



## Chapter 5

### Assessing Structural Similarity of Graphs\*

*Timothy E. Goldsmith and Daniel M. Davenport*

In this chapter, we discuss some measures of the similarity of two graphs. Our work was initially motivated by a need to measure the similarity between Pathfinder networks (Schvaneveldt, Durso, & Dearholt, 1985). The problem of how to compare two different representations, however, exists more generally in the fields of scaling and modeling. A central assumption of our work is that representations have structural properties and comparisons of these representations ought to reflect these structural properties. The aim of this current work, then, is to identify a method of assessing graph similarity that is sensitive to structural information.

We begin by describing a particular view of structure and similarity and its implications for comparing graphs. Next, we identify two basic properties of graphs, paths and neighborhoods, and show how each of these properties can be used as a basis for defining graph similarity. We then describe several related measures for assessing graph similarity, discuss some of their properties, and report results of an initial comparison of the measures. Finally, we offer some generalizations and extensions of the measures.

### Similarity and Structure

The basic problem we wish to address is how to measure the similarity of two graphs. More specifically, we would like to define a function that maps any two graphs onto a real number that reflects the graphs' similarity. The set of such functions is very large because similarity itself is not well-defined. Graph similarity is somewhat akin to making human judgments of similarity. Such judgments are inherently subjective because perceived similarity may depend on a multitude of factors including those characteristics of the objects that are psychologically salient to the perceiver and the beliefs the perceiver has about the purpose of the judgment. Similar concerns arise in defining measures of graph similarity.

Consider, for example, the graphs in Figure 1. Graph *A* is a simple binary tree with seven nodes, and graphs *B* and *C* are deviations of *A*; *B* differs from *A* in three edges, whereas *C* differs in just one edge. Which graph, *B* or *C*, is more similar to *A*? There is, of course, no absolutely right answer, and in fact, we will show shortly that either *B* or *C* can be viewed as more similar to *A*.

Notice in Figure 1 that we have arranged the nodes of graphs *B* and *C* in the same spatial layout as in *A*. We assume that the graphs we compare are always composed of a common set of labeled nodes, and so switching node labels to assess similarity is disallowed. This assumption is realistic for those applications of graph theory where the nodes

\*This work, performed in part at Sandia National Laboratories, was supported by the U.S. Department of Energy under contract No. DE-AC04-76DP00789.



of a graph correspond to some specific set of objects under study and therefore transposing nodes would be meaningless. However, we do assume that the edges of the graphs are unlabeled. This assumption is also reasonable for many applications of graph theory. (We discuss at the end of the chapter the case of labeled edges.) By graph similarity, then, we mean the similarity of the patterns of edges that define how two graphs with common node sets are linked. We take as axiomatic that it is the structure of these edge patterns that we wish to measure.

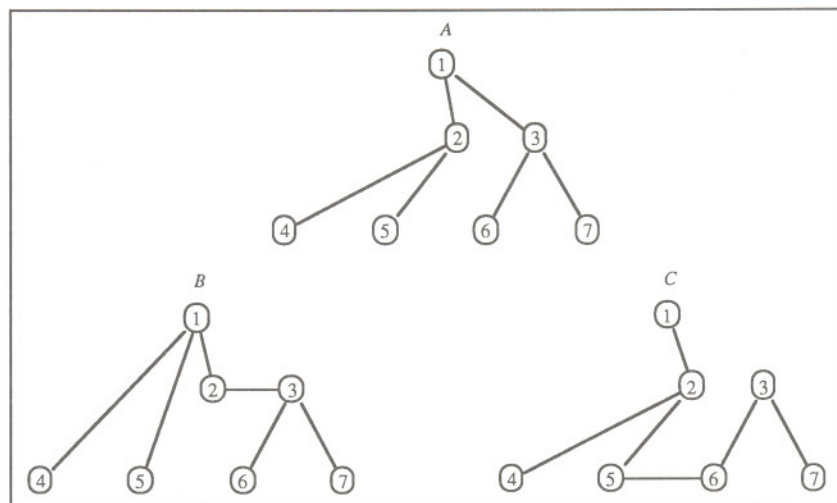


Figure 1. Original Graph A and distortions B and C.

