

How to Normalize Cooccurrence Data? An Analysis of Some Well-Known Similarity Measures

Nees Jan van Eck and Ludo Waltman

*Econometric Institute, Erasmus School of Economics, Erasmus University Rotterdam, The Netherlands
and Centre for Science and Technology Studies, Leiden University, The Netherlands.*

E-mail: {nvaneck, lwaltman}@ese.eur.nl.

In scientometric research, the use of cooccurrence data is very common. In many cases, a similarity measure is employed to normalize the data. However, there is no consensus among researchers on which similarity measure is most appropriate for normalization purposes. In this article, we theoretically analyze the properties of similarity measures for cooccurrence data, focusing in particular on four well-known measures: the association strength, the cosine, the inclusion index, and the Jaccard index. We also study the behavior of these measures empirically. Our analysis reveals that there exist two fundamentally different types of similarity measures, namely, set-theoretic measures and probabilistic measures. The association strength is a probabilistic measure, while the cosine, the inclusion index, and the Jaccard index are set-theoretic measures. Both our theoretical and our empirical results indicate that cooccurrence data can best be normalized using a probabilistic measure. This provides strong support for the use of the association strength in scientometric research.

Introduction

The use of cooccurrence data is very common in scientometric research. Cooccurrence data can be used for a multitude of purposes. Cocitation data, for example, can be used to study relations among authors or journals, coauthorship data can be used to study scientific cooperation, and data on cooccurrences of words can be used to construct so-called co-word maps, which are maps that provide a visual representation of the structure of a scientific field. Usually, when cooccurrence data is used, a transformation is first applied to the data. The aim of such a transformation is to derive similarities from the data or, more specifically, to normalize the data. For example, when researchers study relations among authors based on cocitation data, they typically derive similarities

from the data and then analyze these similarities using multivariate analysis techniques such as multidimensional scaling and hierarchical clustering (e.g., McCain, 1990; White & Griffith, 1981; White & McCain, 1998). Likewise, when researchers use coauthorship data to study scientific cooperation, they typically apply a normalization to the data and then base their analysis on the normalized data (e.g., Glänzel, 2001; Luukkonen, Persson, & Sivertsen, 1992; Luukkonen, Tijssen, Persson, & Sivertsen, 1993).

In this article, our focus is methodological. We study various measures for deriving similarities from cooccurrence data. Basically, there are two approaches that can be taken to derive similarities from cooccurrence data. We refer to these approaches as the direct and the indirect approach, but the approaches are also known as the local and the global approach (Ahlgren, Jarnevning, & Rousseau, 2003; Jarnevning, 2008). Similarity measures that implement the direct approach are referred to as direct similarity measures in this article, while similarity measures that implement the indirect approach are referred to as indirect similarity measures.

The indirect approach to derive similarities from cooccurrence data relies on cooccurrence profiles. The cooccurrence profile of an object is a vector that contains the number of cooccurrences of the object with each other object. Indirect similarity measures determine the similarity between two objects by comparing the cooccurrence profiles of the objects. The indirect approach is mainly used for cocitation data (e.g., McCain, 1990, 1991; White & Griffith, 1981; White & McCain, 1998). From a theoretical point of view, the approach is quite well understood (Ahlgren et al., 2003; Van Eck & Waltman, 2008).

In this article, we focus most of our attention on the direct approach to derive similarities from cooccurrence data. Direct similarity measures determine the similarity between two objects by taking the number of cooccurrences of the objects and adjusting this number for the total number of occurrences or cooccurrences of each of the objects. Researchers use several different direct similarity measures. The cosine and the Jaccard index are especially popular,

Received January 7, 2009; revised February 19, 2009; accepted February 20, 2009

© 2009 ASIS&T • Published online 13 April 2009 in Wiley InterScience (www.interscience.wiley.com). DOI: 10.1002/asi.21075

but other measures are also regularly used. However, relatively little is known about the theoretical properties of the various measures. Also, there is no consensus among researchers on which measure is most appropriate for a particular purpose. In this article, we theoretically analyze some well-known direct similarity measures and we compare their properties. We also study the behavior of the measures empirically. Usually, when a direct similarity measure is applied to cooccurrence data, the purpose is to normalize the data, that is, to correct the data for differences in the total number of occurrences or cooccurrences of objects. The main question that we try to answer in this article is therefore as follows: Which direct similarity measures are appropriate for normalizing cooccurrence data and which are not? An interesting finding is that despite their popularity, the cosine and the Jaccard index turn out not to be appropriate measures for normalization purposes. We argue that an appropriate measure for normalizing cooccurrence data is the association strength (Van Eck & Waltman, 2007; Van Eck, Waltman, Van den Berg, & Kaymak, 2006), also referred to as the proximity index (e.g., Peters & Van Raan, 1993a; Rip & Courtial, 1984) or the probabilistic affinity index (e.g., Zitt, Bassecoulard, & Okubo, 2000). Although this measure is somewhat less well-known, it turns out to have the right theoretical properties for normalizing cooccurrence data.

This article is organized as follows. We first provide an overview of the most popular direct similarity measures. We then analyze these measures theoretically. We also look for empirical relations among the measures. Finally, we answer the question which direct similarity measures are appropriate for normalizing cooccurrence data and which are not.

Overview of Direct Similarity Measures

In this section, we provide an overview of the most popular direct similarity measures. The overview is based on a survey of the scientometric literature.

We first introduce some mathematical notation. Let \mathbf{O} denote an occurrence matrix of order $m \times n$. The columns of \mathbf{O} correspond with the objects of which we want to analyze the cooccurrences. There are n such objects, denoted by $1, \dots, n$. The objects can be, for example, authors (e.g., White & McCain, 1998), countries (e.g., Glänzel, 2001; Zitt et al., 2000), documents (e.g., Gmür, 2003; Klavans & Boyack, 2006b), journals (e.g., Boyack, Klavans, & Börner, 2005; Klavans & Boyack, 2006a), Web pages (e.g., Vaughan, 2006; Vaughan & You, 2006), or words (e.g., Kopcsa & Schiebel, 1998). The rows of \mathbf{O} usually correspond with documents. m then denotes the number of documents on which the cooccurrence analysis is based. Sometimes the rows of \mathbf{O} do not correspond with documents. In Web co-link analysis, for example, the rows of \mathbf{O} correspond with Web pages (e.g., Vaughan, 2006; Vaughan & You, 2006). Throughout this article, however, we assume for simplicity that the rows of \mathbf{O} always correspond with documents. Another assumption that we make is that \mathbf{O} is a binary matrix, that is, each element of \mathbf{O} equals either zero or one. Let o_{ki} denote the element in the

k th row and i th column of \mathbf{O} . o_{ki} equals one if object i occurs in the document that corresponds with the k th row of \mathbf{O} , and it equals zero otherwise. Let \mathbf{C} denote the cooccurrence matrix of the objects $1, \dots, n$. \mathbf{C} is a symmetric non-negative matrix of order $n \times n$. Let c_{ij} denote the element in the i th row and j th column of \mathbf{C} . For $i \neq j$, c_{ij} equals the number of cooccurrences of objects i and j . For $i = j$, c_{ij} equals the number of occurrences of object i . Clearly, for all i and j ,

$$c_{ij} = \sum_{k=1}^m o_{ki} o_{kj}. \quad (1)$$

It follows from this that $\mathbf{C} = \mathbf{O}^T \mathbf{O}$, where \mathbf{O}^T denotes the transpose of \mathbf{O} . Moreover, the assumption that \mathbf{O} is a binary matrix implies that \mathbf{C} is an integer matrix.

As we discussed in the Introduction section, there are two types of measures for determining similarities between objects based on cooccurrence data. We refer to these two types of measures as direct similarity measures and indirect similarity measures. Indirect similarity measures, also known as global similarity measures (Ahlgren et al., 2003; Jarnevich, 2008), determine the similarity between two objects i and j by comparing the i th and the j th row (or column) of the cooccurrence matrix \mathbf{C} . The more similar the cooccurrence profiles in these two rows (or columns) of \mathbf{C} , the higher the similarity between i and j . Indirect similarity measures are especially popular for author cocitation analysis (e.g., McCain, 1990; White & Griffith, 1981; White & McCain, 1998) and journal cocitation analysis (e.g., McCain, 1991). We refer to Ahlgren et al. (2003) and Van Eck and Waltman (2008) for a detailed discussion of the properties of various indirect similarity measures. In this article, we focus most of our attention on direct similarity measures, also known as local similarity measures (Ahlgren et al., 2003; Jarnevich, 2008). Direct similarity measures determine the similarity between two objects i and j by taking the number of cooccurrences of i and j and adjusting this number for the total number of occurrences or cooccurrences of i and the total number of occurrences or cooccurrences of j . We note that in some studies similarities between objects are determined by comparing columns of the occurrence matrix \mathbf{O} (e.g., Leydesdorff & Vaughan, 2006; Schneider, Larsen, & Ingwersen, 2009). In most cases, this approach is mathematically equivalent to the use of a direct similarity measure.¹

Let s_i denote either the total number of occurrences of object i or the total number of cooccurrences of object i . In the first case, we have

$$s_i = c_{ii} = \sum_{k=1}^m o_{ki}, \quad (2)$$

¹Leydesdorff and Vaughan (2006) and Schneider et al. (2009) use the Pearson correlation to compare columns of the occurrence matrix \mathbf{O} . As shown by Guilford (1973), applying the Pearson correlation to a binary occurrence matrix is mathematically equivalent to applying the so-called phi coefficient to the corresponding cooccurrence matrix.

while in the second case, we have

$$s_i = \sum_{j=1, j \neq i}^n c_{ij}. \quad (3)$$

Both definitions are used in scientometric research (see also Leydesdorff, 2008), but the first definition seems to be more popular. We now provide a formal definition of a direct similarity measure.

Definition 1. A *direct similarity measure* is defined as a function $S(c_{ij}, s_i, s_j)$ that has the following three properties:

- The domain of $S(c_{ij}, s_i, s_j)$ equals

$$D_S = \{(c_{ij}, s_i, s_j) \in \mathbf{R}^3 | 0 \leq c_{ij} \leq \min(s_i, s_j) \text{ and } s_i, s_j > 0\}, \quad (4)$$

where \mathbf{R} denotes the set of all real numbers.

- The range of $S(c_{ij}, s_i, s_j)$ is a subset of \mathbf{R} .
- $S(c_{ij}, s_i, s_j)$ is symmetric in s_i and s_j , that is, $S(c_{ij}, s_i, s_j) = S(c_{ij}, s_j, s_i)$ for all $(c_{ij}, s_i, s_j) \in D_S$.

Based on this definition, a number of observations can be made. First, the definition does not require that c_{ij} , s_i , and s_j have integer values. Allowing for non-integer values of c_{ij} , s_i , and s_j simplifies the mathematical analysis of direct similarity measures. Second, although most direct similarity measures take values between zero and one, the definition allows measures to have a different range. And third, because the definition requires direct similarity measures to be symmetric in s_i and s_j , it does not cover asymmetric similarity measures such as those discussed by Egghe and Michel (2002, 2003). As a final observation, we note that Definition 1 is quite general. More specific definitions for special classes of direct similarity measures will be provided later on in this article. We now define the notion of monotonic relatedness of direct similarity measures.

