

# Relating Fuzzy Set Similarity Measures

Valerie Cross<sup>(✉)</sup>

Computer Science and Software Engineering,  
Miami University, Oxford, OH 45056, USA  
crossv@miamioh.edu

**Abstract.** Measuring similarity is an important task in many domains such as psychology, taxonomy, information retrieval, image processing, bioinformatics, and so on. The diversity of domains has led to many different definitions of and methods for determining similarity. Even within fuzzy set theory, how to measure similarity between fuzzy sets presents a wide variety of approaches depending on what characteristic of a fuzzy set is emphasized, for example, set-based, logic-based or geometric-based views of a fuzzy set. First similarity is examined from a psychological viewpoint, and how that perspective might be applicable to fuzzy set similarity measures is explored. Then two fuzzy set similarity measures, one set-based and the other geometric-based, are reviewed, and a comparison is made between the two.

**Keywords:** Fuzzy set similarity · Set-based similarity · Geometric-based similarity · Dissemblance index

## 1 Introduction

Comparing two concepts or objects is a necessary process in many domains such as biology, psychology, taxonomy, statistics and artificial intelligence. This comparison operation attempts to determine a relationship between the two concepts. One such type of relationship that is frequently determined is their similarity. Because of the diversity of domains, the general meaning of similarity is ambiguous with many different definitions and approaches to measuring similarity. As presented in psychological theory [1], a warning on assessing similarity is given, “Like most powerful and widespread ideas, it [similarity] is not amendable to a ready and precise definition; indeed, this very resistance to definition probably goes far to explain its usefulness as a supposed explanatory principle. Ideas that are imprecise are also dangerously versatile when it comes to accounting for the complexities of human behavior.”

Even within one domain such as fuzzy set theory, a wide variety of methods exist for assessing similarity [2], many of which are extensions of similarity measures that are well-known in their respective research domains. The more recent research area of ontological knowledge representation for the Semantic Web has also had a proliferation of semantic similarity measures for various tasks such as ontology alignment, information extraction, and semantic annotation. The objective of a semantic similarity measure, also referred to as an ontological similarity measure, is to calculate the degree to which one concept is similar to another concept within the context of an ontology.

Although several major categories of semantic similarity measures exist such as path-based, information content-based and feature-based, those measures using information content have been the emphasis of much study and evaluation especially in the bioinformatics and biomedical domains. In [3] many of these semantic similarity measures are shown as related to fuzzy set similarity measures if a concept is represented as a fuzzy set consisting of itself and all its ancestor concepts and the membership degrees are based on the information content of each concept within the context of the ontology.

The focus of this paper is that of similarity in fuzzy set theory. This paper examines some “respects for similarity” [4] from the domain of psychological theory and their general applicability to the measurement of similarity in fuzzy set theory. “Respects for similarity” refers to the ways in which two things can be similar. The term frame of reference [5] is also used for respects for similarity. The comparison process has intrinsic factors that determine the respects. As pointed out in [4], asking “How similar are X and Y?” can be viewed as asking a slightly different question, “How are X and Y similar? The process for fixing the respects is a crucial facet of similarity comparisons.

Correspondingly, similarity in other domains is investigated to better understand how fuzzy set similarity measures have been extended from these domains. Two specific similarity measures, one used very early in taxonomy and the other used in calculating distances between intervals on the real line are reviewed and their fuzzy set extensions analyzed. These two similarity measures are compared to determine any relationships between them. To begin, Sect. 2 looks at similarity as an empirical and theoretical psychological construct and attempts to elicit correspondences to fuzzy set similarity. These correspondences might suggest other views and uses for fuzzy set similarity measures. Just like different characteristics of a concept in the context of an ontology are considered important to constructing a semantic similarity measure, a variety of characteristics of a fuzzy set are considered in the construction of a fuzzy set similarity measure. Section 3 presents the taxonomic related fuzzy set similarity measures and its relation to Tversky’s psychological model of similarity. The fuzzy extension of the distance between real number intervals to the similarity between fuzzy set intervals is described in Sect. 4. Section 5 compares and contrasts these two fuzzy set similarity measures and establishes a relationship between them. A summary and plans for future research are provided in Sect. 6.

## 2 Respects for Similarity

In [4] the researchers examine similarity as an explanatory construct in psychological theory where humans are comparing two objects or things. The things being compared ranged from two simple linguistic terms to two visual forms. Their experiments indicate that similarity is highly flexible and in some ways troublingly flexible. Their experimental observations, however, are used to argue that the flexibility is reasonable as long as systematic changes in the process of similarity assessment can be established.

