

# Distributional Measures as Proxies for Semantic Distance: A Survey

Saif Mohammad  
University of Toronto

Graeme Hirst  
University of Toronto

*The ability to mimic human notions of semantic distance has widespread applications. Some measures rely only on raw text (distributional measures) and some rely on knowledge sources such as WordNet. Although extensive studies have been performed to compare WordNet-based measures with human judgment, the use of distributional measures as proxies to estimate semantic distance has received little attention. Even though they have traditionally performed poorly when compared to WordNet-based measures, they lay claim to certain uniquely attractive features, such as their applicability in resource-poor languages and their ability to mimic both semantic similarity and semantic relatedness. Therefore, this paper presents a detailed study of distributional measures. Particular attention is paid to flesh out the strengths and limitations of both WordNet-based and distributional measures, and how distributional measures of distance can be brought more inline with human notions of semantic distance. We conclude with a brief discussion of recent work on hybrid measures.*

## 1. Introduction

**Semantic distance** is a measure of how close or distant the meanings of two units of language are. The units of language may be words, phrases, sentences, paragraphs, or documents. The two nouns *dance* and *choreography*, for example, are closer in meaning than the two nouns *clown* and *bridge*, and so are said to be semantically closer. The semantic distances between words (or more precisely, between concepts) can be used as fundamental building blocks for measuring semantic distance between larger units of language. The ability to mimic human judgments of semantic distance is useful in numerous natural language tasks including machine translation, word sense disambiguation, thesaurus creation, information retrieval, text summarization, automated spelling correction, and identifying discourse structure. This paper describes the state-of-the-art in corpus-based measures of semantic distance between these fundamental units of language. It identifies the significant challenges that existing approaches to semantic distance face and in the process fleshes out questions that lead to a better understanding of why two concepts are considered semantically close. The paper concludes with a discussion of new hybrid approaches, that show the potential to address these challenges.

Units of language, especially words, may have more than one possible meaning. However, their context may be used to determine the intended senses. For example, *star* can mean both CELESTIAL BODY and CELEBRITY; however, *star* in the sentence below refers only to CELESTIAL BODY and is much closer to *sun* than to *famous*:

- (1) *Stars are powered by nuclear fusion.*

Thus, semantic distance between words in context is in fact the distance between their underlying senses or concepts.

Humans consider two concepts to be semantically close if there is a sharing of some meaning. Specifically, two concepts are semantically close if there is a **lexical semantic relation** between the concepts. Putting it differently, the reason why two concepts are considered semantically close can be attributed to a lexical semantic relation that binds them. According to Cruse (1986), a lexical semantic relation is a relation between **lexical units**—a surface form along with a sense. As he points out, the number of semantic relations that bind concepts is innumerable; but certain relations, such as hyponymy, meronymy, antonymy, and troponymy, are more systematic and have enjoyed more attention in the linguistics community. However, as Morris and Hirst (2004) point out, these relations are far out-numbered by others, which they call **non-classical relations**. Here are a few of the kinds of non-classical relations they observed: positive qualities (BRILLIANT, KIND), concepts pertaining to a concept (KIND, CHIVALROUS, FORMAL pertaining to GENTLEMANLY), and commonly co-occurring words (locations such as HOMELESS, SHELTER; problem–solution pairs such as HOMELESS, SHELTER).

### 1.1 Semantic relatedness and semantic similarity

Semantic distance is of two kinds: **semantic similarity** and **semantic relatedness**. The former is a subset of the latter, but the two terms may be used interchangeably in certain contexts, making it even more important to be aware of their distinction. Two concepts are considered to be semantically similar if there is a synonymy (or near-synonymy), hyponymy (hypernymy), antonymy, or troponymy relation between them (examples include APPLES–BANANAS, DOCTOR–SURGEON, DARK–BRIGHT). Two concepts are considered to be semantically related if there is any lexical semantic relation at all between them—classical or non-classical (examples include APPLES–BANANAS, SURGEON–SCALPEL, TREE–SHADE).