Definition 2. Two direct similarity measures $S_1(c_{ij}, s_i, s_j)$ and $S_2(c_{ij}, s_i, s_j)$ are said to be *monotonically related* if and only if

$$\begin{aligned} S_1(c_{ij}, s_i, s_j) &< S_1(c'_{ij}, s'_i, s'_j) \\ \iff S_2(c_{ij}, s_i, s_j) &< S_2(c'_{ij}, s'_i, s'_j) \end{aligned} \quad (5)$$

for all $(c_{ij}, s_i, s_j), (c'_{ij}, s'_i, s'_j) \in D_S$.

Monotonic relatedness of direct similarity measures is important because certain multivariate analysis techniques that are frequently used in scientometric research are insensitive to monotonic transformations of similarities. This is, for example, the case for ordinal or non-metric multidimensional scaling (e.g., Borg & Groenen, 2005) and for single linkage and complete linkage hierarchical clustering (e.g., Anderberg, 1973).

Based on a survey of the literature, we have identified the most popular direct similarity measures in the

field of scientometrics. These measures are defined as follows:

$$S_A(c_{ij}, s_i, s_j) = \frac{c_{ij}}{s_i s_j}, \quad (6)$$

$$S_C(c_{ij}, s_i, s_j) = \frac{c_{ij}}{\sqrt{s_i s_j}}, \quad (7)$$

$$S_I(c_{ij}, s_i, s_j) = \frac{c_{ij}}{\min(s_i, s_j)}, \quad (8)$$

$$S_J(c_{ij}, s_i, s_j) = \frac{c_{ij}}{s_i + s_j - c_{ij}}. \quad (9)$$

We refer to these measures as, respectively, the association strength, the cosine, the inclusion index, and the Jaccard index. Assuming that c_{ij} is an integer, each of the measures takes values between zero and one. Moreover, it is not difficult to see that the measures satisfy

$$\begin{aligned} S_A(c_{ij}, s_i, s_j) &\leq S_J(c_{ij}, s_i, s_j) \\ &\leq S_C(c_{ij}, s_i, s_j) \leq S_I(c_{ij}, s_i, s_j). \end{aligned} \quad (10)$$

We now discuss each of the measures.

The association strength defined in Equation 6 is used by Van Eck and Waltman (2007) and Van Eck et al. (2006).² Under various names, the measure is also used in a number of other studies. Hinze (1994), Leclerc and Gagné (1994), Peters and Van Raan (1993a), and Rip and Courtial (1984) refer to the measure as the proximity index, while Leydesdorff (2008) and Zitt et al. (2000) refer to it as the probabilistic affinity (or activity) index. Luukkonen et al. (1992, 1993) also employ the measure, but in their work, it does not have a name. The association strength is proportional to the ratio between, on the one hand, the observed number of cooccurrences of objects i and j and, on the other hand, the expected number of cooccurrences of objects i and j under the assumption that occurrences of i and j are statistically independent. We will come back to this interpretation later on in this article. The association strength corresponds with the pseudo-cosine measure discussed by Jones and Furnas (1987) and is monotonically related to the (pointwise) mutual information measure used in the field of computational linguistics (e.g., Church & Hanks, 1990; Manning & Schütze, 1999). Measures equivalent to the association strength sometimes also appear outside the field of scientometrics (T.F. Cox & Cox, 2001; M.A.A. Cox & Cox, 2008; Hubálek, 1982).

The cosine defined in Equation 7 equals the ratio between, on the one hand, the number of times that objects i and j are observed together and, on the other hand, the geometric mean of the number of times that object i is observed and the number of times that object j is observed. The measure can be interpreted as the cosine of the angle between the i th and the j th column of the occurrence matrix \mathbf{O} , where the

²The definition of the association strength used in these articles differs slightly from the definition provided in Equation 6. However, since the two definitions are proportional to each other, the difference between them is not important. Throughout this section, direct similarity measures that are proportional to each other will simply be regarded as equivalent.

columns of \mathbf{O} are regarded as vectors in an m -dimensional space (e.g., Salton & McGill, 1983). The cosine seems to be the most popular direct similarity measure in the field of scientometrics. Frequently cited studies in which the measure is used include Braam, Moed, and Van Raan (1991a, 1991b), Klavans and Boyack (2006a), Leydesdorff (1989), Peters and Van Raan (1993b), Peters, Braam, and Van Raan (1995), Small (1994), Small and Sweeney (1985), and Small, Sweeney, and Greenlee (1985). The popularity of the cosine is largely due to the work of Salton in the field of information retrieval (e.g., Salton, 1963; Salton & McGill, 1983). The cosine is therefore sometimes referred to as Salton's measure (e.g., Glänzel, 2001; Glänzel, Schubert, & Czerwon, 1999; Luukkonen et al., 1993; Schubert & Braun, 1990) or as the Salton index (e.g., Morillo, Bordons, & Gómez, 2003). In some studies, a measure called the equivalence index is used (e.g., Callon, Courtial, & Laville, 1991; Kostoff, Eberhart, & Toothman, 1999; Law & Whittaker, 1992; Palmer, 1999). This measure equals the square of the cosine. Outside the fields of scientometrics and information retrieval, the cosine is also known as the Ochiai coefficient (e.g., T.F. Cox & Cox, 2001; M.A.A. Cox & Cox, 2008; Hubálek, 1982; Sokal & Sneath, 1963).

Examples of the use of the inclusion index defined in Equation 8 can be found in the work of Kostoff, Del Río, Humenik, García, and Ramírez (2001), McCain (1995), Peters and Van Raan (1993a), Rip and Courtial (1984), Tijssen (1992, 1993), and Tijssen and Van Raan (1989). We note that a measure somewhat different from the one defined in Equation 8 is sometimes also called the inclusion index (e.g., Braam et al., 1991a; Kostoff et al., 1999; Peters et al., 1995; Qin, 2000). In the field of information retrieval, the inclusion index is referred to as the overlap measure (e.g., Jones & Furnas, 1987; Rorvig, 1999; Salton & McGill, 1983). More in general, the inclusion index is sometimes called the Simpson coefficient (e.g., T.F. Cox & Cox, 2001; M.A.A. Cox & Cox, 2008; Hubálek, 1982).

The Jaccard index defined in Equation 9 equals the ratio between, on the one hand, the number of times that objects i and j are observed together and, on the other hand, the number of times that at least one of the two objects is observed. Small uses the Jaccard index in his early work on cocitation analysis (e.g., Small, 1973, 1981; Small & Greenlee, 1980). Other work in which the Jaccard index is used includes Heimeriks, Hörlesberger, and Van den Besselaar (2003), Kopcsa and Schiebel (1998), Peters and Van Raan (1993a), Peters et al. (1995), Rip and Courtial (1984), Van Raan and Tijssen (1993), Vaughan (2006), and Vaughan and You (2006). As shown by Anderberg (1973), the Jaccard index is monotonically related to the Dice coefficient, which is a well-known measure in information retrieval (e.g., Jones & Furnas, 1987; Rorvig, 1999; Salton & McGill, 1983) and other fields (e.g., T.F. Cox & Cox, 2001; M.A.A. Cox & Cox, 2008; Hubálek, 1982; Sokal & Sneath, 1963).

We note that, in addition to the four direct similarity measures discussed above, many more direct similarity measures have been used in scientometric research. However, the above

TABLE 1. Relations among various direct similarity measures.

Measure	Alternative names	Monotonically related measures
Association strength	Probabilistic affinity index Proximity index Pseudo-cosine	(Pointwise) Mutual information
Cosine	Ochiai coefficient Salton's index/measure	Equivalence index
Inclusion index	Overlap measure Simpson coefficient	
Jaccard index		Dice coefficient

four measures are by far the most popular ones, and we therefore focus most of our attention on them in this article. The relations among various direct similarity measures are summarized in Table 1.

In the field of scientometrics, a number of studies have been performed in which different direct similarity measures are compared with each other. Boyack et al. (2005), Gmür (2003), Klavans and Boyack (2006a), Leydesdorff (2008), Luukkonen et al. (1993), and Peters and Van Raan (1993a) report results of empirical comparisons of different measures. Theoretical analyses of relations between different measures can be found in the work of Egghe (2009) and Hamers et al. (1989). Egghe and Rousseau (2006) also theoretically studied properties of various measures. Schneider and Borlund (2007a, 2007b) provide an extensive discussion of the issue of comparing different measures. Other work that might be of interest has been done in the field of information retrieval. In the information retrieval literature, Chung and Lee (2001) and Rorvig (1999) discuss empirical comparisons of different direct similarity measures, and Jones and Furnas (1987) present a theoretical comparison.³ We further note that general overviews of a large number of direct similarity measures and their properties can be found in the statistical literature (Anderberg, 1973; T.F. Cox & Cox, 2001; M.A.A. Cox & Cox, 2008; Gower, 1985; Gower & Legendre, 1986) and also in the biological literature (Hubálek, 1982; Sokal & Sneath, 1963).

Set-Theoretic Similarity Measures

In this and the next section, we are concerned with two special classes of direct similarity measures. We discuss the class of set-theoretic similarity measures in this section and the class of probabilistic similarity measures in the next section. It turns out that there is a fundamental difference between the

³The results reported by Jones and Furnas are probably not very relevant to scientometric research. This is because Jones and Furnas focus on the effect of term weights on similarity measures. In scientometric research, there is no natural analogue to the term weights used in information retrieval. The reason for this is that the occurrence matrices used in scientometric research contain elements that are usually restricted to zero and one, while the document-term matrices used in information retrieval contain term weights that often do not have this restriction.

cosine, the inclusion index, and the Jaccard index, on the one hand, and the association strength, on the other hand. The first three measures all belong to the class of set-theoretic similarity measures, while the last measure belongs to the class of probabilistic similarity measures. We assume from now on that s_i denotes the total number of occurrences of object i , that is, we assume that the definition of s_i in Equation 2 is adopted. From a theoretical point of view, this definition is more convenient than the definition of s_i in Equation 3. We note that proofs of the theoretical results that we present in this and the next section are provided in the appendix.

Each column of an occurrence matrix can be seen as a representation of a set, namely the set of all documents in which a certain object occurs (cf. Egghe & Rousseau, 2006). Consequently, a natural approach to determine the similarity between two objects i and j seems to be to determine the similarity between, on the one hand, the set of all documents in which i occurs and, on the other hand, the set of all documents in which j occurs. We refer to direct similarity measures that take this approach as set-theoretic similarity measures. In other words, set-theoretic similarity measures are direct similarity measures that are based on the notion of similarity between sets. In this section, we theoretically analyze the properties of set-theoretic similarity measures. We note that these properties are also studied theoretically by Baulieu (1989, 1997), Egghe and Michel (2002, 2003), Egghe and Rousseau (2006), and Janson and Vegelius (1981).

There are a number of properties of which we believe that it is natural to expect that any set-theoretic similarity measure $S(c_{ij}, s_i, s_j)$ has them. Three of these properties are given below.

Property 1. If $c_{ij} = 0$, then $S(c_{ij}, s_i, s_j)$ takes its minimum value.

Property 2. For all $\alpha > 0$, $S(\alpha c_{ij}, \alpha s_i, \alpha s_j) = S(c_{ij}, s_i, s_j)$.

Property 3. If $s'_i > s_i$ and $c_{ij} > 0$, then $S(c_{ij}, s'_i, s_j) < S(c_{ij}, s_i, s_j)$.