The research of Tversky [5] has played a major role in shaping the understanding of similarity in psychological research. Tversky’s research informs the research in [6] where it is noted that “the relative weighting of a feature (as well as the relative

importance of common and distinctive features) varies with the stimulus context and the task; so that there is no unique answer to the questions of how similar is one object to another.” This quote on similarity again emphasizes its “resistance to definition” and that its assessment method is “dangerously versatile.”

The argument [4] is made that instead of viewing similarity assessment as constrained by the perceptual process, similarity assessment is flexible and the comparison process itself methodically sets the respects. Assessing similarity is assumed to be based on matching and mismatching of properties. Things are similar to the degree they share properties and dissimilar to the degree that properties apply to one but not the other. The issue is that two things share a subjective number of properties and likewise they differ in a subjective number of properties. Before similarity can be computed, a prior process must occur that determines what properties are to be used in the similarity computation.

Others [7] argue that the respects for determining how two things are judged as similar are set not by the comparison process but by the goals motivating the comparison process. This view is the result on research to determine the requirements for a similarity measure for use in the automatic generation of textual comparisons. Comparison between objects is categorized into six different types. For example, a clarificatory comparison is a domain-based comparison with the goal of distinguishing one object from another object that is highly similar to it. Domain-based comparisons are used to establish explicit relationships between an object and other objects existing in the same domain.

Regardless of how and when respects are established, most agree that similarity assessment cannot be performed without them. The problem still remains as to the process of selecting the respects as so aptly described by Tversky [5], “When faced with the a particular task (e.g. identification or similarity assessment) we extract and compile from our data base a limited set of relevant features on the basis of which we perform the required task.”

Another important issue discussed in [4] is the effect of context in similarity assessment. Setting the context for comparison contributes to the selection of the respects to which similarity is being assessed. Two objects may be judged less similar when no explicit context is given than when one is given because the context tends to make salient the context-relevant properties to be used in the similarity assessment. The similarity of the two objects is increased based on the degree to which the two objects share values for these now salient properties.

The extension effect of context is also important in similarity assessment. When in one context, properties that are shared by all objects are not useful in similarity assessment; however, if the context is extended or broadened to include objects not sharing these properties, then these properties become more salient. In the extended context, two objects sharing those properties are perceived as more similar than they were in the original context. To summarize, depending on the context, two objects may vary in their similarity, but this variability becomes systematic when incorporated into the specification of a similarity comparison.

Analogy also plays a role in similarity assessment. Instead of focusing on similarity in values for simple properties of objects, it looks for relational or structural similarities. An example given in [4] is “an atom is like the solar system” where the analogy relies on relations such as “revolves around” and not property values such as “hot” or

“yellow”. The research in [4] argues the importance of incorporating relational structure since relational structures can significantly affect the process of determining the correspondences between objects when assessing similarity.

Similarity assessment involves comparing two objects but this comparison process may be directional. The example given in [4] is very informative: “surgeons are like butchers” as compared to “butchers are like surgeons”. The former is critical of surgeons whereas the latter is favorable of butchers. In the contrast model of similarity [5], the less salient or less prominent object is compared to the more salient or prominent object as evidence by the results of human experiments where the less salient object is considered more similar to the more salient object than vice versa.

The direction of comparison also affects the properties selected for assessing similarity as shown in their experimental results [4]. The selected properties may be more closely related to the base object to which a comparison is being made. The common properties used in comparing two objects may vary as a function of the direction of a comparison and the bias is to select properties more strongly associated with the base object. To summarize, similarity is more than identity since similarity comparisons may encompass properties of one object becoming the candidate properties of the other in performing the similarity assessment.

Since similarity assessment usually involves multiple properties, research in cognitive psychology has concentrated on how multiple pieces of information are integrated into a single assessment of similarity. Similarity assessment is affected by both the selection of the applicable properties and the constraints that the integration method places on the process. As part of the integration method, weighting may be involved that favors certain properties over others. In [4] experiments have shown that this weighting procedure is not independent of the outcome of the comparison process.

One last interesting aspect brought out in [4] is the notion of experience and learning affecting the process of similarity assessment. Children, for example, judge similarity in a more holistic manner and are less likely to analyze individual components, but as they mature, they base their similarity judgements more on abstract, relational, and less on superficial properties.