### 1.2 Human judgments of semantic distance

Humans are adept at estimating semantic distance; but consider the following questions: How strongly will two people agree/disagree on distance estimates? Will the agreement vary over different sets of concepts? Are we equally good at estimating semantic similarity and semantic relatedness? In our minds, is there a clear distinction between related and unrelated concepts or are concept-pairs spread across the whole range from synonymous to unrelated? Some of the earliest work that begins to address these questions is by Rubenstein and Goodenough (1965). They conducted quantitative experiments with human subjects (51 in all) who were asked to rate 65 English word pairs on a scale from 0.0 to 4.0 as per their semantic distance. The word pairs chosen ranged from almost synonymous to unrelated. However, they were all noun pairs and those that were semantically close were also semantically similar; the dataset did not contain word pairs that are semantically related but not semantically similar. The subjects repeated the annotation after two weeks and the new distance values had a Pearson's correlation  $r$  of 0.85 with the old ones. Miller and Charles (1991) also conducted a similar study on 30 word pairs taken from the Rubenstein–Goodenough pairs. These annotations had a high correlation ( $r = 0.97$ ) with the mean annotations of Rubenstein and Goodenough (1965). Resnik (1999) repeated these experiments and found the inter-annotator correlation ( $r$ ) to be 0.90.

**Table 1**

Different datasets that are manually annotated with distance values. Pearson’s correlation coefficient ( $r$ ) was used to determine inter-annotator correlation (last column).

Dataset	Year	Language	# pairs	PoS	# subjects	$r$
Rubenstein and Goodenough	1965	English	65	N	51	-
Miller and Charles	1991	English	30	N	38	.90
Resnik and Diab	2000	English	27	V	10	.76 and .79
Gurevych	2005	German	65	N	24	.81
Zesch and Gurevych	2006	German	350	N, V, A	8	.69

Note: Rubenstein and Goodenough (1965) do not report inter-subject correlation, but determine intra-subject correlation to be 0.85 for 36 (out of the 65) word pairs for which similarity judgments were made again by 15 (of the 51) subjects after a span of two weeks.

Resnik and Diab (2000) conducted annotations of 48 verb pairs and found inter-annotator correlation ( $r$ ) to be 0.76 (when the verbs were presented without context) and 0.79 (when presented in context). Gurevych (2005) and Zesch, Gurevych, and Mühlhäuser (2007) asked native German speakers to mark two different sets of German word pairs with distance values. Set 1 was a German translation of the Rubenstein and Goodenough (1965) dataset. It had 65 noun–noun word pairs. Set 2 was a larger dataset containing 350 word pairs made up of nouns, verbs, and adjectives. The semantically close word pairs in the 65-word set were mostly synonyms or hypernyms (hyponyms) of each other, whereas those in the 350-word set had both classical and non-classical relations with each other. Details of these **semantic distance benchmarks** are summarized in Table 1. Inter-subject correlations (last column in Table 1) are indicative of the degree of ease in annotating the datasets.

It should be noted here that even though the annotators were presented with word-pairs and not concept-pairs, it is reasonable to assume that they were annotated as per their closest senses. For example, given the noun pair *bank* and *interest*, most if not all will identify it as semantically related even though both words have more than one sense and many of the sense–sense combinations are unrelated (for example, the RIVER BANK sense of *bank* and the SPECIAL ATTENTION sense of *interest*). The high agreement and correlation values suggest that humans are quite good and consistent at estimating semantic distance of noun-pairs; however, annotating verbs and adjectives and a combination of parts of speech is harder. This also means that estimating semantic relatedness is harder than estimating semantic similarity.

Apart from showing that humans can indeed estimate semantic distance, these datasets act as “gold standards” to evaluate automatic distance measures. However, lack of large amounts of data from human subject experimentation limits the reliability of this mode of evaluation. Therefore automatic distance measures are also evaluated by their usefulness in natural language tasks such as those mentioned earlier.