By structure we mean the organization of an object's constituent parts as viewed from the perspective of the whole object. Fundamental to structure is the idea of a relation. An object is defined by the primitive relations on the parts that make it up. An object's structure, in contrast, exists at a level of higher-order relations, or relations on relations. A family tree, for example, is defined by the primitive family relations *son of* and *daughter of*. A higher-order relation on a family tree is the relation *descendant of*. This relation, by our definition, reveals more of the structure of the family tree than either *son of* or *daughter of*. Higher-order relations necessarily appear at a global rather than primitive level. In this regard, structure is an emergent property.

Another way that structure can be viewed is as a collection of subobjects of an object rather than higher-order relations on an object. In this sense structure is viewed as an entity (specifically, a collection of subobjects). Although this distinction in the term may seem subtle, we believe it is important for defining structural similarity. If structure is viewed as an entity, then structural similarity of objects should be assessable by identifying the objects' structural subobjects. On the other hand, if structure is viewed as a property, a measure of structural similarity should compare the objects' structural properties.

Consider first structure as an entity. To define graph similarity under this view, we might begin by identifying specific subentities to serve as the basis for comparing graphs. Subgraphs, such as cycles, stars, and cliques, which are already well-defined and

intuitively represent graph structures, could be used to define an index of similarity by counting the common substructures of two graphs. Graham (1987) describes an approach of Ulam's using a technique based on this idea for measuring graph similarity. Ulam's method is to partition the edges of two graphs into pairwise isomorphic subgraphs; the smaller the number of subgraphs needed to decompose the original graphs, the more similar the graphs. Two graphs with the same number of edges, no matter how dissimilar, can always be partitioned into sets of identical subgraphs by letting each subgraph be a single edge. However, graphs that are more structurally similar will decompose into fewer identical subgraphs, with the extreme case being graphs that are isomorphic, for which *no* partitioning is needed. Ulam's measure of similarity seems appropriate for abstract graphs, however it is problematic for defining similarity for graphs as representations. First, his method assumes that the graphs are unlabeled, whereas the applications of graph theory we consider usually deal with labeled graphs, and second, the task of finding the minimum Ulam decomposition for two graphs is a very difficult computational problem in itself for which there is no known tractable solution.

Consider next the view of structure as a property. Under this perspective, graph similarity might be assessed by comparing two graphs' higher-order property relations, such as distance between two nodes or the constituents of a neighborhood about a node. This is the approach that we adopt, and in the next section we employ these two higher-order relations of graphs to define graph similarity.

### Graph Properties and Similarity

First, we begin with a few basic definitions of graphs and their properties. We limit our discussion to undirected graphs, without loops, and with a common set of labeled nodes. An edge relation is a binary relation defined on pairs of nodes and is the most primitive relation of a graph. All higher-order graph properties are derived from the edge relation. One such property is a path. The distance between two nodes,  $v$  and  $v'$ , is the minimum path length for all paths between  $v$  and  $v'$  provided such a path exists. Because a path is a relation defined on the edge relation, it is a higher-order relation. A graph can be completely described by either its edge relations (i.e., adjacency matrix) or its path lengths (i.e., distance matrix). In the case of undirected graphs, without loops, each of these matrices can be reduced to a vector consisting of only the upper triangular cell entries, and so a graph with  $n$  nodes may be represented by a vector of  $(n^2-n)/2$  distances.

A second higher-order graph property is a neighborhood. A *neighborhood* about some node,  $v$ , is defined to be the set of nodes that are within distance one from  $v$ , excluding  $v$  itself. Referring back to Figure 1, the neighborhoods for Node 1 in graphs A, B, and C are {2, 3}, {2, 4, 5}, and {2}, respectively. By excluding  $v$  from the set, which diverges from the normal definition of a neighborhood, we can simplify our definitions of graph similarity. Notice that a neighborhood is also a relation on the edge relation, and thus a higher-order relation. So, both path length and neighborhood content are, by our definition, structural properties of graphs. Further, these structural properties are sufficiently general such that every graph may be described in terms of them. Next we show how both of these properties can be used to define graph similarity.