To summarize the research in [4] for the domain of psychology, similarity assessment is dynamic and highly variable but connected to the details of the comparison process. The details that are focused on in their research are the fixing of the properties or respects to which objects are similar, the context, the direction of the comparison, the kind of properties whether simple attributes or relational structures, the integration of multiple information and the weighting of this information in the process and human experience. Many of these details of the comparison process in similarity assessment can be found in fuzzy set similarity measurements.

In the following two sections, a set based fuzzy similarity measure from taxonomy related its related measures and a then a geometric based fuzzy set similarity measure from distance between real line intervals are described. Their details are examined from the viewpoint of similarity assessment in the domain of psychology.

### 3 Set-Based Similarity

One of the early set based similarity measures for crisp sets is the Jaccard index [8], which was used in taxonomic classification. A specimen is represented by a set of attributes describing it. Two specimens are judged to be similar based on the similarity between their set of attributes. In taxonomy, the Jaccard index has also been referred to as the “coefficient of similarity” [9] and in psychology, it is the unparameterized ratio model of similarity [5]. Its formula where  $X$  and  $Y$  are sets is expressed as

$$S_{jaccard}(X, Y) = \frac{f(X \cap Y)}{f(X \cup Y)}. \quad (1)$$

The function  $f$  is an additive function and is typically the cardinality of the set. The Jaccard index is easily extended when  $X$  and  $Y$  are fuzzy sets by using fuzzy set operators to perform the intersection and the union on the two fuzzy sets and the function  $f$  is fuzzy set cardinality, which is simply the sum of the membership degrees for all elements in the fuzzy set. A fuzzy Jaccard dissimilarity measure can be derived by subtracting the Jaccard similarity from 1, i.e.,  $D_J = 1 - S_{jaccard}(X, Y)$ .

From the psychological analysis of similarity assessment, the fuzzy sets  $X$  and  $Y$  are being compared based not only on the elements making up each set but the degree of membership of each element in the set. The selection of properties in this similarity measure is natural; that is, all elements in the support of a fuzzy set describe it. The selected properties for the comparison process, therefore, include both the support of  $X$  and the support of  $Y$ .

The correspondence or alignment between the properties of the two fuzzy sets is automatic since each element in the fuzzy set is considered a property and the constraint on a fuzzy intersection is an exact match on each element in the intersection. The weighting, however, for a property (element) in this comparison process is its degree of membership or agreement with the fuzzy concept being represented by the fuzzy set. In addition to the required exact match on the aligned property values is the constraint on the integration between their two membership degrees using a fuzzy set intersection operator, which is typically *min*. The result is that multiple pieces of information exist since there are multiple elements (properties) and further integration, referred to as aggregation in fuzzy set theory, must occur to assess the overall similarity of the two fuzzy sets. With the Jaccard index, the aggregation operator is summation, that is, the cardinality of the fuzzy set intersection.

The numerator of the Jaccard index provides an assessment of the agreement of properties between the two fuzzy sets but does not take into consideration, properties in one fuzzy set that are not contained in the other fuzzy set and vice versa. The denominator, which is the union of the fuzzy sets, typically using the *max* operator, does consider this and thus normalizes the overall similarity assessment in  $[0, 1]$ .

Psychological similarity considers direction of comparison as a critical aspect in the process. The Jaccard index does not account for presupposing a direction for the comparison. An inclusion index, however, can and is a version of the parameterized ratio model of similarity [5], which is given as

$$S_{Tversky-ratio}(X, Y) = \frac{f(X \cap Y)}{f(X \cap Y) + \alpha f(X - Y) + \beta f(Y - X)}. \quad (2)$$

where  $(X - Y)$  is set difference operator. Setting  $\alpha = 1$ ,  $\beta = 1$  produces the Jaccard index. Setting  $\alpha = 1$ ,  $\beta = 0$  produces the degree of inclusion for  $X$ , that is, the proportion of  $X$  overlapping with  $Y$ , given as

$$S_{inclusion}(X, Y) = \frac{f(X \cap Y)}{f(X)}. \quad (3)$$

In the parameterized ratio model, the value  $f(X)$  for object  $x$  is considered a measure of the overall salience of that object. In psychology, the factors adding to an object's salience include "intensity, frequency, familiarity, good form, and informational content" [5]. Although the cardinality of a fuzzy set is a very simple way to measure the "salience" of a fuzzy set, i.e., the larger the cardinality, the less salience, other ways might be more useful depending on the application. Both fuzzy entropy [10] and a function of the distance of a set to its complement [11] have been used as fuzziness measures. One could consider that a fuzzy set is more salient than another fuzzy set if it has less fuzziness.