### 1.3 Automatic measures of semantic distance

Automatic measures of semantic distance quantify the semantic distance between word pairs. They give values within a certain range (for example, 0 to 1), such that one end of this range represent maximally close or synonymous, while the other end represents maximally distant. Depending on what the two ends represent, measures of semantic

distance can be classified as **measures of distance** (larger values indicate greater distance and less close) and **measures of closeness** (larger values indicate shorter distance and more close).<sup>1</sup> A measure of closeness can be easily converted to a measure of distance by applying a suitable inverse function, such as  $(1 - x)^{-1}$  or  $e^{-(1-x)}$ , and the other way round.

Two classes of automatic methods have been traditionally used to determine semantic distance. **Knowledge-rich measures of concept-distance**, such as those of Jiang and Conrath (1997), Leacock and Chodorow (1998), and Resnik (1995), rely on the structure of a knowledge source, such as WordNet, to determine the distance between two concepts defined in it.<sup>2</sup> **Distributional measures of word-distance (knowledge-lean measures)**, such as cosine and  $\alpha$ -skew divergence (Lee 2001), rely on the **distributional hypothesis**, which states that two words tend to be semantically close if they occur in similar contexts (Firth 1957). Distributional measures rely simply on text (and possibly some shallow syntactic processing) and can give the distance between any two words that occur at least a few times.

The various WordNet-based measures have been widely studied (Budanitsky and Hirst 2006; Patwardhan, Banerjee, and Pedersen 2003). The study of distributional measures on the whole has received much less attention.<sup>3</sup> Even though, as Weeds (2003) and Mohammad and Hirst (2006b) show, they perform poorly when compared to WordNet-based measures, the distributional measures of word-distance have many attractive features, including their ability to measure both semantic similarity and semantic relatedness. Further, they are not dependent on costly knowledge sources that do not exist for most languages. This paper therefore focuses on distributional measures and analyzes their strengths and limitations. Particular attention is paid to the different kinds of distributional measures and their components. The motivation is that a better understanding of distributional measures will lead to bringing them more inline with human notions of semantic distance, while still maintaining their applicability to resource-poor languages and their ability to mimic both semantic similarity and semantic distance.

## 2. Knowledge-rich approaches to semantic distance

Before we begin our examination of distributional measures, we look briefly at the resource-based measures. In some ways they are complementary to distributional measures and so the discussion will set the context for the analysis of distributional measures.

Creation of electronically available ontologies and semantic networks like WordNet has allowed their use to help solve numerous natural language problems including the measurement of semantic distance. Budanitsky and Hirst (2006), Hirst and Budanitsky (2005), and Patwardhan, Banerjee, and Pedersen (2003) have done an extensive survey of the various WordNet-based measures, their comparisons with human judgment on

---

1 A note about terminology: In many contexts, the term *distance measures* refers to the complete set of measures (irrespective of what the different ends of its range signify). In certain other contexts (as in this paragraph), *distance measures* refers only to those measures that give larger values to signify greater distance. The context, usually by its reference to this numeric property or lack thereof, make clear the intended meaning of the term.

2 The nodes in WordNet (synsets) represent concepts and edges between nodes represent semantic relations such as hyponymy and meronymy.

3 See Curran (2004), Weeds, Weir, and McCarthy (2004) for other work that compares various distributional measures.

selected word pairs, and their usefulness in applications such as real-word spelling correction and word sense disambiguation. Hence, this section provides only a brief summary of the major knowledge-rich measures of semantic distance.

## 2.1 Measures that exploit WordNet’s semantic network

A number of WordNet-based measures consider two concepts to be close if they are close to each other in WordNet. One of the earliest and simplest measures is Rada et al.’s (1989) **edge-counting** method. The shortest path in the network between the two target concepts (**target path**) is determined. The more edges there are between two words, the more distant they are. Elegant as it may be, the measure hinges on the largely incorrect assumption that all the network edges correspond to identical semantic distance.