First, two graphs may be compared by computing the correlation coefficient of the the two graphs' distance vectors. A correlation coefficient assesses the shared pool of variability between two sets of numbers, standardized by a measure of the total pool of variability. The idea of forming a ratio of shared attributes to total attributes seems intuitively



appealing for measuring similarity, and in fact comparisons of Pathfinder networks with this approach have proved meaningful (e.g., Schvaneveldt, Durso, Goldsmith, Breen, Cooke, Tucker, & DeMaio, 1985).

Turning next to neighborhoods, two graphs may be compared by assessing the similarity of their neighborhoods for corresponding nodes. As with the correlation coefficient, we would like to measure the degree of shared elements relative to some total pool of elements. This may be accomplished with sets by examining their intersection and union. More specifically, an index of similarity for a common node in two graphs is the cardinality of the intersection of the node neighborhoods divided by the cardinality of the union of the neighborhoods. One measure of overall graph similarity is the mean of these  $n$  values. This measure will vary from zero to one with higher values indicating greater similarity.

We now have two ways of defining graph similarity, one comparing path distances with a correlation coefficient and the other assessing neighborhood regions with simple set operations. The next logical question is whether these two measures actually differ in their assessments of graph similarity. Consider again the graphs in Figure 1 and the question of which graph,  $B$  or  $C$ , is more similar to  $A$ . We now have a means for answering this question, and the answer is that  $B$  is closer to  $A$  in terms of path lengths, but  $C$  is closer to  $A$  in terms of neighborhoods. The correlation of path lengths between  $A$  and  $B$  is .79 and between  $A$  and  $C$  is .42, whereas the neighborhood similarities between  $A$  and  $B$  is .43 and between  $A$  and  $C$  is .74. So we come to exactly opposite conclusions about the graph's relative similarity with these two measures. Although there are undoubtedly cases where both approaches would agree, we believe that, in general, path lengths and neighborhoods offer qualitatively different ways of assessing structural similarity of graphs. In the following section, we describe some closely related similarity measures and their properties. (A more formal treatment of these definitions and their properties is given in Appendix A.) Following this we describe the results of a study comparing these various measures.

### Definitions and Properties of Graph Similarity Measures

In this section we describe a class of related graph similarity measures. Let  $C_i(A, B)$  be the similarity between graphs  $A$  and  $B$  with a common labeled node set as measured by  $C_i$  for  $i = 1$  to 8. The first four measures are based on neighborhoods and the second four are based on path lengths. The measures  $C_1$  and  $C_8$  are simply the neighborhood and path length measures, respectively, described above.

Two other measures,  $C_2$  and  $C_3$ , are similar to  $C_1$  and measure the average similarity of neighborhoods. In each case, the cardinality of the intersection of the neighborhoods is divided by some number which normalizes the index.  $C_2$  and  $C_3$  differ from  $C_1$  only in their normalizations.  $C_4$  is the number of edges that match between  $A$  and  $B$  divided by the number of possible matches (i.e., the number of node pairs). It is also one minus the mean absolute difference between entries in the adjacency vectors of  $A$  and  $B$ .

$C_5$  is the mean of the ratios of smallest to largest values for corresponding entries in the two graphs' distance vectors.  $C_6$  is one minus the mean absolute difference between entries in the graphs' distance vectors normalized by the sum of these distances. Finally,  $C_7$  is the correlation coefficient of the graphs' adjacency vectors. The interested reader may refer to Appendix A for formal definitions.

We next evaluate the above measures with respect to several desirable properties for a graph similarity measure. First, the measure should be independent of the size of the node set or the density of the graph. Some of the  $C_i$ 's appear to meet this criterion better than

others, a point that comes up later.  $C_1$  through  $C_6$  are normalized over the range  $[0, 1]$  where 0 is *least similar* and 1 is *most similar*.  $C_7$  and  $C_8$  are normalized over the range  $[-1, 1]$  where -1 is *least similar* and 1 is *most similar*. Each measure takes on its maximum value for identical graphs.  $C_1$  through  $C_4$  will take on their minimum value when comparing complementary graphs. Finally, algorithms computing the various  $C_i$ 's are easily coded in standard computer languages which run in  $O(n^2)$  to  $O(n^3)$  time in the number of nodes.