Property 1 is based on the idea that the similarity between two sets should be minimal if the sets are disjoint, that is, if they have no elements in common. Property 2 is based on the idea that the similarity between two sets should remain unchanged in the case of a proportional increase or decrease in both the number of elements of each of the sets and the number of elements of the intersection of the sets. Egghe and Rousseau (2006) refer to this idea as replication invariance. It underlies the notion of Lorenz similarity that is studied by Egghe and Rousseau. Janson and Vegelius (1981) also use a similar idea, calling it homogeneity. Property 3 is based on the idea that the similarity between two sets should decrease if an element is added to one of the sets and this element does not belong to the other set. Baulieu (1989, 1997) uses a similar idea. It is not difficult to see that Properties 1, 2, and 3 are independent of each other, that is, none of the properties is implied by the others. We regard Properties 1, 2, and 3 as the characterizing properties of set-theoretic similarity measures. This is formally stated in the following definition.

Definition 3. A *set-theoretic similarity measure* is defined as a direct similarity measure $S(c_{ij}, s_i, s_j)$ that has Properties 1, 2, and 3.

This definition implies that the cosine defined in Equation 7 and the Jaccard index defined in Equation 9 are set-theoretic similarity measures. The association strength defined in Equation 6 does not have Property 2 and is therefore not a set-theoretic similarity measure. The inclusion index defined in Equation 8 is also not a set-theoretic similarity measure. This is because the inclusion index does not have Property 3. However, the inclusion index does have the following property, which is a weakened version of Property 3.

Property 4. If $s'_i > s_i$ and $c_{ij} > 0$, then $S(c_{ij}, s'_i, s_j) \leq S(c_{ij}, s_i, s_j)$.

This property naturally leads to the following definition.

Definition 4. A *weak set-theoretic similarity measure* is defined as a direct similarity measure $S(c_{ij}, s_i, s_j)$ that has Properties 1, 2, and 4.

It follows from this definition that the inclusion index is a weak set-theoretic similarity measure. We note that our definition of a set-theoretic similarity measure seems to be more restrictive than the definition of a Lorenz similarity function that is provided by Egghe and Rousseau (2006). This is because a Lorenz similarity function need not have Properties 1 and 3.

In addition to Properties 1, 2, and 3, there are some other properties that we consider indispensable for any set-theoretic similarity measure $S(c_{ij}, s_i, s_j)$. Four of these properties are given below.

Property 5. If $S(c_{ij}, s_i, s_j)$ takes its minimum value, then $c_{ij} = 0$.

Property 6. If $c_{ij} = s_i = s_j$, then $S(c_{ij}, s_i, s_j)$ takes its maximum value.

Property 7. If $S(c_{ij}, s_i, s_j)$ takes its maximum value, then $c_{ij} = s_i = s_j$.

Property 8. For all $\alpha > 0$, if $c_{ij} < s_i$ or $c_{ij} < s_j$, then $S(c_{ij} + \alpha, s_i + \alpha, s_j + \alpha) > S(c_{ij}, s_i, s_j)$.

Properties 5, 6, and 7 are based on the idea that the similarity between two sets should be minimal only if the sets are disjoint and that it should be maximal if and only if the sets are equal. Property 8 is based on the idea that the similarity between two sets should increase if the same element is added to both sets. It turns out that Properties 5, 6, 7, and 8 are implied by Properties 1, 2, and 3. This is stated by the following proposition.

Proposition 1. All set-theoretic similarity measures $S(c_{ij}, s_i, s_j)$ have Properties 5, 6, 7, and 8.

We note that weak set-theoretic similarity measures need not have Properties 5, 7, and 8. They do have Property 6.

TABLE 2. Summary of the properties of a number of direct similarity measures. If a measure has a certain property, this is indicated using a \times symbol.

	Property												
	1	2	3	4	5	6	7	8	9	10	11	12	13
Association strength	\times		\times	\times	\times		\times		\times			\times	\times
Cosine	\times			\times									
Inclusion index	\times	\times		\times	\times	\times					\times	\times	
Jaccard index	\times			\times									

We now consider the following two properties.

Property 9. If $s'_i s'_j > s_i s_j$ and $c_{ij} > 0$, then $S(c_{ij}, s'_i, s'_j) < S(c_{ij}, s_i, s_j)$. If $s'_i s'_j = s_i s_j$, then $S(c_{ij}, s'_i, s'_j) = S(c_{ij}, s_i, s_j)$.

Property 10. If $s'_i + s'_j > s_i + s_j$ and $c_{ij} > 0$, then $S(c_{ij}, s'_i, s'_j) < S(c_{ij}, s_i, s_j)$. If $s'_i + s'_j = s_i + s_j$, then $S(c_{ij}, s'_i, s'_j) = S(c_{ij}, s_i, s_j)$.

It is easy to see that these properties both imply Property 3. Hence, Properties 9 and 10 are both stronger than Property 3. It can further be seen that the cosine has Property 9 and that the Jaccard index has Property 10. The following two propositions indicate the importance of Properties 9 and 10.

Proposition 2. All set-theoretic similarity measures $S(c_{ij}, s_i, s_j)$ that have Property 9 are monotonically related to the cosine defined in Equation 7.

Proposition 3. All set-theoretic similarity measures $S(c_{ij}, s_i, s_j)$ that have Property 10 are monotonically related to the Jaccard index defined in Equation 9.

It follows from Proposition 2 that Properties 1, 2, and 9 characterize the class of all set-theoretic similarity measures that are monotonically related to the cosine. Likewise, it follows from Proposition 3 that Properties 1, 2, and 10 characterize the class of all set-theoretic similarity measures that are monotonically related to the Jaccard index. We now apply a similar idea to the inclusion index. The inclusion index has the following property.

Property 11. If $\min(s'_i, s'_j) > \min(s_i, s_j)$ and $c_{ij} > 0$, then $S(c_{ij}, s'_i, s'_j) < S(c_{ij}, s_i, s_j)$. If $\min(s'_i, s'_j) = \min(s_i, s_j)$, then $S(c_{ij}, s'_i, s'_j) = S(c_{ij}, s_i, s_j)$.

This property implies Property 4. Together with Properties 1 and 2, Property 11 characterizes the class of all weak set-theoretic similarity measures that are monotonically related to the inclusion index. This is an immediate consequence of the following proposition.

Proposition 4. All weak set-theoretic similarity measures $S(c_{ij}, s_i, s_j)$ that have Property 11 are monotonically related to the inclusion index defined in Equation 8.

In the above discussion, we have introduced a large number of properties that a direct similarity measure may or may not have. For convenience, in Table 2, we summarize for the association strength, the cosine, the inclusion index, and the Jaccard index, which of these the properties they have and which they do not have. We note that the last two properties in the table will be introduced in the next section.

To provide some additional insight into the relations among various (weak and non-weak) set-theoretic similarity measures, we now introduce what we call the generalized similarity index (for a similar idea, see Warrens, 2008).

Definition 5. The *generalized similarity index* is defined as a direct similarity measure that is given by

$$S_G(c_{ij}, s_i, s_j; p) = \frac{2^{1/p} c_{ij}}{(s_i^p + s_j^p)^{1/p}}, \quad (11)$$

where p denotes a parameter that takes values in $\mathbf{R} \setminus \{0\}$.

For all values of the parameter p , the generalized similarity index takes values between zero and one. The index equals the ratio between, on the one hand, the number of times that objects i and j are observed together and, on the other hand, a power mean of the number of times that object i is observed and the number of times that object j is observed. (Power means, also known as generalized means or Hölder means, are a generalization of arithmetic, geometric, and harmonic means.) An interesting property of the generalized similarity index is that, for various values of p , the index reduces to a well-known (weak or non-weak) set-theoretic similarity measure. More specifically, it can be seen that

$$\lim_{p \rightarrow -\infty} S_G(c_{ij}, s_i, s_j; p) = \frac{c_{ij}}{\min(s_i, s_j)}, \quad (12)$$

$$S_G(c_{ij}, s_i, s_j; -1) = \frac{1}{2} \left(\frac{c_{ij}}{s_i} + \frac{c_{ij}}{s_j} \right), \quad (13)$$

$$\lim_{p \rightarrow 0} S_G(c_{ij}, s_i, s_j; p) = \frac{c_{ij}}{\sqrt{s_i s_j}}, \quad (14)$$

$$S_G(c_{ij}, s_i, s_j; 1) = \frac{2 c_{ij}}{s_i + s_j}, \quad (15)$$

$$S_G(c_{ij}, s_i, s_j; 2) = \frac{\sqrt{2}c_{ij}}{\sqrt{s_i^2 + s_j^2}}, \quad (16)$$

$$\lim_{p \rightarrow \infty} S_G(c_{ij}, s_i, s_j; p) = \frac{c_{ij}}{\max(s_i, s_j)}, \quad (17)$$

where Equations 12, 14, and 17 follow from the properties of power means as discussed by, for example, Hardy, Littlewood, and Pólya (1952). Equations 12 and 13 indicate that for $p \rightarrow -\infty$ the generalized similarity index equals the inclusion index and that for $p = -1$ it equals the so-called joint conditional probability measure that is used by McCain (1995). The latter measure is more generally known as one of the Kulczynski coefficients (e.g., T.F. Cox & Cox, 2001; M.A.A. Cox & Cox, 2008; Hubálek, 1982; Sokal & Sneath, 1963). It is easy to see that this measure is a set-theoretic similarity measure. Equations 14 and 15 indicate that for $p \rightarrow 0$ the generalized similarity index equals the cosine and that for $p = 1$ it equals the Dice coefficient. It follows from Equations 9 and 15 that

$$S_G(c_{ij}, s_i, s_j; 1) = \frac{2S_J(c_{ij}, s_i, s_j)}{S_J(c_{ij}, s_i, s_j) + 1}, \quad (18)$$

which implies that for $p = 1$ the generalized similarity index is monotonically related to the Jaccard index. Equations 16 and 17 indicate that for $p = 2$ and $p \rightarrow \infty$ the generalized similarity index equals, respectively, the measures N and O_2 that are studied by Egghe and Michel (2002, 2003) and Egghe and Rousseau (2006). It is clear that N is a set-theoretic similarity measure and that O_2 is a weak set-theoretic similarity measure. Measures equivalent to Equation 17 are also discussed by T.F. Cox and Cox (2001), M.A.A. Cox and Cox (2008) and Hubálek (1982).

The following proposition points out an important property of the generalized similarity index.

Proposition 5. For all finite values of the parameter p , the generalized similarity index defined in Equation 11 is a set-theoretic similarity measure.

This proposition states that the generalized similarity index describes an entire class of set-theoretic similarity measures. Each member of this class corresponds with a particular value of p . Only in the limit case in which $p \rightarrow \pm\infty$, the generalized similarity index is not a set-theoretic similarity measure. In this limit case, the generalized similarity index is a weak set-theoretic similarity measure.

Probabilistic Similarity Measures

In the previous section, we discussed the class of set-theoretic similarity measures. The cosine, the inclusion index, and the Jaccard index turned out to be (weak or non-weak) set-theoretic similarity measures. The association strength, however, turned out not to belong to the class of set-theoretic similarity measures. In this section, we discuss the class of

probabilistic similarity measures. This is the class to which the association strength turns out to belong.

We are interested in direct similarity measures $S(c_{ij}, s_i, s_j)$ that have the following two properties.

Property 12. If $s_1 = s_2 = \dots = s_n$, then $S(c_{ij}, s_i, s_j) = \alpha c_{ij}$ for all $i \neq j$ and for some $\alpha > 0$.

Property 13. For all $\alpha > 0$, $S(\alpha c_{ij}, \alpha s_i, s_j) = S(c_{ij}, s_i, s_j)$.

