem.* **2006**, *39*, 363–376.
- (9) Restrepo, G.; Llanos, E. J.; Mesa, H. Topological Space of the Chemical Elements and Its Properties. *J. Math. Chem.* **2006**, *39*, 401–416.
- (10) Restrepo, G.; Mesa, H.; Llanos, E. J.; Villaveces, J. L. Topological Study of the Periodic System. In *The Mathematics of the Periodic Table*; King, R. B., Rouvray, D. H., Eds.; Nova: New York, 2006; pp 75–100.
- (11) Daza, M. C.; Restrepo, G.; Uribe, E. A.; Villaveces, J. L. Quantum Chemical and Chemotopological Study of Fourth Row Monohydrides. *Chem. Phys. Lett.* **2006**, *428*, 55–61.
- (12) Mesa, H.; Restrepo, G. On Dendrograms and Topologies. *J. Math. Chem.* **2007**, Submitted.
- (13) Felsenstein, J. The Number of Evolutionary Trees. *Syst. Zool.* **1978**, *27*, 27–33.
- (14) Estabrook, G. F.; McMorris, F. R.; Meacham, C. A. Comparison of Undirected Phylogenetic Trees Based on Subtrees of Four Evolutionary Units. *Syst. Zool.* **1985**, *34*, 193–200.
- (15) Felsenstein, J. *Inferring Phylogenies*; Sinauer Associates: Sunderland, MA, 2004; Chapter 30, pp 528–535.
- (16) Day, W. H. E. Optimal Algorithms for Comparing Trees with Labeled Leaves. *J. Class.* **1985**, *2*, 7–28.
- (17) Penny, D.; Hendy, M. D. The Use of Tree Comparison Metrics. *Syst. Zool.* **1985**, *34*, 75–82.
- (18) Steel, M. A.; Penny, D. Distributions of Tree Comparison Metrics—Some New Results. *Syst. Biol.* **1993**, *42*, 126–141.
- (19) Penny, D.; Foulds, L. R.; Hendy, M. D. Testing the Theory of Evolution by Comparing Phylogenetic Trees Constructed from Five Different Protein Sequences. *Nature* **1982**, *297*, 197–200.
- (20) Critchlow, D. E.; Pearl, D. K.; Qian, C. The Triples Distance for Rooted Bifurcating Phylogenetic Trees. *Syst. Biol.* **1996**, *45*, 323–334.
- (21) Robinson, D. F. Comparison of Labeled Trees with Valency Three. *J. Comb. Theory* **1971**, *11*, 105–119.
- (22) Waterman, M. S.; Smith, T. F. On the Similarity of Dendrograms. *J. Theor. Biol.* **1978**, *73*, 789–800.
- (23) Boorman, S. A.; Olivier, D. C. Metrics on Spaces of Finite Trees. *J. Math. Psychol.* **1973**, *10*, 26–59.
- (24) Rouvray, D. H. Definition and Role of Similarity Concepts in the Chemical and Physical Sciences. *J. Chem. Inf. Comput. Sci.* **1992**, *32*, 580–586.
- (25) Restrepo, G.; Brüggemann, R. Ranking Regions through Cluster Analysis and Posets. *WSEAS Trans. Inf. Sci. Appl.* **2005**, *2*, 976–981.
- (26) Chartrand, G.; Lesniak, L. *Graphs & Digraphs*; Wadsworth & Brooks/Cole: Monterey, CA, 1986; pp 77–83.
- (27) Deza, M. M.; Deza, E. *Dictionary of Distances*; Elsevier: Amsterdam, 2006; p 6.
- (28) Page, R. D. M. *COMPONENT User's Manual (Release 1.5)*; University of Auckland: Auckland, New Zealand, 1989; Chapter 4, pp 4.1–4.7. URL: <http://taxonomy.zoology.gla.ac.uk/rod/cpw.html> (accessed Apr 2007).
- (29) Slowinski, J. B. Review. *Cladistics* **1993**, *9*, 351–353.
- (30) Felsenstein, J. PHYLIP - Phylogeny Inference Package (Version 3.2). *Cladistics* **1989**, *5*, 164–166.