### Comparison of the Measures

We next compare and contrast the various measures by applying them to a common set of graphs. Consider again the graphs in Figure 1. The similarities between each pair of graphs are given in Table 1 for all eight measures. Realize that the various measures are not directly comparable (except  $C_7$  and  $C_8$ ) because each one occurs on a unique scale. Even  $C_1$  through  $C_3$ , which employ similar set-theoretic definitions, are scaled differently because of different normalizations. Therefore, only relative differences of their values are meaningful.

Table 1. The similarity between each pair of graphs  $A$ ,  $B$ ,  $C$ , in Figure 1 as measured by  $C_1$  through  $C_8$  and rank orders in parentheses from most (1) to least (3) similar.

Similarity Measure	$C(A, B)$	$C(A, C)$	$C(B, C)$
$C_1$	0.43 (2)	0.74 (1)	0.39 (3)
$C_2$	0.50 (2)	0.83 (1)	0.48 (3)
$C_3$	0.50 (3)	0.87 (1)	0.52 (2)
$C_4$	0.71 (2.5)	0.91 (1)	0.71 (2.5)
$C_5$	0.79 (1)	0.78 (2)	0.71 (3)
$C_6$	0.87 (1)	0.85 (2)	0.81 (3)
$C_7$	0.30 (2.5)	0.77 (1)	0.30 (2.5)
$C_8$	0.79 (1)	0.42 (3)	0.45 (2)

Notice first that  $C_1$  through  $C_4$  and  $C_7$  all agree that graphs  $A$  and  $C$  are most similar, whereas  $C_5$ ,  $C_6$ , and  $C_8$  show that  $A$  and  $B$  are most similar. However, the rank orders of similarity are not identical for these five measures. Notice also that  $C_5$  and  $C_6$  show an identical pattern of ranks but  $C_8$  has a different pattern. Hence, with these simple graphs there appear to be important similarities and differences in what the measures are assessing.

We turn next to a similar analysis of more complex graphs. The graphs are Pathfinder solutions to relatedness ratings of 30 course-relevant concepts given by 20 students and one instructor. For each of the 21 datasets, four classes of graphs varying in *graph density* were derived. Each graph was then compared with every other graph (210 comparisons) in its class using all eight measures. Graph-theoretic distances were used by all of the measures. The resulting similarity values for each measure were then correlated with every other measure.  $C_1$  through  $C_3$  correlated very highly across all of the graphs, as did  $C_5$



and  $C_6$ . For this reason we do not report the results of  $C_2$ ,  $C_3$ , and  $C_6$ . The resulting correlation matrices for each of the four sets of graphs are shown in Tables 2 through 5.

Notice first, as we also saw in Table 1, that  $C_1$  correlates highly with  $C_7$ , and this holds for both sparse and dense graphs. Apparently the method of comparing (set-theoretic functions or correlation) edge relation information has little effect on the resulting similarity indices.  $C_4$  correlates highly with  $C_1$  and  $C_7$  for sparser graphs, but less so for denser graphs. Recall that  $C_4$  and  $C_7$  compare edge relations with mean absolute difference and correlation, respectively. We speculate that the normalization of  $C_4$  is not as good as  $C_7$ 's or  $C_1$ 's, and this weakness becomes especially apparent with denser graphs.

Table 2. Correlations on results of  $C$ 's applied to 30 node graphs (mean graph density = 0.067).

	$C_4$	$C_5$	$C_7$	$C_8$
$C_1$	0.97	0.59	0.97	0.46
$C_4$		0.62	1.00	0.51
$C_5$			0.62	0.73
$C_7$				0.50

Table 3. Correlations on results of  $C$ 's applied to 30 node graphs (mean graph density = 0.124).

	$C_4$	$C_5$	$C_7$	$C_8$
$C_1$	0.93	0.84	0.98	0.73
$C_4$		0.81	0.97	0.78
$C_5$			0.85	0.84
$C_7$				0.75

Table 4. Correlations on results of  $C$ 's applied to 30 node graphs (mean graph density = 0.188).

	$C_4$	$C_5$	$C_7$	$C_8$
$C_1$	0.49	0.76	0.96	0.83
$C_4$		0.17	0.66	0.60
$C_5$			0.67	0.61
$C_7$				0.88

Table 5. Correlations on results of  $C$ 's applied to 30 node graphs (mean graph density = 0.278).