Nodes in a network may be connected by different kinds of lexical relations such as hyponymy, meronymy, and so on. Edge counts apart, Hirst and St-Onge’s (1998) measure takes into account the fact that if the target path consists of edges that belong to many different relations, then the target concepts are likely more distant. The idea is that if we start from a particular node  $c_1$  and take a path via a particular relation (say, hyponymy), to a certain extent the concepts reached will be semantically related to  $c_1$ . However, if during the way we take edges belonging to different relations (other than hyponymy), very soon we may reach words that are unrelated. Hirst and St-Onge’s measure of semantic relatedness is:

$$HS(c_1, c_2) = C - path\ length - k \times d \tag{1}$$

where  $c_1$  and  $c_2$  are the target concepts,  $d$  is the number of times an edge pertaining to a relation different from that of the preceding edge is taken, and  $C$  and  $k$  are empirically determined constants. More recently, Yang and Powers (2005) proposed a weighted edge-counting method to determine semantic relatedness using the hypernymy/hyponymy, holonymy/meronymy, and antonymy links in WordNet.

Leacock and Chodorow (1998) used just one relation (hyponymy) and modified the path length formula to reflect the fact that edges lower down in the *is-a* hierarchy correspond to smaller semantic distance than the ones higher up. For example, synsets pertaining to *sports car* and *car* (low in the hierarchy) are much more similar than those pertaining to *transport* and *instrumentation* (higher up in the hierarchy) even though both pairs of nodes are separated by exactly one edge in WordNet’s *is-a* hierarchy. Their formula is:

$$LC(c_1, c_2) = -\log \frac{len(c_1, c_2)}{2D} \tag{2}$$

where  $D$  is the maximum depth of the taxonomy.

Resnik (1995) suggested a measure that combines corpus statistics with WordNet. He proposed that since the **lowest common subsumer** or **lowest superordinate (lso)** of the target nodes represents what is similar between them, the semantic similarity between the two concepts is directly proportional to how specific the lso is. The more general the lso is, the larger the semantic distance between the target nodes. This specificity is measured by the formula for information content (IC)—the negative logarithm of the probability of the lso:

$$Res(c_1, c_2) = IC(lso(c_1, c_2)) = -\log P(lso(c_1, c_2)) \tag{3}$$

Observe that using information content has the effect of inherently scaling the semantic similarity measure by the depth of the taxonomy. Usually, the lower the lowest super-ordinate, the lower the probability of occurrence of the lso and the concepts subsumed by it, and hence, the higher its information content.

As per Resnik's formula, given a particular lowest super-ordinate, the exact positions of the target nodes below it in the hierarchy do not have any effect on the semantic similarity. Intuitively, we would expect that word pairs closer to the lso are more semantically similar than those that are distant. Jiang and Conrath (1997) and Lin (1997) incorporate this notion into their measures which are arithmetic variations of the same terms. The Jiang and Conrath (1997) measure (*JC*) determines how dissimilar each target concept is from the lso ( $IC(c_1) - IC(lso(c_1, c_2))$  and  $IC(c_2) - IC(lso(c_1, c_2))$ ). The final semantic distance between the two concepts is then taken to be the sum of these differences. Lin (1997) (like Resnik) points out that the lso is what is common between the two target concepts and that its information content is the common information between the two concepts. His formula (*Lin*) can be thought of as taking the Dice coefficient of the information in the two target concepts.

$$JC(c_1, c_2) = 2 \log p(lso(c_1, c_2)) - (\log(p(c_1)) + (\log(p(c_2)))) \quad (4)$$

$$Lin(c_1, c_2) = \frac{2 \times \log p(lso(c_1, c_2))}{\log(p(c_1)) + (\log(p(c_2)))} \quad (5)$$

Budanitsky and Hirst (2006) showed that the Jiang-Conrath measure has the highest correlation (0.850) with the Miller and Charles noun pairs and performs better than all these measures in a spelling correction task. Patwardhan, Banerjee, and Pedersen (2003) achieved similar results using the measure for word sense disambiguation.