For fuzzy rule-based reasoning systems, salience of the two fuzzy sets being compared is not relevant. One approach that is used is to set the comparison direction from the observation fuzzy set as compared to the rule antecedent fuzzy set, which becomes the base for comparison to. The objective is to determine how certain is it that the observation satisfies the antecedent. The more the observation is included within the antecedent, the more certain that the antecedent is satisfied. If the observation fuzzy set is a subset of the antecedent fuzzy set, the inclusion measure produces a one. Not every fuzzy rule base system, however, uses an inclusion measure to assess agreement between the rule antecedent and the observation fuzzy sets.

In fuzzy applications that are to mimic human directional comparison judgments, the use of the more salient fuzzy set as the base for comparison might be more appropriate. Here the properties of the more salient fuzzy set  $S$  become the selected properties for the comparison process, and those properties in the less salient fuzzy  $L$  set that are not in  $S$  are simply ignored. Here similarity is more than an identify as described in [4]. In this use of similarity, the inclusion index measures the proportion of the properties of  $S$  found in  $L$  to all the properties of  $S$  and is given as

$$S_{inclusion}(S, L) = \frac{f(S \cap L)}{f(S)}. \quad (4)$$

Image processing applications [12] using fuzzy set theory tested two different versions of the inclusion index with other fuzzy set similarity measures in a shape classification experiment. The denominator of the inclusion index is replaced by either  $\min(f(X), f(Y))$  or  $\max(f(X), f(Y))$ . In the experimental results, the error rate for the  $\max$  version of the denominator were much smaller than that of the  $\min$  version. The comparison direction and how to choose that direction makes a difference and is application dependent.

For an example of the extension effect discussed with psychological research on similarity, consider another fuzzy set similarity measure used in [12] which follows the formula of the Jaccard index but instead measures the similarity between the complements of the fuzzy sets, i.e.  $X'$  and  $Y'$ . The more the complements of the two sets are similar, then the more the two fuzzy sets are similar. With this approach, if the context or the universe of discourse for the two fuzzy sets is extended, i.e., its size increased, then the Jaccard index for the two fuzzy sets would not be affected by the extension since the properties considered salient would still be those in the union of the two fuzzy sets. The Jaccard index as measured using the complements of the two fuzzy sets, however, would be affected and would be greater in the extended context than in the original context. Intuitively, in the extended context the complements of the fuzzy sets share more properties.

## 4 Geometric-Based Similarity

Geometric based similarity relies on the dissemblance index, which provides a normalized distance between two real intervals. If  $V = [v_1, v_2]$  and  $W = [w_1, w_2]$ , the dissemblance index is given as

$$D(V, W) = \frac{(|v_1 - v_2| + |w_1 - w_2|)}{2 * (\beta_2 - \beta_1)}. \quad (5)$$

where  $[\beta_1, \beta_2]$  is the smallest interval that contains both the  $V$  and  $W$  intervals. The factor  $2 * (\beta_2 - \beta_1)$  is necessary to produce a normalized dissemblance in  $[0, 1]$ .

The dissemblance index consists of two components, the left and right distance between the two intervals and may be generalized to fuzzy intervals. A pair of boundary functions  $L_N$  and  $R_N$  and parameters  $(r_1, r_2, \lambda, \rho)$  define a fuzzy interval. The core of  $N$  is  $[r_1, r_2]$  and  $\lambda$  and  $\rho$  are parameters of the boundary functions  $L_N$  and  $R_N$  such that the support of  $N$  is in the interval  $[r_1 - \lambda, r_2 + \rho]$ . If  $L_N$  and  $R_N$  are positively and negatively sloping linear functions, respectively, then  $N$  is represented by a trapezoidal fuzzy set membership function. Figure 1 illustrates two fuzzy trapezoidal fuzzy sets  $X$  and  $Y$  and labels for left and right boundaries.