	$C_4$	$C_5$	$C_7$	$C_8$
$C_1$	0.66	0.83	0.94	0.86
$C_4$		0.46	0.86	0.79
$C_5$			0.74	0.69
$C_7$				0.92

$C_8$  correlates most highly with  $C_5$  for sparser, but lowest with denser graphs. Also,  $C_8$ 's correlations with  $C_1$  and  $C_7$  steadily increase with increasing graph density. A likely reason for this is because as graph density increases, the difference between the adjacency and distance vectors decreases, until for completely connected graphs the two are identical. In summary, although certain of the measures appear quite similar, overall there are some important and systematic differences. Additional analyses employing other types of graphs and graphs from other applications are needed before more general conclusions can be reached.

The above analyses simply examined relative agreement of the various measures. Evidence that one measure is actually "better" than another for assessing graph similarity would require some external index of similarity. As reported elsewhere in this volume (Goldsmith & Johnson, Chapter 17), there is evidence that the similarity between a student's PFNET of course concepts and the instructor's PFNET predicts final course grades better for similarity measured by  $C_1$  than by  $C_8$ . This finding was interpreted as support for the idea that a neighborhood comparison of graphs is more sensitive to configural (i.e., structural) information, and it is this type of information that is important in knowledge structures.

### Generalizations and Extensions of the Measures

In this section, we briefly describe some generalizations and extensions of the neighborhood measures. First, there may be occasions when we want to compare two unlabeled graphs. For example, we might want to find a *best fit* of a test graph with some target graph by permuting the nodes of the test graph. Brute force methods are useless here due to the explosive growth of the number of permutations as node size increases. Further, the problem does not seem to lend itself to a linear programming approach. Instead, we could employ a particular  $C_i$  as an optimization function in a simulated annealing process (Kirkpatrick, Gelatt, & Vecchi, 1983). Here, we would attempt to find that permutation of nodes for the test graph giving a minimum value for  $C$  when compared to the target graph. When the two graphs are not isomorphic, the annealing technique yields a *best fit* as defined by the particular  $C$  being used. Initial results based on a visual comparison of the target graph with permuted test graph are promising.

Finally, with some simple modifications we may extend the definitions of  $C$  to work for edge- and node-weighted graphs, as well as edge- and node-colored graphs. This will allow  $C$  to compare, for example, two chemical molecules by representing each with a node-colored, edge-weighted graph, whose nodes represent individual atoms and whose

edge weights represent bond valence, or perhaps, bond energy. Interest in assessing similarity of molecular structures has recently led chemists to define graph measures of similarity. For example, Herndon (1988) has developed a technique for computing a quantitative index of similarity between molecular graphs that requires first translating each graph into a string of symbols and then employing string comparison methods to determine the original graphs' structural similarity. In Appendix B we describe in more detail some modifications of a neighborhood measure of  $C$  that may be useful in comparing graph representations of molecules.

## Conclusions

We began by assuming that a measure of graph similarity should assess structural information. We attempted to argue that structure is best viewed as a property of graphs rather than as an entity. We then defined structural property as a higher-order relation and identified two distinct types of structural properties in graphs: paths and neighborhoods. Several related measures of graph similarity employing either path lengths or neighborhoods were defined and some of their properties noted. An initial analysis of these measures indicated that use of path lengths and neighborhoods to determine similarity assess different characteristics of graphs.

What lies behind these differences? Path distances describe how far away nodes are located; neighborhoods describe which nodes are linked. Path distances employ node pairs as the unit of comparison; neighborhoods use nodes. The path distance approach first converts the information contained in a path to a real number (the path length) and then compares corresponding path lengths between graphs; the neighborhood approach first compares corresponding neighborhoods with set-theoretic measures and then converts this result to a real number. Whether one approach is better than the other of course ultimately depends on its functional utility within a particular application. If the phenomena being represented by graphs is inherently described by distance information, then path distances will likely be better. However, if what is being represented is more accurately reflected by associations within neighborhoods, then the neighborhood approach should prove better.