Property 12 requires that, if all objects occur equally frequently, the similarity between two objects is proportional to the number of cooccurrences of the objects. Property 13 requires that the similarity between two objects remains unchanged in the case of a proportional increase or decrease in, on the one hand, the number of cooccurrences of the objects and, on the other hand, the number of occurrences of one of the objects. (Notice the difference between this property and Property 2.) We regard Properties 12 and 13 as the characterizing properties of probabilistic similarity measures. This results in the following definition.

Definition 6. A *probabilistic similarity measure* is defined as a direct similarity measure $S(c_{ij}, s_i, s_j)$ that has Properties 12 and 13.

The cosine, the inclusion index, and the Jaccard index do not have Property 13 and therefore are not probabilistic similarity measures. The association strength, on the other hand, is a probabilistic similarity measure because it has both Property 12 and Property 13. In this respect, the association strength is quite unique, as the following proposition indicates.

Proposition 6. All probabilistic similarity measures are proportional to the association strength defined in Equation 6.

This proposition states that the class of probabilistic similarity measures consists only of the association strength and of measures that are proportional to the association strength. There are no other measures that belong to the class of probabilistic similarity measures. The following result is an immediate consequence of Proposition 6.

Corollary 7. A direct similarity measure cannot be both a (weak or non-weak) set-theoretic similarity measure and a probabilistic similarity measure.

This result makes clear that there is a fundamental difference between set-theoretic similarity measures and probabilistic similarity measures. In other words, there is a fundamental difference between measures such as the cosine, the inclusion index, and the Jaccard index, on the one hand, and the association strength, on the other hand. We will come back to this difference later on in this article.

We now explain the rationale for Properties 12 and 13. To do so, we first discuss why direct similarity measures are applied to cooccurrence data. The number of cooccurrences of two objects can be seen as the result of two independent effects. We refer to these effects as the similarity effect and

the size effect.⁴ The similarity effect is the effect that, other things being equal, more similar objects have more cooccurrences. The size effect is the effect that, other things being equal, an object that occurs more frequently has more cooccurrences with other objects. If one is interested in the similarity between two objects, the number of cooccurrences of the objects is in general not an appropriate measure. This is because, due to the size effect, the number of cooccurrences is likely to give a distorted picture of the similarity between the objects (see also Waltman & Van Eck, 2007). Two frequently occurring objects, for example, may have a large number of cooccurrences and may therefore look very similar. However, it is quite well possible that the large number of cooccurrences of the objects is completely due to their high frequency of occurrence (i.e., the size effect) and has nothing to do with their similarity. Usually, when a direct similarity measure is applied to cooccurrence data, the aim is to correct the data for the size effect.

Based on the above discussion, the idea underlying Property 12 can be explained as follows. Property 12 is concerned with the behavior of a direct similarity measure in the special case in which all objects occur equally frequently. In this special case, the size effect is equally strong for all objects, which means that unlike in the more general case, the number of cooccurrences of two objects is an appropriate measure of the similarity between the objects. Taking this into account, it is natural to expect that in the special case considered by Property 12 a direct similarity measure does not transform the cooccurrence frequencies of objects in any significant way. Property 12 implements this idea by requiring that, if all objects occur equally frequently, the similarity between two objects is proportional to the number of cooccurrences of the objects.

We now consider Property 13. The idea underlying this property is best clarified by means of an example. Consider an arbitrary object i , and suppose that the total number of occurrences of i doubles. It can then be expected that the total number of cooccurrences of i also double, at least approximately. Suppose that the total number of cooccurrences of i indeed doubles and that the new cooccurrences of i are distributed over the other objects in the same way as the old cooccurrences of i . This simply means that the number of cooccurrences of i with each other object doubles. We believe that this increase in the number of occurrences and cooccurrences of i should not have any influence on the similarities between i and the other objects. This is because the number of occurrences of i and the number of cooccurrences of i with each other object have all increased proportionally, namely by a factor of two. Hence, relatively speaking, the frequency with which i cooccurs with each other object has not changed. This means that the increase in the number of cooccurrences of i with each other object is completely due to the size effect and has not been caused by the similarity effect. Taking this into

account, it is natural to expect that the similarities between i and the other objects remain unchanged. Property 13 implements this idea. It does so not only for the case in which the number of occurrences and cooccurrences of an object doubles but more generally for any proportional increase or decrease in the number of occurrences and cooccurrences of an object. We note that the idea underlying Property 13 is not new. Ahlgren et al. (2003) and Van Eck and Waltman (2008) study properties of indirect similarity measures. A property that turns out to be particularly important is the so-called property of coordinate-wise scale invariance. It is interesting to note that this property relies on exactly the same idea as Property 13. Hence, direct similarity measures that have Property 13 and indirect similarity measures that have the property of coordinate-wise scale invariance are based on similar principles.

Finally, we discuss the probabilistic interpretation of probabilistic similarity measures (see also Leclerc & Gagné, 1994; Luukkonen et al., 1992, 1993; Zitt et al., 2000). Let p_i denote the probability that object i occurs in a randomly chosen document. It is clear that $p_i = s_i/m$. If two objects i and j occur independently of each other, the probability that they cooccur in a randomly chosen document equals $p_{ij} = p_i p_j$. The expected number of cooccurrences of i and j then equals $e_{ij} = mp_{ij} = mp_i p_j = s_i s_j/m$. A natural way to measure the similarity between i and j is to calculate the ratio between on the one hand the observed number of cooccurrences of i and j and on the other hand the expected number of cooccurrences of i and j under the assumption that i and j occur independently of each other (for a similar argument in a more general context, see De Solla Price, 1981). This results in a measure that equals c_{ij}/e_{ij} . This measure has a straightforward probabilistic interpretation. If $c_{ij}/e_{ij} > 1$, i and j cooccur more frequently than would be expected by chance. If, on the other hand, $c_{ij}/e_{ij} < 1$, i and j cooccur less frequently than would be expected by chance. It is easy to see that $c_{ij}/e_{ij} = mS_A(c_{ij}, s_i, s_j)$. Hence, the measure c_{ij}/e_{ij} is proportional to the association strength and, consequently, belongs to the class of probabilistic similarity measures. Since probabilistic similarity measures are all proportional to each other (this follows from Proposition 6), they all have a similar probabilistic interpretation as the measure c_{ij}/e_{ij} .

Empirical Comparison

In the previous two sections, the differences between a number of well-known direct similarity measures were analyzed theoretically. It turned out that some measures have fundamentally different properties than others. An obvious question now is whether in practical applications there is much difference between the various measures. This is the question with which we are concerned in this section.

Leydesdorff (2008) reports the results of an empirical comparison of a number of direct and indirect similarity measures (for a theoretical explanation for some of the results, see Egghe, 2009). The measures are applied to a data set comprising the cocitation frequencies of 24 authors, 12 from

⁴The similarity effect and the size effect can be seen as analogous to what statisticians call, respectively, interaction effects and main effects.

TABLE 3. Main characteristics of the author data set, the journal data set, and the term data set.

	Author Data set	Journal Data set	Term Data set
No. objects	100	389	332
No. documents	5 463	24 106	6 235
No. occurrences	7 768	32 697	26 211
No. co-occurrences	22 520	13 378	60 640
Zeros in co-occurrence matrix (%)	26%	93%	74%

the field of information retrieval and 12 from the field of scientometrics.⁵ It turns out that the direct similarity measures are strongly correlated with each other. The Spearman rank correlations between the association strength (referred to as the probabilistic affinity or activity index), the cosine, and the Jaccard index are all above .98. Hence, for the particular data set studied by Leydesdorff, there does not seem to be much difference between various direct similarity measures.

In this section, we examine whether the results reported by Leydesdorff hold more generally. To do so, we study three data sets, one comprising cocitation frequencies of authors, one comprising cocitation frequencies of journals, and one comprising cooccurrence frequencies of terms. We refer to these data sets as, respectively, the author data set, the journal data set, and the term data set. The author data set comprises the cocitation frequencies of 100 authors in the field of information science in the 1988–1995 period. White and McCain extensively studied the data set in a well-known article (1998; see also White, 2003), and the data set is also used in one of our earlier articles (Van Eck & Waltman, 2008). The journal data set has not been studied before. The data set comprises the cocitation frequencies of 389 journals belonging to at least one of the following five subject categories of Thomson Reuters: *Business*, *Business-Finance*, *Economics*, *Management*, and *Operations Research & Management Science*. The cocitation frequencies of the journals were determined based on citations in articles published between 2005 and 2007 to articles published in 2005. The term data set comprises the cooccurrence frequencies of 332 terms in the field of computational intelligence in the 1996–2000 period. Cooccurrences of terms were counted in abstracts of articles published in important journals and conference proceedings in the computational intelligence field. For a more detailed description of the term data set, we refer to an earlier article (Van Eck & Waltman, 2007). In Table 3, we summarize the main characteristics of the three data sets that we study.

To examine how the association strength, the cosine, the inclusion index, and the Jaccard index are empirically related to each other, we analyzed each of the three data sets as follows. We first calculated for each combination of two objects the value of each of the four similarity measures. For each

combination of two similarity measures, we then drew a scatter plot that shows how the values of the two measures are related to each other. The scatter plots obtained for the author data set and the term data set are shown in Figures 1 and 2, respectively. The scatter plots obtained for the journal data set look very similar to the ones obtained for the term data set and are therefore not shown. After drawing the scatter plots, we determined for each combination of two similarity measures how strongly the values of the measures are correlated with each other. We calculated both the Pearson correlation and the Spearman correlation. The Pearson correlation was used to measure the degree to which the values of two measures are linearly related, while the Spearman correlation was used to measure the degree to which the values of two measures are monotonically related. When calculating the Pearson and Spearman correlations between the values of two measures, we only took into account values above zero.⁶ The correlations obtained for the three data sets are reported in Tables 4, 5, and 6. In each table, the values in the upper right part are Pearson correlations, while the values in the lower left part are Spearman correlations.

The scatter plots in Figures 1 and 2 clearly show that in practical applications there can be substantial differences between different direct similarity measures, which is confirmed by the correlations in Tables 4, 5, and 6. These results differ from the ones reported by Leydesdorff (2008), who finds no substantial differences between different direct similarity measures. The difference between our results and the results of Leydesdorff is probably due to the unusual nature of the data set studied in Leydesdorff, in particular the small number of objects in the data set (24 authors) and the division of the objects into two strongly separated groups (the information retrieval researchers and the scientometricians). When looking in more detail at the scatter plots in Figures 1 and 2, it can be seen that the similarity measures that are strongest related to each other are the cosine and the Jaccard index. The same observation can be made in Tables 4, 5, and 6. The relatively strong relation between the cosine and the Jaccard index has been observed before and is discussed by Egghe (2009), Hamers et al. (1989), and Leydesdorff (2008). Apart from the relation between the cosine and the Jaccard index, the relations between the different similarity measures are quite weak. This is especially the case for the relations between the association strength and the other three measures. Consider, for example, how the association strength and the inclusion index are related to each other in the term data set. As can be seen in Figure

⁵The same data set is also studied by Ahlgren et al. (2003), Leydesdorff and Vaughan (2006), and Waltman and Van Eck (2007).