To calculate the fuzzy dissemblance index between two fuzzy intervals  $X$  and  $Y$ , the formula uses integration over the  $\alpha$ -cuts of the fuzzy intervals as

$$fD(X, Y) = \frac{1}{2(\beta_2 - \beta_1)} \int_0^1 (|L_X(\alpha) - L_Y(\alpha)| + |R_X(\alpha) - R_Y(\alpha)|) d\alpha. \quad (6)$$

where  $[\beta_1, \beta_2]$  is the smallest interval that contains both the support of the  $X$  and  $Y$  fuzzy intervals.  $fD$  calculates a fuzzy dissimilarity measure between two fuzzy intervals based on a normalized distance and can be converted into a fuzzy similarity measure as  $S_{fD}(X, Y) = 1 - fD(X, Y)$ .

With the fuzzy similarity measure  $S_{fD}$ , also referred to as a geometric fuzzy similarity [2], the alignment between properties is not based on identical property values as for the Jaccard fuzzy similarity measure but on identical  $\alpha$  values. The comparison is measured between the property values at the identical  $\alpha$  values for the left and the right

components of the fuzzy interval. This geometric similarity differs from the Jaccard fuzzy similarity measure since the comparison is done on the  $\alpha$  values and resolved using a fuzzy set intersection operator. Correspondingly, both have a normalizing factor that includes the support of both  $X$  and  $Y$ .

Some similarity research have been proposed to approximate  $fD$  to avoid the computationally expensive integration over  $\alpha$ , the value [13]. These approximations use only the distance obtained from a single  $\alpha$ -cut, for example, only the distance between the core intervals of the fuzzy sets. This approximation does not incorporate information about the proximity of the support intervals. Thus, the approximation result may be much smaller than  $fD$ . A summarization technique was introduced in [14]. First, the distance between the support intervals is determined as

$$fD_0(X, Y) = \frac{1}{2(\beta_2 - \beta_1)} (|L_X(0) - L_Y(0)| + |R_X(0) - R_Y(0)|) \quad (7)$$

and similarly for the core intervals,  $fD_1$ . The summarized distance is the average core and support distances given as

$$fD_{\textcircled{a}} = \frac{fD_0 + fD_1}{2}. \quad (8)$$

For trapezoidal fuzzy sets in which  $L_X$  does not intersect  $L_Y$  and  $R_X$  does not intersect  $R_Y$ , this summarization technique produces equivalent results as  $fD$ . When  $L_X$  does intersect  $L_Y$  at  $a_L$  the left distance for the support interval must be factored by  $a_L$  and the left distance for the core interval must be factored by  $(1 - a_L)$  and similarly if  $R_X$  does intersect  $R_Y$  at  $a_R$ . This factor represents the height of the triangle created at the intersections.

The geometric fuzzy similarity measure  $S_{\text{diss}}(X, Y) = 1 - fD(X, Y)$  and its use in fuzzy reasoning is presented in [14] since using this distance based measure allows a fuzzy conclusion to be determined using the left and right distances between the fuzzy rule antecedent and the fuzzy observation even when there is no overlap between the two. The details of this fuzzy reasoning approach are not examined here but instead a relationship between the fuzzy Jaccard similarity and the fuzzy geometric similarity measures are explored.

## 5 Relating Set and Geometric Similarity

When extending the similarity measures of psychology and taxonomy to similarity measures for fuzzy sets, it is natural to see how features of objects are replaced by elements of the fuzzy sets, crisp set cardinality replaced with fuzzy set cardinality and set operators replaced with fuzzy set operators. However, not all equalities using crisp set operators are true for all possible fuzzy set operators. For example, when  $X$  and  $Y$  are crisp sets and not disjoint,  $f(X \cup Y) = f(X) + f(Y) - f(X \cap Y)$ . This equality is true for fuzzy sets only when members of Frank's family of dual t-norms and t-conorms [15] are selected for the union and intersection operators. A more methodical method of



creating a framework for fuzzy set similarity measures is based on developing a set of properties that they should satisfy. In order to develop the relationship between the Jaccard and geometric fuzzy similarity measures the theoretical foundation for the fuzzy Jaccard similarity measure is first presented [16].

One of the properties established for a fuzzy set similarity measure between  $X$  and  $Y$  is that  $S(X, Y) = 1$  if and only if the symmetric difference between the two,  $(X\Delta Y)$  is the empty set. Another property is if  $X$  and  $Y$  have disjoint sets, then  $S(X, Y) = 0$ . To meet these conditions a fuzzy set similarity measure is derived using relative cardinality on the negation of the symmetric difference between  $X$  and  $Y$ ,  $g((X\Delta Y)')$  where  $g$  is relative cardinality and

$$X\Delta Y = (X \cup Y) \cap (X' \cap Y') = (X \cap Y') \cap (X' \cap Y)$$

Here is another example of an equality being true for crisp sets but only true for fuzzy sets when minimum is used for intersection and maximum is used for union.