- (31) Swofford, D. L. *PAUP: Phylogenetic Analysis Using Parsimony, Version 3.1*; Illinois Natural History Survey: Champaign, IL, 1993; p 53.
- (32) Gusfield, D. Efficient Algorithms for Inferring Evolutionary Trees. *Networks* **1991**, *21*, 19–28.
- (33) Golumbic, M. C.; Trenk, A. N. *Tolerance Graphs*; Cambridge University Press: Cambridge, U. K., 2004; Chapter 1, p 4.
- (34) Bollobás, B. *Combinatorics: Set Systems, Hypergraphs, Families of Vectors, and Combinatorial Probability*; Cambridge University Press: Cambridge, U. K., 1988; p 1.
- (35) Gross, J.; Yellen, J. *Graph Theory and Its Applications*; CRC Press: Boca Raton, FL, 1999; Chapter 2, pp 64–65.
- (36) Chartrand, G.; Lesniak, L. *Graphs & Digraphs*; Chapman & Hall: London, 1996; pp 50–51.
- (37) Kratsch, D.; Hemaspaandra, L. A. On the Complexity of Graph Reconstruction. *Math. Syst. Theory* **1994**, *27*, 257–273.
- (38) Kelly, P. J. A Congruence Theorem for Trees. *Pac. J. Math.* **1957**, *7*, 961–968.
- (39) Restrepo, G.; Brüggemann, R. Ranking Regions Through Cluster Analysis and Posets. *WSEAS Trans. Inf. Sci. Appl.* **2005**, *2*, 976–981.
- (40) Gower, J. C. A General Coefficient of Similarity and Some of Its Properties. *Biometrics* **1971**, *27*, 857–874.
- (41) The NCI database is available at URL <http://dtp.nci.nih.gov/> (accessed Apr 2007).
- (42) Barnard, J. M.; Downs, G. M. Clustering of Chemical Structures on the Basis of Two-Dimensional Similarity Measures. *J. Chem. Inf. Comput. Sci.* **1992**, *32*, 644–649.
- (43) Barnard, J. M.; Downs, G. M. Chemical Fragment Generation and Clustering Software. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 141–142.
- (44) The BCI software is available from Digital Chemistry Ltd. at URL <http://www.digitalchemistry.co.uk> (accessed Apr 2007).
- (45) Sneath, P. H. A.; Sokal, R. R. *Numerical Taxonomy*; WH Freeman: San Francisco, CA, 1973.
- (46) Gower, J. C. Measures of Similarity, Dissimilarity and Distance. In *Encyclopaedia of Statistical Sciences*; Kotz, S., Johnson, N. L., Read, C. B., Eds.; Wiley: Chichester, U. K., 1982; pp 397–405.
- (47) Brüggemann, R.; Halfon, E.; Welzl, G.; Voigt, K.; Steinberg, C. E. W. Applying the Concept of Partially Ordered Sets on the Ranking

- of Near-Shore Sediments by a Battery of Tests. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 918–925.
- (48) Adamson, G. W.; Bawden, D. Comparison of Hierarchical Cluster Analysis Techniques for Automatic Classification of Chemical Structures. *J. Chem. Inf. Comput. Sci.* **1981**, *21*, 204–209.
- (49) Downs, G. M.; Willett, P. *Similarity Searching in Databases of Chemical Structures*; Lipkowitz, K. B., Boyd, D. B., Eds.; Wiley-VCH: New York, U.S.A., 1996; Vol. 7, pp 1–66.
- (50) MacCuish, J.; Nicolaou, C.; MacCuish, N. E. Ties in Proximity and Clustering Compounds. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 134–146.
- (51) Kvasnicka, V.; Pospichal, J. Fast Evaluation of Chemical Distance by Tabu Search Algorithm. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 1109–1112.
- (52) Diudea, M. V. Molecular Topology. 16. Layer Matrices in Molecular Graphs. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 1064–1071.
- (53) Rücker, C.; Rücker, G.; Meringer, M. Exploring the Limits of Graph Invariant- and Spectrum-Based Discrimination of (Sub)structures. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 640–650.

CI6005189