⁶If two objects have zero cooccurrences, all four similarity measures have a value of zero. Cooccurrence matrices usually contain a large number of zeros (see Table 3). This leads to high correlations (close to one) between the values of the four similarity measures. We regard these high correlations as problematic because they do not properly reflect how the similarity measures are related to each other in the case of objects with a non-zero number of cooccurrences. To avoid the problem of the high correlations, we only took into account values above zero when calculating correlations between the values of the four similarity measures.

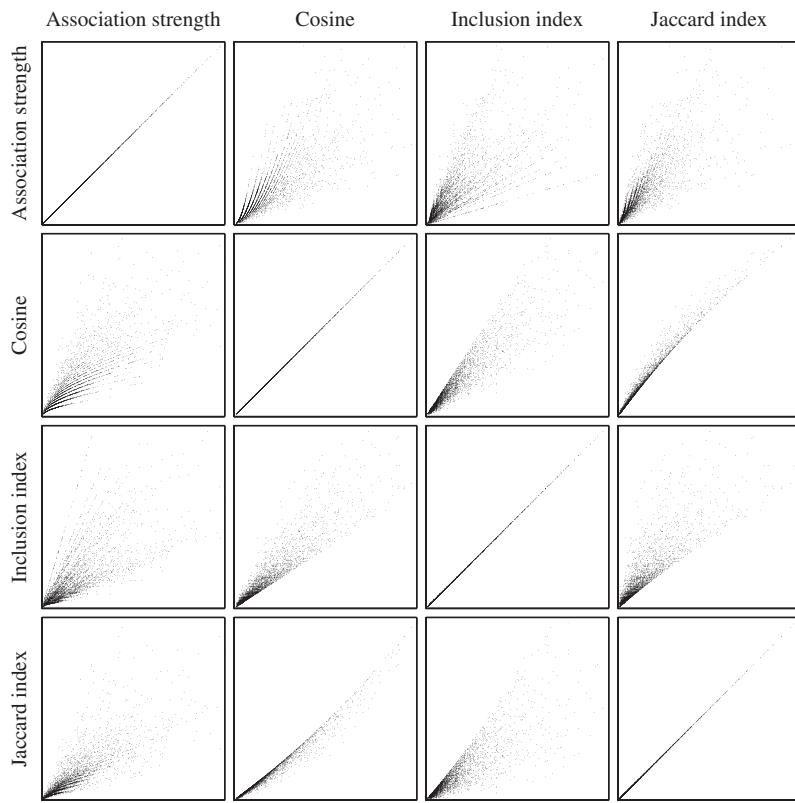


FIG. 1. Scatter plots obtained for the author data set. In each plot, the lower left corner corresponds with the origin. The scales used for the different similarity measures are not the same.

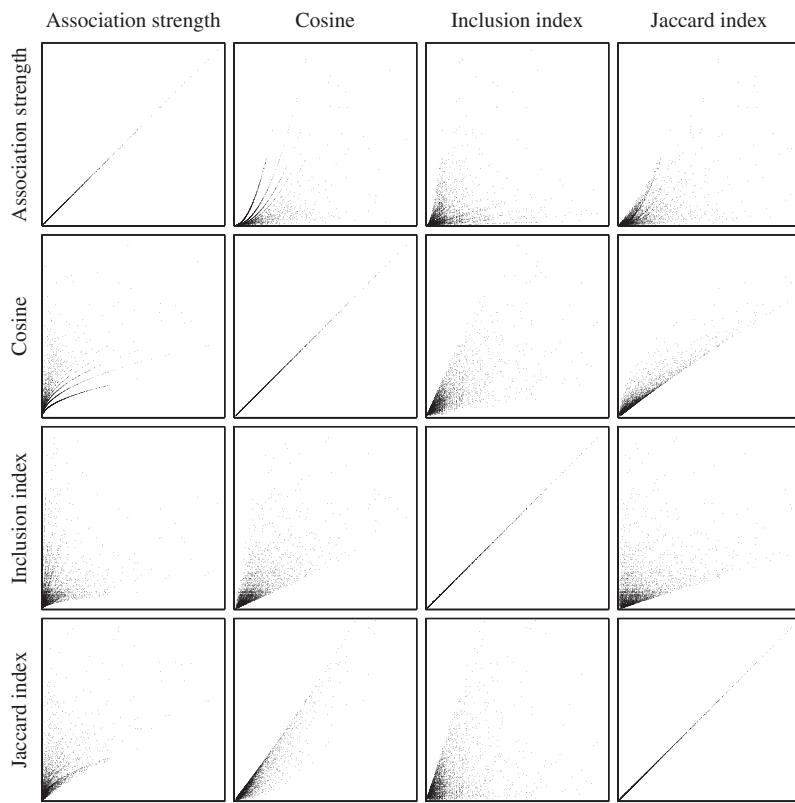


FIG. 2. Scatter plots obtained for the term data set. In each plot, the lower left corner corresponds with the origin. The scales used for the different similarity measures are not the same.

TABLE 4. Correlations obtained for the author data set.

	Association strength	Cosine	Inclusion index	Jaccard index
Association strength		.824	.721	.823
Cosine	.913		.929	.987
Inclusion index	.847	.964		.866
Jaccard index	.920	.994	.931	

TABLE 5. Correlations obtained for the journal data set.

	Association strength	Cosine	Inclusion index	Jaccard index
Association strength		.602	.556	.554
Cosine	.892		.800	.971
Inclusion index	.808	.881		.644
Jaccard index	.832	.952	.708	

TABLE 6. Correlations obtained for the term data set.

	Association strength	Cosine	Inclusion index	Jaccard index
Association strength		.653	.347	.688
Cosine	.786		.736	.950
Inclusion index	.562	.799		.511
Jaccard index	.776	.916	.520	

2, a low value of the association strength sometimes corresponds with a high value of the inclusion index and, the other way around, a low value of the inclusion index sometimes corresponds with a high value of the association strength. This clearly indicates that the relation between the two measures is rather weak, which is confirmed by the correlations in Table 6. It is further interesting to compare our empirical results with the theoretical results presented by Egghe (2009). Egghe mathematically studies relations between various (weak and non-weak) set-theoretic similarity measures under the simplifying assumption that the ratio of the number of occurrences of two objects is fixed. He proves that, under this assumption, there exist simple monotonic (often linear) relations between many measures. However, especially for the inclusion index, the scatter plots in Figures 1 and 2 do not show such relations. Our empirical results therefore seem to indicate that the practical relevance of the theoretical results presented by Egghe might be somewhat limited.

The general conclusion that can be drawn from our empirical analysis is that there are quite significant differences between various direct similarity measures and, hence, that in practical applications, it is important to use the measure that is most appropriate for one's purposes. In the next section, we discuss how an appropriate similarity measure can be chosen based on sound theoretical considerations. We focus in particular on the case in which a similarity measure is used for normalization purposes.

How to Normalize Cooccurrence Data?

As we discussed in the previous sections, there are various ways in which similarities between objects can be determined based on cooccurrence data. The different types of similarity measures that can be used are shown in Figure 3. The first decision that one has to make is whether to use a direct or an indirect similarity measure. If one decides to use a direct similarity measure, then one has to decide whether to use a probabilistic or a set-theoretic similarity measure.

We first briefly discuss the use of indirect similarity measures. As pointed out by Schneider and Borlund (2007a), from a statistical perspective, the use of an indirect similarity measure is a quite unconventional approach.⁷ However, despite being unconventional, we do not believe that the approach has any fundamental statistical problems.⁸ Appropriate indirect similarity measures include the Bhattacharyya distance, the cosine,⁹ and the Jensen-Shannon distance. These measures are known to have good theoretical properties (Van Eck & Waltman, 2008). A very popular indirect similarity measure, especially for author cocitation analysis (e.g., McCain, 1990; White & Griffith, 1981; White & McCain, 1998), is the Pearson correlation. However, this measure does not have good theoretical properties and should therefore not be used (Ahlgren et al., 2003; Van Eck & Waltman, 2008). The chi-squared distance, which is proposed as an indirect similarity measure by Ahlgren et al. (2003), also does not have all the theoretical properties that we believe an appropriate indirect similarity measure should have (Van Eck & Waltman, 2008). We note that theoretical studies of indirect similarity measures can also be found in the psychometric literature (e.g., Zegers & Ten Berge, 1985). In this literature, the cosine is referred to as Tucker's congruence coefficient.

The notions of direct and indirect similarity are fundamentally different. Direct and indirect similarity measures may therefore lead to significantly different results (e.g., Schneider et al., 2009). In general, we believe the notion of direct similarity to be closer to the intuitive idea of similarity. Consider two objects that do not cooccur at all but that have quite similar cooccurrence profiles. The direct similarity between the objects will be very low, while the indirect similarity between the objects will be quite high. However, a high similarity between two objects that do not cooccur can

⁷A similar approach is sometimes taken in psychological research (e.g., Rosenberg & Jones, 1972; Rosenberg, Nelson, & Vivekananthan, 1968). In the psychological literature, there is some discussion about the advantages and disadvantages of this approach (Drasgow & Jones, 1979; Simmen, 1996; Van der Kloot & Van Herk, 1991).

⁸One of the issues that is sometimes raised is how the diagonal of a cooccurrence matrix should be treated. From a theoretical point of view, there are in our opinion two satisfactory solutions. One solution is to treat diagonal elements as missing values. The other solution is to set diagonal elements equal to the number of times objects occur at least twice in the same document (see also Ahlgren et al., 2003).

⁹There are two different similarity measures, a direct and an indirect one, that are both referred to as the cosine. Here we mean the cosine as discussed by, for example, Ahlgren et al. (2003) and Van Eck and Waltman (2008). This is a different measure than the one defined in Equation 7.

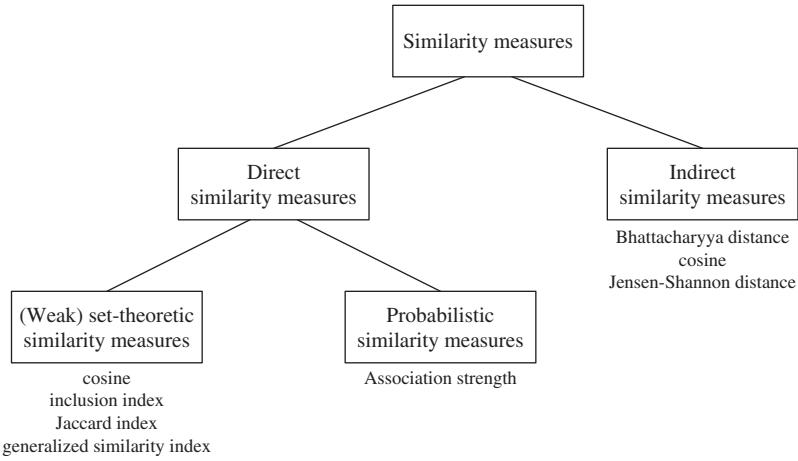


FIG. 3. Different types of similarity measures.

be rather counterintuitive, at least in certain contexts. For that reason, we believe that in general the notion of direct similarity is more natural than the notion of indirect similarity. We note, however, that indirect similarity measures may also have an advantage over direct similarity measures. Compared with direct similarity measures, indirect similarity measures are calculated based on a larger amount of data and most likely they therefore involve less statistical uncertainty.