The fuzzy similarity measure should be in the interval  $[0, 1]$  so the range for  $g((X\Delta Y)')$  must be found to produce a normalized value. The maximum value for  $g(X\Delta Y)$  occurs when the two fuzzy sets are disjoint, which is  $g(X \cup Y)$ . The minimum value for  $g((X\Delta Y)')$ , therefore, occurs for  $g((X \cup Y)')$ . The range for  $g((X\Delta Y)')$  is  $[(g((X \cup Y)'), 1]$ . The fuzzy similarity measure can be derived as

$$S(X, Y) = \frac{g((X\Delta Y)') - g((X \cup Y)')}{1 - g((X \cup Y)')}$$

This equation can be rewritten as

$$S(X, Y) = \frac{g(X \cup Y) - g(X\Delta Y)}{g(X \cup Y)} = 1 - \frac{g(X\Delta Y)}{g(X \cup Y)}$$

since  $g(X) = 1 - g(X')$  for relative cardinality. From the above equation, the fuzzy similarity measure produces a 0 if and only if  $g(X\Delta Y) = g(X \cup Y)$ , that is the fuzzy sets  $X$  and  $Y$  are disjoint. The fuzzy set similarity measure produces a 1 if and only if the symmetric difference produces the empty set. When  $X$  and  $Y$  are crisp and  $X = Y$ , all the symmetric difference operators produce an empty set. When  $X$  and  $Y$  are fuzzy sets, however, the only symmetric difference operator to produce an empty set when  $X = Y$  is derived using  $(X \cap Y') \cup (X' \cap Y)$  with bold intersection,  $\max(0, u_X(v) + u_Y(v) - 1)$  and bold union,  $\min(1, u_X(v) + u_Y(v))$ .

Using this symmetric difference operator and replacing relative cardinality with fuzzy set cardinality since the cardinality of the universe of discourse may be cancelled out in the numerator and denominator

$$\frac{g(X\Delta Y)}{g(X \cup Y)} = \frac{\sum_v \min(1, (\max(0, u_X(v) - u_Y(v)) + \max(0, u_Y(v) - u_X(v))))}{|X \cup Y|}$$

Since the differences in the membership degrees cannot be larger than 1 the min operation can be removed to produce

$$\frac{g(X\Delta Y)}{g(X\cup Y)} = \frac{\sum_v \max(0, u_X(v) - u_Y(v)) + \max(0, u_Y(v) - u_X(v))}{|X\cup Y|}$$

Since either the membership of  $v$  in  $X$  is greater than or equal to its membership in  $Y$ ,

$$\frac{g(X\Delta Y)}{g(X\cup Y)} = \frac{\sum u_X(v) - \min(u_X(v), u_Y(v)) + u_Y(v) - \min(u_X(v), u_Y(v))}{|X\cup Y|}$$

Now rewriting by distributing the summation operator over each component in the summation and using set intersection  $\cap$  for minimum produces

$$\frac{g(X\Delta Y)}{g(X\cup Y)} = \frac{(|X| + |Y| - 2|X\cap Y|)}{|X\cup Y|} = \frac{(|X\cup Y| - |X\cap Y|)}{|X\cup Y|} = 1 - \frac{|X\cap Y|}{|X\cup Y|}$$

since for the maximum and minimum operators,  $|X| + |Y| = |X\cup Y| + |X\cap Y|$ , therefore, resulting in

$$S(X, Y) = 1 - \left(1 - \frac{g(X\cap Y)}{g(X\cup Y)}\right) = \frac{g(X\cap Y)}{g(X\cup Y)}$$

which is the original proposed ‘‘similarity of coefficient’’ used in taxonomic classification. If the fuzzy Jaccard similarity measure is converted to a dissimilarity measure by subtracting from 1, then

$$D_J(X, Y) = \frac{g(X\Delta Y)}{g(X\cup Y)}$$

which also incorporates the symmetric difference.

There is a strong relationship between the fuzzy dissemblance measure and the Jaccard dissimilarity measure. The fuzzy distances calculated for the left and right components of dissemblance dissimilarity measure when added together include the symmetric difference between  $X$  and  $Y$ .