All of the approaches described above rely heavily (if not solely) on the hypernymy/hyponymy network in WordNet; they are designed for, and evaluated on, noun-noun pairs. However, more recently, Resnik and Diab (2000) and Yang and Powers (2006) developed measures aimed at verb-verb pairs. Resnik and Diab (2000) ported several measures which are traditionally applied on the noun hypernymy/hyponymy network (edge counting, Resnik (1995), and Lin (1997)) to the relatively shallow verb troponymy network. The two information content-based measures ranked a carefully chosen set of 48 verbs best in order of their semantic distance.<sup>4</sup> Yang and Powers (2006) ported their earlier work on nouns (Yang and Powers 2005) to verbs. In order to compensate for the relatively shallow verb troponymy hierarchy and the lack of a corresponding holonymy/meronymy hierarchy, they proposed several back-off models—the most useful one being the distance between a noun pair that has the same lexical form as the verb pair. However, the approach has too many tuned parameters (9 in all) and performed poorly on a set of 36 TOEFL word-choice questions involving verb targets and alternatives.

## 2.2 Measures that rely on dictionaries and thesauri

Lesk (1986) introduced a method to perform word sense disambiguation using word glosses (definitions). The glosses of the senses of a target word are compared with those of its context and the number of word overlaps is determined. The sense with

---

<sup>4</sup> Only those verbs were selected which require a theme, and the sub-categorization frames of verb pairs had to match.

the most number of overlaps is chosen as the intended sense of the target. Inspired by this approach, Banerjee and Pedersen (2003) proposed a semantic relatedness measure that deems two concepts to be more semantically related if there is more overlap in their glosses. Notably, they overcome the problem of short glosses by considering the glosses of concepts related to the target concepts through the WordNet lexical semantic relations such as hyponymy/hypernymy. They also give more weight to larger overlap sequences. Patwardhan and Pedersen (2006) proposed another gloss-based semantic relatedness measure which performed slightly worse than the extended gloss overlap measure in a word sense disambiguation task, but markedly better at ranking the Miller and Charles (1991) word pairs. They create **aggregate co-occurrence vectors** for a WordNet sense by adding the co-occurrence vectors of the words in its WordNet gloss. The distance between two senses is then determined by the cosine of the angle between their aggregate vectors. Such aggregate co-occurrence vectors are expected to be noisy because they are created from data that is not sense-annotated.

Jarmasz and Szpakowicz (2003) use the taxonomic structure of *Roget's Thesaurus* to determine semantic similarity. Two words are considered maximally similar if they occur in the same semicolon group in the thesaurus. Then on, decreasing in similarity are word pairs in the same paragraph, words pairs in different paragraphs belonging to the same part of speech and within the same category, word pairs in the category, and so on until word pairs which have nothing in common except that they are in the thesaurus (maximally distant). They show that this simple approach performs remarkably well at ranking word pairs and determining the correct answer in sets of TOEFL, ESL, and *Reader's Digest* word choice problems.

## 2.3 Challenges

In this section, we review some of the shortcomings of resource-based measures, in order to motivate and to compare them with distributional measures that we will introduce in Section 3.

**2.3.1 Lack of high-quality WordNet-like knowledge sources.** Ontologies, WordNets, and semantic networks are available for a few languages such as English, German, and Hindi. Creating them requires human experts and it is time intensive. Thus, for most languages, we cannot use WordNet-based measures simply due to the lack of a WordNet in that language. Further, even if created, updating an ontology is again expensive and there is usually a lag between the current state of language usage/comprehension and the semantic network representing it. Further, the complexity of human languages makes creation of even a near-perfect semantic network of its concepts impossible. Thus in many ways the ontology-based measures are only as good as the networks on which they are based.

On the other hand, distributional measures require only text. Large corpora, billions of words in size, may now be collected by a simple web crawler. Large corpora of more-formal writing are also available (for example, the *Wall Street Journal* or the *American Printing House for the Blind (APHB)* corpus). This makes distributional measures very attractive.

**2.3.2 Poor estimation of semantic relatedness.** As Morris and Hirst (2004) pointed out, a large number of concept pairs, such as STRAWBERRY-CREAM and DOCTOR-SCALPEL, have a non-classical relation between them (STRAWBERRIES are usually eaten with CREAM and a DOCTOR uses a SCALPEL to make an incision). These words are not

semantically similar, but rather semantically related. An ontology- or WordNet-based measure will correctly identify the amount of semantic relatedness only if such relations are explicitly coded into the knowledge source. Further, the most accurate WordNet-based measures rely only on its extensive is-a hierarchy. This is because networks of other lexical-relations such as meronymy are much less developed. Further, the networks for different parts of speech are not well connected. All this means that, while WordNet-based measures accurately estimate semantic similarity between nouns, their estimation of semantic relatedness, especially in pairs other than noun–noun, is at best poor and at worse non-existent. On the other hand, distributional measures can be used to determine both semantic relatedness and semantic similarity.