In the rest of this section, we focus our attention on the use of direct similarity measures. Direct similarity measures determine the similarity between two objects by taking the number of cooccurrences of the objects and adjusting this number for the total number of occurrences of each of the objects. In scientometric research, when a direct similarity measure is applied to cooccurrence data, the aim usually is to normalize the data, that is, to correct the data for differences in the number of occurrences of objects. This brings us to the main question of this article: How should cooccurrence data be normalized? Or, in other words, which direct similarity measures are appropriate for normalizing cooccurrence data and which are not? We argue that cooccurrence data should always be normalized using a probabilistic similarity measure. Other direct similarity measures are not appropriate for normalization purposes. In particular, set-theoretic similarity measures should not be used to normalize cooccurrence data.

To see why probabilistic similarity measures have the right properties for normalizing cooccurrence data, recall that the number of cooccurrences of two objects can be seen as the result of two independent effects, the similarity effect and the size effect. As we discussed earlier in this article, probabilistic similarity measures correct for the size effect. This follows from Property 13. Set-theoretic similarity measures do not have this property, and they therefore do not properly correct for the size effect. As a consequence, set-theoretic similarity measures have, on average, higher values for objects that occur more frequently (see also Luukkonen et al., 1993; Zitt et al., 2000). The values of probabilistic similarity measures, on the other hand, do not depend on how frequently objects occur. This difference between

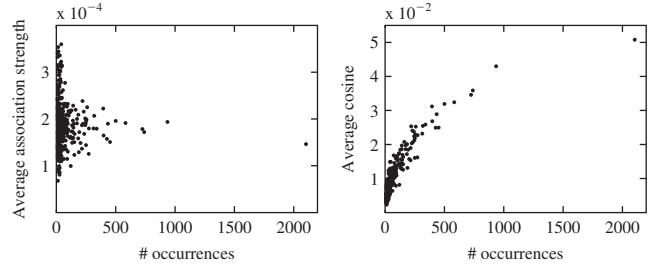


FIG. 4. Relation between on the one hand the number of occurrences of a term and on the other hand the average similarity of a term with other terms. In the left panel, similarities are determined using the association strength. In the right panel, similarities are determined using the cosine.

set-theoretic and probabilistic similarity measures can easily be demonstrated empirically. In Figure 4, this is done for the term data set discussed in the previous section. (The author data set and the journal data set yield similar results.) The figure shows the relation between, on the one hand, the number of occurrences of a term and, on the other hand, the average similarity of a term with other terms. In the left panel of the figure, similarities are determined using a probabilistic similarity measure, namely the association strength. In this panel, there is no substantial correlation between the number of occurrences of a term and the average similarity of a term ($r = -0.069$, $\rho = -0.029$). This is very different in the right panel, in which similarities are determined using a set-theoretic similarity measure, namely the cosine. (The inclusion index and the Jaccard index yield similar results.) In the right panel, there is a strong positive correlation between the number of occurrences of a term and the average similarity of a term ($r = 0.839$, $\rho = 0.882$). Results such as those shown in the right panel clearly indicate that set-theoretic similarity measures do not properly correct for the size effect and, consequently, do not properly normalize cooccurrence data. It follows from this observation that one should be very careful with the interpretation of similarities that have been derived from cooccurrence data using a set-theoretic

similarity measure (see also Luukonen et al., 1993; Zitt et al., 2000). Moreover, when such similarities are analyzed using multivariate analysis techniques such as multidimensional scaling or hierarchical clustering, one should pay special attention to possible artifacts in the results of the analysis. When using multidimensional scaling, for example, it is our experience that frequently occurring objects tend to cluster together in the center of a solution.

To provide some additional insight why probabilistic similarity measures are more appropriate for normalization purposes than set-theoretic similarity measures, we now compare the main ideas underlying these two types of measures. Suppose that we are performing a co-word analysis and that we want to determine the similarity between two words, word i and word j . We consider two hypothetical scenarios to which we refer as scenario 1 and scenario 2. The scenarios are summarized in Table 7, and they are illustrated graphically in the left and right panels of Figure 5. In each panel of the figure, the rectangle represents the set of all documents used in the co-word analysis, the dark gray circle represents the set of all documents in which word i occurs, and the striped circle represents the set of all documents in which word j occurs. The area of a rectangle or circle is proportional to the number of documents in the corresponding set.

As can be seen in Table 7 and Figure 5, in scenario 1 words i and j both occur quite frequently, while in scenario 2, they both occur relatively infrequently. In both scenarios, however, the relative overlap of the set of documents in which word i occurs and the set of documents in which word j occurs is the same. That is, in both scenarios word i occurs in 30% of the documents in which word j occurs and, the other way around, word j occurs in 30% of the documents in which word i occurs. Because the relative overlap is the

same, set-theoretic similarity measures, such as the cosine, the inclusion index, and the Jaccard index, yield the same similarity between words i and j in both scenarios (see Table 7). This is a consequence of Property 2. At first sight, it might seem a natural result to have the same similarity between words i and j in both scenarios. However, we argue that this result is far from natural, at least for normalization purposes.

We first consider scenario 1 in more detail. In this scenario, words i and j each occur in 30% of all documents. If there is no special relation between words i and j and if, as a consequence, occurrences of the two words are statistically independent, one would expect the two words to cooccur in approximately $30\% \times 30\% = 9\%$ of all documents. As can be seen in Table 7, words i and j cooccur in exactly 9% of all documents. Hence, occurrences of words i and j seem to be statistically independent, at least approximately, and there seems to be no strong relation between the two words.

We now consider scenario 2. In this scenario, words i and j each occur in 2% of all documents. If occurrences of words i and j are statistically independent, one would expect the two words to cooccur in approximately 0.04% of all documents. However, words i and j cooccur in 0.6% of all documents, that is, they cooccur 15 times more frequently than would be expected under the assumption of statistical independence. Hence, there seems to be a quite strong relation between words i and j , definitely much stronger than in scenario 1.

It is clear that set-theoretic similarity measures yield results that do not properly reflect the difference between scenario 1 and scenario 2. This is because set-theoretic similarity measures are based on the idea of measuring the relative overlap of sets instead of the idea of measuring the deviation from statistical independence. Probabilistic similarity measures, such as the association strength, are based on the latter idea, and they therefore yield results that do properly reflect the difference between scenario 1 and scenario 2. As can be seen in Table 7, the association strength indicates that in scenario 2 the similarity between words i and j is 15 times higher than in scenario 1. This reflects that in scenario 2 the cooccurrence frequency of words i and j is 15 times higher than would be expected under the assumption of statistical independence, while in scenario 1, the cooccurrence frequency of the two words equals the expected cooccurrence frequency under the independence assumption.

Conclusions

We have studied the application of direct similarity measures to cooccurrence data. Our survey of the scientometric literature has indicated that the most popular direct similarity measures are the association strength, the cosine, the inclusion index, and the Jaccard index. We have therefore focused most of our attention on these four measures. To make a well-considered decision which measure is most appropriate for one's purposes, we believe it to be indispensable to have a good theoretical understanding of the properties of the various measures. In this article, we have analyzed these

TABLE 7. Summary of two hypothetical scenarios in a co-word analysis.

	Scenario 1	Scenario 2
m	1 000	1 000
s_i	300	20
s_j	300	20
c_{ij}	90	6
Association strength	.001	.015
Cosine	.300	.300
Inclusion index	.300	.300
Jaccard index	.176	.176

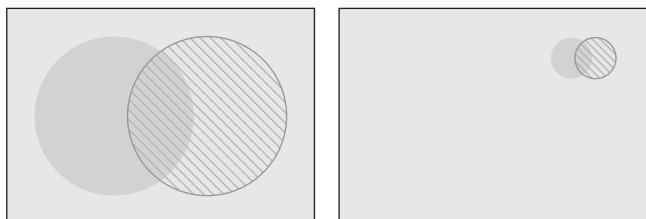


FIG. 5. Graphical illustration of two hypothetical scenarios in a co-word analysis. Scenario 1 is shown in the left panel; scenario 2 is shown in the right panel.

properties in considerable detail. Our analysis has revealed that there are two fundamentally different types of direct similarity measures. On the one hand, there are set-theoretic similarity measures, which can be interpreted as measures of the relative overlap of two sets. On the other hand, there are probabilistic similarity measures, which can be interpreted as measures of the deviation of observed cooccurrence frequencies from expected cooccurrence frequencies under an independence assumption. The cosine, the inclusion index, and the Jaccard index are examples of set-theoretic similarity measures, while the association strength is an example of a probabilistic similarity measure. Set-theoretic and probabilistic similarity measures serve different purposes, and it therefore makes no sense to argue that one measure is always better than another. In scientometric research, however, similarity measures are usually used for normalization purposes, and we have argued that in that specific case probabilistic similarity measures are much more appropriate than set-theoretic ones. Consequently, for most applications of direct similarity measures in scientometric research, we advise against the use of set-theoretic similarity measures and we recommend the use of a probabilistic similarity measure.

In addition to our theoretical analysis, we have also performed an empirical analysis of the behavior of various direct similarity measures. The analysis has shown that in practical applications the differences between various direct similarity measures can be quite large. This indicates that the issue of choosing an appropriate similarity measure is not only of theoretical interest but also has a high practical relevance. Another empirical observation that we have made is that set-theoretic similarity measures yield systematically higher values for frequently occurring objects than for objects that occur only a limited number of times. This confirms our theoretical finding that set-theoretic similarity measures do not properly correct for size effects. Probabilistic similarity measures do not have this problem.

There is one final comment that we would like to make. Above, we have argued in favor of the use of probabilistic similarity measures in scientometric research. Since probabilistic similarity measures are all proportional to each other, it does not really matter which probabilistic similarity measure one uses. In this article, we have focused most of our attention on one particular probabilistic similarity measure, namely the association strength defined in Equation 6. This measure shares with many other direct similarity measures the property that it takes values between zero and one. For practical purposes, however, it may be convenient not to use the measure in Equation 6 directly but instead to multiply this measure by the number of documents m (e.g., Van Eck & Waltman, 2007; Van Eck et al., 2006). This results in a slight variant of the association strength. We have pointed out that this variant has the appealing property that it equals one if the observed cooccurrence frequency of two objects equals the cooccurrence frequency that would be expected under the assumption that occurrences of the objects are statistically independent. It takes a value above or below one if the observed cooccurrence frequency is, respectively, higher or

lower than the expected cooccurrence frequency under the independence assumption.