## Appendix A

### Formal Definitions and Properties of Some Graph Similarity Measures

We identify a class of similarity functions for comparing graphs as  $\mathcal{C}$ . We denote a particular function from  $\mathcal{C}$  as  $C$  and use subscripts to distinguish different measures. A graph  $G = (V, E)$  is a finite set  $V$  of  $n$  nodes and a set  $E$  of edges where  $E$  is a subset of  $V \times V$  (Carre, 1979). We say two graphs  $G_1(V_1, E_1)$  and  $G_2(V_2, E_2)$  have a common node set if  $V_1 = V_2$ . Given two undirected and labeled graphs  $A$  and  $B$  with common node set  $V$ ,  $C(A, B)$  is the similarity between  $A$  and  $B$  as measured by  $C$ . We define a neighborhood as a region about a particular node in a graph. Let  $\delta_G(v, v')$  be the graph distance between nodes  $v$  and  $v'$ . Define  $\alpha_G(v, v')$  to be 1 if  $\delta_G(v, v') = 1$  and 0 otherwise. We denote by  $G_v$  the set of nodes  $v'$  such that  $\alpha_G(v, v') = 1$ . This set is the neighborhood about  $v$ .

We define the following similarity measures. In cases where the denominator of a summand is zero, we employ the convention that if both neighborhoods for an element are empty then the summand for that element is one, but if only one of the neighborhoods is empty, then the summand for that element is zero.

$$C_1(A, B) = \frac{1}{n} \sum_{v \in V} \frac{|A_v \cap B_v|}{|A_v \cup B_v|}$$

$$C_2(A, B) = \frac{2}{n} \sum_{v \in V} \frac{|A_v \cap B_v|}{|A_v| + |B_v|}$$

$$C_3(A, B) = \frac{1}{2n} \sum_{v \in V} |A_v \cap B_v| \left( \frac{1}{|A_v|} + \frac{1}{|B_v|} \right)$$

Notice that if we write  $C_3$  as

$$C_3(A, B) = \frac{1}{2} \left[ \frac{1}{n} \sum_{v \in V} \frac{|A_v \cap B_v|}{|A_v|} + \frac{1}{n} \sum_{v \in V} \frac{|A_v \cap B_v|}{|B_v|} \right]$$

we see that it is the average of two other similarity measures. These other measures are interesting in their own right and are reminiscent of conditional probabilities. In assessing similarity between  $A$  and  $B$  in one case, the measure is sensitive to those edges in  $A$  omitted



by  $B$ , and in the other case it is sensitive to edges in  $B$  not found in  $A$ . As the average of the two,  $C_3$  captures both cases. There may be applications where one would want to use only one of the individual measures. For example, if we assume that graph  $A$  represents a prototype of some sort, then it may be meaningful to assess the similarity of an exemplar ( $B$ ) to the prototype ( $A$ ). We would expect a different result when the similarity of  $A$  to  $B$  is computed.

Next, we define, for two nonnegative real numbers  $a$  and  $b$ ,  $a \theta b$  to be one if  $a = b$ ,  $a/b$  if  $a < b$ , and  $b/a$  if  $b < a$ . Also, let  $A \oplus B$  be the symmetric difference between sets  $A$  and  $B$ . Then,  $C_4$  is defined as follows:

$$\begin{aligned} C_4(A, B) &= 1 - \frac{1}{n^2 - n} \sum_{v \neq v'} |\alpha_A(v, v') - \alpha_B(v, v')| \\ &= \frac{1}{n^2 - n} \sum_{v \neq v'} \alpha_A(v, v') \theta \alpha_B(v, v') \\ &= 1 - \frac{1}{n} \sum_{v \in V} |A_v \oplus B_v| \end{aligned}$$

Interestingly,  $C_4$  is also one minus the average of the symmetric differences of the neighborhoods. This may be seen by noting that

$$\sum_{v \neq v'} \alpha_A(v, v') \cdot \alpha_B(v, v') = |A_v \cap B_v|$$

and

$$|\alpha_A(v, v') - \alpha_B(v, v')| = \alpha_A(v, v') + \alpha_B(v, v') - 2(\alpha_A(v, v') \cdot \alpha_B(v, v'))$$

We now define two other measures  $C_5$  and  $C_6$  which are based on path lengths. Notice that  $C_5$  is similar to  $C_4$  except that the function  $\theta$  is applied to actual graph distances for  $C_5$ , but to edge relations for  $C_4$ .