To establish the relationship between the two fuzzy dissimilarity measures, first consider two cases, (1) the fuzzy sets  $X$  and  $Y$  do not intersect and (2) the fuzzy sets  $X$  and  $Y$  do intersect. Case 1 is easier since when they do not intersect,  $D_J(X, Y) = 1$  because the symmetric difference produces the same as the union of the two sets. Thus  $fD(X, Y) \leq D_J(X, Y)$ . Case 2 has two subcases: (1) the cores of the fuzzy sets intersect and (2) the cores of the fuzzy sets do not intersect.

Subcase 1 is easier since with overlap in the cores of the fuzzy sets, the dissemblance dissimilarity only includes the symmetric difference as in  $D_J$ . Thus,  $fD(X, Y) \leq D_J(X, Y)$  since both have the same numerator  $g(X\Delta Y)$  but the normalization factor in the denominator for  $fD$  is  $2 * (\beta_2 - \beta_1)$  which is always greater than or equal to  $|X\cup Y|$ .

Subcase 2 is most difficult since when  $R_X(\alpha)$  intersect  $L_Y(\alpha)$  at  $\alpha_I$  there is a distance between  $[R_X(I), L_Y(I)]$ . Figure 1 illustrates this. The dissemblance dissimilarity measure in addition to the symmetric difference, includes this distance as twice the area of the top triangle  $T$  with base of  $(L_Y(I) - R_X(I))$  and height of  $(1 - \alpha_I)$  value since  $\alpha_I$  is the point of intersection. The  $(1 - \alpha_I)$  value represents the height of triangle since the triangle is formed above the  $\alpha_I$  intersection point. This triangle area is included twice because both the distance between the left boundary functions of  $X$  and  $Y$  and between the right boundary functions are include this triangle area. Rewriting the fuzzy dissemblance measure and using symbol  $T$  in the equation,

$$fD(X, Y) = \frac{g(X\Delta Y) + 2 * T}{2 * (\beta_2 - \beta_1)}$$

To analyze this, the starting point is when  $R_X(\alpha)$  intersect  $L_Y(\alpha)$  at  $\alpha_I = 0$ . Since  $X$  and  $Y$  are disjoint,  $fD(X, Y) \leq D_f(X, Y)$ , the case 1 scenario. When  $R_X(\alpha)$  intersects  $L_Y(\alpha)$  at  $\alpha_I$ , two triangles are formed the top triangle  $T$  and the bottom triangle  $B$ . The area of  $B$  is  $|X \cap Y|$ . The area of  $T$  is at a maximum when  $\alpha_I = 0$  since its height, therefore, would be 1. However, this is case 1 and  $fD(X, Y) \leq D_f(X, Y)$  for this case. As  $\alpha_I$  increases, the area of  $T$  shrinks. As the area of the intersection grows, the corresponding area of  $T$  shrinks. In comparing to  $D_f(X, Y)$ , even at the maximum area for  $T$ , the fuzzy dissemblance similarity is still smaller than the fuzzy Jaccard dissimilarity measure. Twice the area of triangle  $T$  cannot produce a large enough value to cause  $fD(X, Y)$  to surpass  $D_f(X, Y)$ .

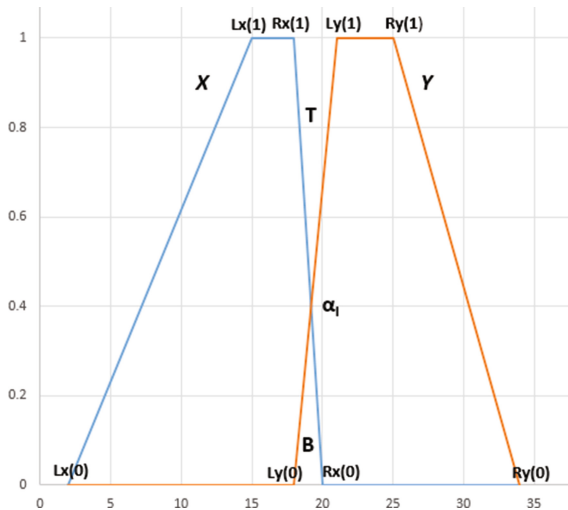


Fig. 1. Trapezoidal fuzzy sets X and Y with B intersection area and T dissemblance overlap.