References

- Ahlgren, P., Jarneving, B., & Rousseau, R. (2003). Requirements for a cocitation similarity measure, with special reference to Pearson's correlation coefficient. *Journal of the American Society for Information Science and Technology*, 54(6), 550–560.
- Anderberg, M.R. (1973). *Cluster analysis for applications*. New York: Academic Press.
- Baulieu, F.B. (1989). A classification of presence/absence based dissimilarity coefficients. *Journal of Classification*, 6(1), 233–246.
- Baulieu, F.B. (1997). Two variant axiom systems for presence/absence based dissimilarity coefficients. *Journal of Classification*, 14(1), 159–170.
- Borg, I., & Groenen, P.J.F. (2005). *Modern multidimensional scaling* (2nd ed.). Springer.
- Boyack, K.W., Klavans, R., & Börner, K. (2005). Mapping the backbone of science. *Scientometrics*, 64(3), 351–374.
- Braam, R.R., Moed, H.F., & Van Raan, A.F.J. (1991a). Mapping of science by combined co-citation and word analysis. I. Structural aspects. *Journal of the American Society for Information Science*, 42(4), 233–251.
- Braam, R.R., Moed, H.F., & Van Raan, A.F.J. (1991b). Mapping of science by combined Cocitation and word analysis. II. Dynamical aspects. *Journal of the American Society for Information Science*, 42(4), 252–266.
- Callon, M., Courtial, J.P., & Laville, F. (1991). Co-word analysis as a tool for describing the network of interactions between basic and technological research: The case of polymer chemistry. *Scientometrics*, 22(1), 155–205.
- Chung, Y.M., & Lee, J.Y. (2001). A corpus-based approach to comparative evaluation of statistical term association measures. *Journal of the American Society for Information Science and Technology*, 52(4), 283–296.
- Church, K.W., & Hanks, P. (1990). Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1), 22–29.
- Cox, M.A.A., & Cox, T.F. (2008). Multidimensional scaling. In C. Chen, W. Härdle, & A. Unwin (Eds.), *Handbook of data visualization* (pp. 315–347). Springer.
- Cox, T.F., & Cox, M.A.A. (2001). *Multidimensional scaling* (2nd ed.). Chapman & Hall/CRC.
- De Solla Price, D. (1981). The analysis of scientometric matrices for policy implications. *Scientometrics*, 3(1), 47–53.
- Drasgow, F., & Jones, L.E. (1979). Multidimensional scaling of derived dissimilarities. *Multivariate Behavioral Research*, 14(2), 227–244.
- Egghe, L. (2009). New relations between similarity measures for vectors based on vector norms. *Journal of the American Society for Information Science and Technology*, 60(2), 232–239.
- Egghe, L., & Michel, C. (2002). Strong similarity measures for ordered sets of documents in information retrieval. *Information Processing and Management*, 38(6), 823–848.
- Egghe, L., & Michel, C. (2003). Construction of weak and strong similarity measures for ordered sets of documents using fuzzy set techniques. *Information Processing and Management*, 39(5), 771–807.
- Egghe, L., & Rousseau, R. (2006). Classical retrieval and overlap measures satisfy the requirements for rankings based on a Lorenz curve. *Information Processing and Management*, 42(1), 106–120.
- Glänzel, W. (2001). National characteristics in international scientific co-authorship relations. *Scientometrics*, 51(1), 69–115.
- Glänzel, W., Schubert, A., & Czerwon, H.-J. (1999). A bibliometric analysis of international scientific cooperation of the European Union (1985–1995). *Scientometrics*, 45(2), 185–202.
- Gmür, M. (2003). Co-citation analysis and the search for invisible colleges: A methodological evaluation. *Scientometrics*, 57(1), 27–57.
- Gower, J.C. (1985). Measures of similarity, dissimilarity, and distance. In S. Kotz & N.L. Johnson (Eds.), *Encyclopedia of statistical sciences* (Vol. 5, pp. 397–405). Wiley.
- Gower, J.C., & Legendre, P. (1986). Metric and Euclidean properties of dissimilarity coefficients. *Journal of Classification*, 3(1), 5–48.

- Guilford, J.P. (1973). *Fundamental statistics in psychology and education* (5th ed.). McGraw-Hill.
- Hamers, L., Hemeryck, Y., Herweyers, G., Janssen, M., Keters, H., Rousseau, R., et al. (1989). Similarity measures in scientometric research: The Jaccard index versus Salton's cosine formula. *Information Processing and Management*, 25(3), 315–318.
- Hardy, G.H., Littlewood, J.E., & Pólya, G. (1952). *Inequalities* (2nd ed.). Cambridge University Press.
- Heimeriks, G., Hörlesberger, M., & Van den Besselaar, P. (2003). Mapping communication and collaboration in heterogeneous research networks. *Scientometrics*, 58(2), 391–413.
- Hinze, S. (1994). Bibliographical cartography of an emerging interdisciplinary discipline: The case of bioelectronics. *Scientometrics*, 29(3), 353–376.
- Hubálek, Z. (1982). Coefficients of association and similarity, based on binary (presence-absence) data: An evaluation. *Biological Reviews*, 57(4), 669–689.
- Janson, S., & Vegelius, J. (1981). Measures of ecological association. *Oecologia*, 49(3), 371–376.
- Jarneving, B. (2008). A variation of the calculation of the first author cocitation strength in author cocitation analysis. *Scientometrics*, 77(3), 485–504.
- Jones, W.P., & Furnas, G.W. (1987). Pictures of relevance: A geometric analysis of similarity measures. *Journal of the American Society for Information Science*, 38(6), 420–442.
- Klavans, R., & Boyack, K.W. (2006a). Identifying a better measure of relatedness for mapping science. *Journal of the American Society for Information Science and Technology*, 57(2), 251–263.
- Klavans, R., & Boyack, K.W. (2006b). Quantitative evaluation of large maps of science. *Scientometrics*, 68(3), 475–499.
- Kopcsa, A., & Schiebel, E. (1998). Science and technology mapping: A new iteration model for representing multidimensional relationships. *Journal of the American Society for Information Science*, 49(1), 7–17.
- Kostoff, R.N., Del Rio, J.A., Humenik, J.A., García, E.O., & Ramírez, A.M. (2001). Citation mining: Integrating text mining and bibliometrics for research user profiling. *Journal of the American Society for Information Science and Technology*, 52(13), 1148–1156.
- Kostoff, R.N., Eberhart, H.J., & Toothman, D.R. (1999). Hypersonic and supersonic flow roadmaps using bibliometrics and database tomography. *Journal of the American Society for Information Science*, 50(5), 427–447.
- Law, J., & Whittaker, J. (1992). Mapping acidification research: A test of the co-word method. *Scientometrics*, 23(3), 417–461.
- Leclerc, M., & Gagné, J. (1994). International scientific cooperation: The continentalization of science. *Scientometrics*, 31(3), 261–292.
- Leydesdorff, L. (1989). Words and co-words as indicators of intellectual organization. *Research Policy*, 18(4), 209–223.
- Leydesdorff, L. (2008). On the normalization and visualization of author Cocitation data: Salton's cosine versus the Jaccard index. *Journal of the American Society for Information Science and Technology*, 59(1), 77–85.
- Leydesdorff, L., & Vaughan, L. (2006). Cooccurrence matrices and their applications in information science: Extending ACA to the Web environment. *Journal of the American Society for Information Science and Technology*, 57(12), 1616–1628.
- Luukkonen, T., Persson, O., & Sivertsen, G. (1992). Understanding patterns of international scientific collaboration. *Science, Technology, and Human Values*, 17(1), 101–126.
- Luukkonen, T., Tijssen, R.J.W., Persson, O., & Sivertsen, G. (1993). The measurement of international scientific collaboration. *Scientometrics*, 28(1), 15–36.
- Manning, C.D., & Schütze, H. (1999). *Foundations of statistical natural language processing*. MIT Press.
- McCain, K.W. (1990). Mapping authors in intellectual space: A technical overview. *Journal of the American Society for Information Science*, 41(6), 433–443.
- McCain, K.W. (1991). Mapping economics through the journal literature: An experiment in journal cocitation analysis. *Journal of the American Society for Information Science*, 42(4), 290–296.
- McCain, K.W. (1995). The structure of biotechnology R & D. *Scientometrics*, 32(2), 153–175.
- Morillo, F., Bordons, M., & Gómez, I. (2003). Interdisciplinarity in science: A tentative typology of disciplines and research areas. *Journal of the American Society for Information Science and Technology*, 54(13), 1237–1249.
- Palmer, C.L. (1999). Structures and strategies of interdisciplinary science. *Journal of the American Society for Information Science*, 50(3), 242–253.
- Peters, H.P.F., Braam, R.R., & Van Raan, A.F.J. (1995). Cognitive resemblance and citation relations in chemical engineering publications. *Journal of the American Society for Information Science*, 46(1), 9–21.
- Peters, H.P.F., & Van Raan, A.F.J. (1993a). Co-word-based science maps of chemical engineering. Part I: Representations by direct multidimensional scaling. *Research Policy*, 22(1), 23–45.
- Peters, H.P.F., & Van Raan, A.F.J. (1993b). Co-word-based science maps of chemical engineering. Part II: Representations by combined clustering and multidimensional scaling. *Research Policy*, 22(1), 47–71.
- Qin, J. (2000). Semantic similarities between a keyword database and a controlled vocabulary database: An investigation in the antibiotic resistance literature. *Journal of the American Society for Information Science*, 51(3), 166–180.
- Rip, A., & Courtial, J.-P. (1984). Co-word maps of biotechnology: An example of cognitive scientometrics. *Scientometrics*, 6(6), 381–400.
- Rorvig, M. (1999). Images of similarity: A visual exploration of optimal similarity metrics and scaling properties of TREC topic-document sets. *Journal of the American Society for Information Science*, 50(8), 639–651.
- Rosenberg, S., & Jones, R. (1972). A method for investigating and representing a person's implicit theory of personality: Theodore Dreiser's view of people. *Journal of Personality and Social Psychology*, 22(3), 372–386.
- Rosenberg, S., Nelson, C., & Vivekananthan, P.S. (1968). A multidimensional approach to the structure of personality impressions. *Journal of Personality and Social Psychology*, 9(4), 283–294.
- Salton, G. (1963). Associative document retrieval techniques using bibliographic information. *Journal of the ACM*, 10(4), 440–457.
- Salton, G., & McGill, M.J. (1983). *Introduction to modern information retrieval*. McGraw-Hill.
- Schneider, J.W., & Borlund, P. (2007a). Matrix comparison, part 1: Motivation and important issues for measuring the resemblance between proximity measures or ordination results. *Journal of the American Society for Information Science and Technology*, 58(11), 1586–1595.
- Schneider, J.W., & Borlund, P. (2007b). Matrix comparison, part 2: Measuring the resemblance between proximity measures or ordination results by use of the Mantel and Procrustes statistics. *Journal of the American Society for Information Science and Technology*, 58(11), 1596–1609.
- Schneider, J.W., Larsen, B., & Ingwersen, P. (2009). A comparative study of first and all-author co-citation counting, and two different matrix generation approaches applied for author Cocitation analyses. *Scientometrics*, 80(1), 105–132.
- Schubert, A., & Braun, T. (1990). International collaboration in the sciences, 1981–1985. *Scientometrics*, 19(1–2), 3–10.
- Simmen, M.W. (1996). Multidimensional scaling of binary dissimilarities: Direct and derived approaches. *Multivariate Behavioral Research*, 31(1), 47–67.
- Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24(4), 265–269.
- Small, H. (1981). The relationship of information science to the social sciences: A co-citation analysis. *Information Processing and Management*, 17(1), 39–50.
- Small, H. (1994). A SCI-Map case study: Building a map of AIDS research. *Scientometrics*, 30(1), 229–241.
- Small, H., & Greenlee, E. (1980). Citation context analysis of a co-citation cluster: Recombinant-DNA. *Scientometrics*, 2(4), 277–301.