$$\begin{aligned} C_5(A, B) &= \frac{1}{n^2 - n} \sum_{v \neq v'} \delta_A(v, v') \theta \delta_B(v, v') \\ C_6(A, B) &= 1 - \frac{1}{n^2 - n} \sum_{v \neq v'} \frac{|\delta_A(v, v') - \delta_B(v, v')|}{\delta_A(v, v') + \delta_B(v, v')} \end{aligned}$$

For similarity measures  $C_1$  through  $C_6$  we may derive distance measures,  $D$ , by letting  $D = 1 - C$ . Some of the measures have interesting forms as distances.

$$D_1(A, B) = \frac{1}{n} \sum_{v \in V} \frac{|A_v \oplus B_v|}{|A_v \cup B_v|}$$

$$D_2(A, B) = \frac{1}{n} \sum_{v \in V} \frac{|A_v \oplus B_v|}{|A_v| + |B_v|}$$

$$D_4(A, B) = \frac{1}{n} \sum_{v \in V} |A_v \oplus B_v|$$

Also, from the definitions of  $C$  we may discover several interesting properties. Since

$$2 \frac{|A_v| |B_v|}{|A_v| + |B_v|} \leq \frac{1}{2} (|A_v| + |B_v|) \leq |A_v \cup B_v| \leq n-1$$

we get  $C_1 \leq C_2 \leq C_3$ ,  $C_1 \leq C_4$ , and  $C_5 \leq C_6$ . In general,  $C_4$  is incomparable to  $C_2$  and  $C_3$ . Similar relations hold for the distance measures as well. Also, it is possible (but tedious) to show that the distance indices for all of the similarity measures, except  $C_3$ , are metrics on the space of subsets of  $V \times V$ .

We can think of graphs with common node set  $V$  as subsets of  $V \times V$ . As such, the set of all such graphs is a Boolean ring whose multiplicative identity,  $I$ , is the completely connected graph and whose zero is the empty graph. In this ring, multiplication is given by intersection and addition is given by symmetric difference. If we choose  $D_1$  for a metric on

this space, we may define a norm (which we call the  $C_1$  norm) on a graph  $G$  by  $\|G\| = 1 - D_1(G, I) = C_1(G, I)$ . With this definition, the following properties hold for all graphs  $A$ ,  $B$ , and  $C$  over node set  $V$ :

$$\|A \cup B\| + \|A \cap B\| = \|A\| + \|B\|$$

$$1 - \|\bar{A}\| = \|A\|$$

$$\|A \oplus B\| \leq \|A\| + \|B\|$$

$$\|A \cap B\| \leq \|A\|$$

$$\|A \oplus B\| + \|A \oplus C\| \leq \|B \oplus C\|$$

where  $\bar{A}$  is the *complement* of  $A$  with respect to  $I$ . This last inequality allows us to define a metric  $D(A, B) = \|A \oplus B\|$ , which turns out surprisingly to be  $D_4$ .

## Appendix B

### *An Extension of Neighborhood Similarity Measures*

The weight of an edge  $(v, v')$  in a graph  $G$  will be denoted by  $w_G(v, v')$  and its color by  $x_G(v, v')$ . The weight of a node  $v$  in  $G$  will be denoted  $w_G(v)$  and its color by  $x_G(v)$ . Next we define a Kronecker delta function on the set of colors in a graph by  $\delta(c, c') = 1$ , if  $c$  is the same color as  $c'$  and is 0 otherwise. All of our measures then take the form

$$\frac{1}{v_1} \sum_{v \in V} (w_A(v) \theta w_B(v)) (\delta(x_A(v), x_B(v)))$$

$$\frac{1}{v_2(v)} \sum_{v \neq v'} (w_A(v, v') \theta w_B(v, v')) (\delta(x_A(v, v'), x_B(v, v')))$$

where  $v_1$  and  $v_2$  are normalizing functions. For example

$$\frac{1}{n} \sum_{v \in V} \delta(x_A(v), x_B(v)) \left[ \frac{1}{|A_v \cup B_v|} \sum_{v \neq v'} (w_A(v, v') \theta w_B(v, v')) \right]$$

is a measure of similarity between two node-colored, edge-weighted graphs. Notice that this measure reduces to  $C_1$  for graphs that are not colored or weighted. With careful normalization, a number of potentially useful measures are definable.