## 6 Conclusions

Measuring similarity is an important task in many domains such as psychology and taxonomy. Typically, similarity assessment is performed on crisp sets. In fuzzy reasoning, it is performed on fuzzy sets. Natural extensions to crisp set similarity have been made for fuzzy set similarity. It is important to understand how the research in psychology regarding important factors affecting human similarity judgements might apply to fuzzy set similarity assessments. Examples of such important issue in psychological research include selecting the “respects of similarity”, determining context of comparison, understanding comparison direction, considering the difference in properties whether simple attributes or relational structures, integrating multiple information and the weighting of this information, and taking into account human experience. Many of these issues of the comparison process in human similarity assessment can be found in fuzzy set similarity measurements.

Two different measurements of similarity and correspondingly dissimilarity have been reviewed and some of their theoretical foundations presented. The fuzzy Jaccard similarity measure based on taxonomy’s “coefficient of similarity” and Tversky’s non-parameterized ratio model of similarity and the fuzzy dissemblance dissimilarity measure have been compared. In this paper, a fuzzy similarity measure can be used to produce a fuzzy dissimilarity measure by subtracting it from 1.

The relationship between the fuzzy Jaccard dissimilarity measure  $D_f(X, Y)$  and the fuzzy dissemblance measure  $fD(X, Y)$  is the fuzzy dissemblance measure always produces a value less than or equal to the fuzzy Jaccard dissimilarity measure. Correspondingly the fuzzy Jaccard similarity measure  $S_f(X, Y)$  always produces a value less than or equal to the fuzzy dissemblance similarity measure  $S_{diss}(X, Y)$ . Both  $D_f(X, Y)$  and  $fD(X, Y)$  use the symmetric difference between X and Y.

## References

1. Gregson, R.M.: Psychometrics of Similarity. Academic Press, New York (1975)
2. Cross, V., Sudkamp, T.: Similarity and Compatibility in Fuzzy Set Theory: Assessment and Applications. Physica-Verlag, New York (2002)
3. Cross, V.: Ontological similarity. In: Popescu, M., Xu, D. (eds.) Data Mining in Biomedicine Using Ontologies, pp. 23–43. Artech House, Norwood (2009)
4. Medin, D.L., Goldstone, R.L., Gentner, D.: Respects for similarity. Psychol. Rev. **100**, 254–278 (1993)
5. Tversky, A.: Features of similarity. Psychol. Rev. **84**, 327–352 (1977)
6. Murphy, G.L., Medin, D.L.: The role of theories in conceptual coherence. Psychol. Rev. **92**, 289–316 (1985)
7. Milosavljevic, M.: Comparison purpose and the respects for similarity. In: Slezak, P. (ed.) Proceedings of the Joint International Conference on Cognitive Science with the Australasian Society for Cognitive Science. University of New South Wales, Sydney (2003)
8. Jaccard, P.: Nouvelles recherches sur la distribution florale. Bull. Soc. Vaud Sci. Nat. **44**, 223 (1908)
9. Sneath, P.H.A., Sokal, R.R.: Numerical Taxonomy. W. H. Freeman and Company, San Francisco (1973)

10. De Luca, A., Termini, S.: A definition of non-probabilistic entropy in the setting of fuzzy set theory. *Inf. Control* **20**, 301–312 (1972)
11. Yager, R.R.: On the measure of fuzziness and negation, Part I: membership in the unit interval. *Int. J. Gen. Syst.* **5**, 221–229 (1979)
12. der Weken, D.V., Nachtegael, M., Kerre, E.E.: Using similarity measures and homogeneity for the comparison of images. *Image Vis. Comput.* **22**, 695–702 (2004)
13. Zwick, R., Carlstein, E., Budescu, D.: Measures of similarity among fuzzy concepts: a comparative analysis. *Int. J. Approx. Reason.* **1**, 221–242 (1987)
14. Cross, V., Sudkamp, T.: Geometric compatibility modification. *Fuzzy Sets Syst.* **84**, 283–299 (1996)
15. Frank, M.J.: On the simultaneous associativity of  $F(x,y)$  and  $x + y - F(x,y)$ . *Aequationes Math.* **19**, 194–226 (1979)
16. Dubois, D., Prade, H.: A unifying view of comparison indices in a fuzzy set-theoretic framework. In: Yager, R. (ed.) *Fuzzy Sets and Possibility Theory: Recent Developments*, pp. 3–13. Pergamon Press, New York (1982)