- Small, H., & Sweeney, E. (1985). Clustering the science citation index using co-citations. I. A comparison of methods. *Scientometrics*, 7(3–6), 391–409.
- Small, H., Sweeney, E., & Greenlee, E. (1985). Clustering the science citation index using co-citations. II. Mapping science. *Scientometrics*, 8(5–6), 321–340.
- Sokal, R.R., & Sneath, P.H.A. (1963). *Principles of numerical taxonomy*. San Francisco: Freeman.
- Tijssen, R.J.W. (1992). A quantitative assessment of interdisciplinary structures in science and technology: Co-classification analysis of energy research. *Research Policy*, 21(1), 27–44.
- Tijssen, R.J.W. (1993). A scientometric cognitive study of neural network research: Expert mental maps versus bibliometric maps. *Scientometrics*, 28(1), 111–136.
- Tijssen, R.J.W., & Van Raan, A.F.J. (1989). Mapping co-word structures: A comparison of multidimensional scaling and LEXIMAPPE. *Scientometrics*, 15(3–4), 283–295.
- Van der Kloot, W.A., & Van Herk, H. (1991). Multidimensional scaling of sorting data: A comparison of three procedures. *Multivariate Behavioral Research*, 26(4), 563–581.
- Van Eck, N.J., & Waltman, L. (2007). Bibliometric mapping of the computational intelligence field. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 15(5), 625–645.
- Van Eck, N.J., & Waltman, L. (2008). Appropriate similarity measures for author co-citation analysis. *Journal of the American Society for Information Science and Technology*, 59(10), 1653–1661.
- Van Eck, N.J., Waltman, L., Van den Berg, J., & Kaymak, U. (2006). Visualizing the computational intelligence field. *IEEE Computational Intelligence Magazine*, 1(4), 6–10.
- Van Raan, A.F.J., & Tijssen, R.J.W. (1993). The neural net of neural network research: An exercise in bibliometric mapping. *Scientometrics*, 26(1), 169–192.
- Vaughan, L. (2006). Visualizing linguistic and cultural differences using Web co-link data. *Journal of the American Society for Information Science and Technology*, 57(9), 1178–1193.
- Vaughan, L., & You, J. (2006). Comparing business competition positions based on Web co-link data: The global market vs. the Chinese market. *Scientometrics*, 68(3), 611–628.
- Waltman, L., & Van Eck, N.J. (2007). Some comments on the question whether cooccurrence data should be normalized. *Journal of the American Society for Information Science and Technology*, 58(11), 1701–1703.
- Warrens, M.J. (2008). Similarity coefficients for binary data. Doctoral dissertation, Leiden University.
- White, H.D. (2003). Pathfinder networks and author cocitation analysis: A remapping of paradigmatic information scientists. *Journal of the American Society for Information Science and Technology*, 54(5), 423–434.
- White, H.D., & Griffith, B.C. (1981). Author co-citation: A literature measure of intellectual structure. *Journal of the American Society for Information Science*, 32(3), 163–171.
- White, H.D., & McCain, K.W. (1998). Visualizing a discipline: An author co-citation analysis of information science, 1972–1995. *Journal of the American Society for Information Science*, 49(4), 327–355.
- Zegers, F.E., & Ten Berge, J.M.F. (1985). A family of association coefficients for metric scales. *Psychometrika*, 50(1), 17–24.
- Zitt, M., Bassecoulard, E., & Okubo, Y. (2000). Shadows of the past in international cooperation: Collaboration profiles of the top five producers of science. *Scientometrics*, 47(3), 627–657.

Appendix

In this appendix, we prove the theoretical results presented in the article.

Proof of Proposition 1.

We prove each property separately.
 (Property 5) This property follows from Property 3. Property 3 implies that, if $c_{ij} > 0$, $S(c_{ij}, s_i, s_j) > S(c_{ij}, s_i + 1, s_j)$. Hence, if $c_{ij} > 0$, $S(c_{ij}, s_i, s_j)$ cannot take its minimum value. This means that $S(c_{ij}, s_i, s_j)$ can take its minimum value only if $c_{ij} = 0$. This proves Property 5.

(Property 6) This property follows from Properties 1, 2, and 3. Suppose that $c_{ij} = s_i = s_j$. For all $(c'_{ij}, s'_i, s'_j) \in D_S$ such that $c'_{ij} = 0$, Property 1 implies that $S(c'_{ij}, s'_i, s'_j) \leq S(c_{ij}, s_i, s_j)$. For all $(c'_{ij}, s'_i, s'_j) \in D_S$ such that $c'_{ij} > 0$, Property 3 implies that $S(c'_{ij}, s'_i, s'_j) \leq S(c'_{ij}, c'_{ij}, c'_{ij})$ and Property 2 implies that $S(c'_{ij}, c'_{ij}, c'_{ij}) = S(c_{ij}, s_i, s_j)$. Hence, for all $(c'_{ij}, s'_i, s'_j) \in D_S$, $S(c'_{ij}, s'_i, s'_j) \leq S(c_{ij}, s_i, s_j)$. This means that, if $c_{ij} = s_i = s_j$, $S(c_{ij}, s_i, s_j)$ takes its maximum value. This proves Property 6.

(Property 7) This property follows from Properties 1, 3, and 5. Properties 1 and 5 imply that, if $c_{ij} = 0$, $S(c_{ij}, s_i, s_j)$ cannot take its maximum value. Property 3 implies that, if $0 < c_{ij} < s_i$ or $0 < c_{ij} < s_j$, $S(c_{ij}, s_i, s_j) < S(c_{ij}, c_{ij}, c_{ij})$. Hence, if $0 < c_{ij} < s_i$ or $0 < c_{ij} < s_j$, $S(c_{ij}, s_i, s_j)$ cannot take its maximum value. It now follows that $S(c_{ij}, s_i, s_j)$ can take its maximum value only if $c_{ij} = s_i = s_j$. This proves Property 7.

(Property 8) This property follows from Properties 1, 2, 3, and 5. If $c_{ij} = 0$, the property follows trivially from Properties 1 and 5. We, therefore, focus on the case in which $c_{ij} > 0$. Suppose, without loss of generality, that $0 < c_{ij} < s_i$. Consider an arbitrary constant $\alpha > 0$, and let $\beta = (c_{ij} + \alpha)/c_{ij}$. Property 2 implies that $S(\beta c_{ij}, \beta s_i, \beta s_j) = S(c_{ij}, s_i, s_j)$. Moreover, because $\beta c_{ij} = c_{ij} + \alpha$, $\beta s_i > s_i + \alpha$, and $\beta s_j \geq s_j + \alpha$, Property 3 implies that $S(\beta c_{ij}, \beta s_i, \beta s_j) < S(c_{ij} + \alpha, s_i + \alpha, s_j + \alpha)$. It now follows that $S(c_{ij} + \alpha, s_i + \alpha, s_j + \alpha) > S(c_{ij}, s_i, s_j)$. This proves Property 8.

Proof of Proposition 2. Let $S(c_{ij}, s_i, s_j)$ denote an arbitrary set-theoretic similarity measure that has Property 9. We start by showing that for all $(c_{ij}, s_i, s_j), (c'_{ij}, s'_i, s'_j) \in D_S$ the properties of set-theoretic similarity measures together with Property 9 are sufficient to determine whether $S(c_{ij}, s_i, s_j)$ is greater than, less than, or equal to $S(c'_{ij}, s'_i, s'_j)$. Suppose first that $c_{ij}, c'_{ij} > 0$. Let $\alpha = c_{ij}/c'_{ij}$. Property 2 implies that $S(\alpha c'_{ij}, \alpha s'_i, \alpha s'_j) = S(c'_{ij}, s'_i, s'_j)$. Moreover, taking into account that $c_{ij} = \alpha c'_{ij}$, it can be seen that Property 9 determines whether $S(c_{ij}, s_i, s_j)$ is greater than, less than, or equal to $S(\alpha c'_{ij}, \alpha s'_i, \alpha s'_j)$. Hence, if $c_{ij}, c'_{ij} > 0$, Properties 2 and 9 are sufficient to determine whether $S(c_{ij}, s_i, s_j)$ is greater than, less than, or equal to $S(c'_{ij}, s'_i, s'_j)$. Suppose next that $c_{ij} = 0$ or $c'_{ij} = 0$. Property 1 implies that $S(c_{ij}, s_i, s_j) = S(c'_{ij}, s'_i, s'_j)$ if $c_{ij} = c'_{ij} = 0$. Furthermore, Properties 1 and 5 imply that $S(c_{ij}, s_i, s_j) > S(c'_{ij}, s'_i, s'_j)$ if $c_{ij} > c'_{ij} = 0$ and, conversely, that $S(c_{ij}, s_i, s_j) < S(c'_{ij}, s'_i, s'_j)$ if $c'_{ij} > c_{ij} = 0$. Hence, if $c_{ij} = 0$ or

$c'_{ij} = 0$, Properties 1 and 5 are sufficient to determine whether $S(c_{ij}, s_i, s_j)$ is greater than, less than, or equal to $S(c'_{ij}, s'_i, s'_j)$. It now follows that for all $(c_{ij}, s_i, s_j), (c'_{ij}, s'_i, s'_j) \in D_S$ the properties of set-theoretic similarity measures together with Property 9 are sufficient to determine whether $S(c_{ij}, s_i, s_j)$ is greater than, less than, or equal to $S(c'_{ij}, s'_i, s'_j)$. This implies that all set-theoretic similarity measures that have Property 9 are monotonically related to each other. One of these measures is the cosine defined in Equation 7. Hence, all set-theoretic similarity measures that have Property 9 are monotonically related to the cosine. This completes the proof of the proposition.

Proof of Proposition 3. The proof is analogous to the proof of Proposition 2 provided above.

Proof of Proposition 4. Let $S(c_{ij}, s_i, s_j)$ denote an arbitrary weak set-theoretic similarity measure that has Property 11. Property 11 implies that, if $c_{ij} > 0$, $S(c_{ij}, s_i, s_j) > S(c_{ij}, s_i + 1, s_j + 1)$. Hence, if $c_{ij} > 0$, $S(c_{ij}, s_i, s_j)$ cannot take its minimum value. This means that $S(c_{ij}, s_i, s_j)$ can take its minimum value only if $c_{ij} = 0$. In other words, $S(c_{ij}, s_i, s_j)$ has Property 5. This shows that all weak set-theoretic similarity measures $S(c_{ij}, s_i, s_j)$ that have Property 11 also have Property 5. The rest of the proof is now analogous to the proof of Proposition 2 provided above.

Proof of Proposition 5. It is easy to see that for all finite values of the parameter p the generalized similarity index defined in Equation 11 has Properties 1, 2, and 3. Hence, it follows from Definition 3 that for all finite values of the parameter p the generalized similarity index is a set-theoretic similarity measure. This completes the proof of the proposition.

Proof of Proposition 6. Let $S(c_{ij}, s_i, s_j)$ denote an arbitrary probabilistic similarity measure. Furthermore, let $c'_{ij} = c_{ij}/(s_i s_j)$ for all $i \neq j$, and let $s'_i = 1$ for all i . It follows from Property 13 that $S(c_{ij}, s_i, s_j) = S(c'_{ij}, s'_i, s'_j)$ for all $i \neq j$, and it follows from Property 12 that $S(c'_{ij}, s'_i, s'_j) = \alpha c'_{ij}$ for all $i \neq j$ and for some $\alpha > 0$. Hence, for all $i \neq j$ and for some $\alpha > 0$, $S(c_{ij}, s_i, s_j) = S(c'_{ij}, s'_i, s'_j) = \alpha c'_{ij} = \alpha c_{ij}/(s_i s_j) = \alpha S_A(c_{ij}, s_i, s_j)$. In other words, $S(c_{ij}, s_i, s_j)$ is proportional to the association strength defined in Equation 6. This completes the proof of the proposition.

Proof of Corollary 7. The association strength defined in Equation 6 does not have Property 2 and is therefore not a (weak or non-weak) set-theoretic similarity measure. The same is true for all measures that are proportional to the association strength. Consequently, it follows from Proposition 6 that a probabilistic similarity measure cannot also be a (weak or non-weak) set-theoretic similarity measure. This completes the proof of the corollary.