**2.3.3 Inability to cater to specific domains.** Given a concept pair, measures that rely only on WordNet and no text, such as that of Rada et al. (1989), give just one distance value. However, two concepts may be very close in a certain domain but not so much in another. For example, *SPACE* and *TIME* are close in the domain of quantum mechanics but not so much in most others. Ontologies have been made for specific domains, which may be used to determine semantic similarity specific to these domains. However, the number of such ontologies is very limited. Some of the more successful WordNet-based measures, such as Jiang and Conrath (1997), that rely also on text, do indeed capture domain-specificity to some extent, but the distance values are still largely shaped by the underlying network, which is not domain-specific. On the other hand, distributional measures rely primarily (if not completely) on text, and large amounts of corpora specific to particular domains can easily be collected.

**2.3.4 Computational complexity and storage requirements.** As applications for linguistic distance become more sophisticated and demanding, it becomes attractive to pre-compute and store the distance values between all possible pairs of words or senses. However both WordNet-based and distributional measures have large space requirements to do this, requiring matrices of size  $N \times N$ , where  $N$  is very large. In case of WordNet-based measures,  $N$  is the number of senses (81,000 just for nouns). In case of distributional measures,  $N$  is the size of the vocabulary (at least 100,000 for most languages). Given that the above matrices tend to be sparse<sup>5</sup> and that computational capabilities are continuing to improve, the above limitation may not seem hugely problematic, but as we see more and more natural language applications in embedded systems and hand-held devices, such as cell phones, iPods, and medical equipment, memory and computational power become serious constraints.

**2.3.5 Reluctance to cross the language barrier.** Both WordNet-based and distributional measures have largely been used in a monolingual framework. Even though semantic distance seems to hold promise in tasks such as machine translation and multi-lingual text summarization that inherently involve two or more languages, automatic measures of semantic distance have rarely been applied. With the development of the EuroWordNet, involving interconnected networks of seven different languages, it is possible that we shall see more cross-lingual work using WordNet-based measures in the future. However, such an interconnected network will be very hard to create for more-different language pairs such as English and Chinese or English and Arabic.

---

<sup>5</sup> Even though WordNet-based and distributional measures give non-zero closeness values to a large number of term pairs, values below a suitable threshold can be reset to 0.

### 3. Knowledge-lean, distributional approaches to semantic distance

#### 3.1 The distributional hypotheses: the original and the revised

**Distributional measures** are inspired by the maxim “You shall know a word by the company it keeps” (Firth 1957). These measures rely simply on raw text and possibly some shallow syntactic processing. They are much less resource-hungry than the semantic measures, but they measure the distance between words rather than word-senses or concepts. Two words are considered close if they occur in similar contexts. The context of a target word is usually taken to be the set of words within a certain window around it, for example,  $\pm 5$  words or the complete sentence. The set of contexts of a target word is usually represented by the set of words in these contexts, their strength of association (SoA) with the target word, and possibly their syntactic relation with the target, for example verb–object, subject–verb, and so on. The strength of co-occurrence association between the target and another word quantifies how much more (or less) than chance the two words occur together in text. Commonly used measures of association are conditional probability (CP) and pointwise mutual information (PMI). The distance between the sets of contexts of two target words can be used as a proxy for their semantic distance as words found in similar contexts tend to be semantically similar—the **distributional hypothesis** (Firth 1957; Harris 1968).

The hypothesis makes intuitive sense, as Budanitsky and Hirst (2006) point out: If two words have many co-occurring words in common, then similar things are being said about both of them and so they are likely to be semantically similar. Conversely, if two words are semantically similar, then they are likely to be used in a similar fashion in text and thus end up with many common co-occurrences. For example, the semantically similar *bug* and *insect* are expected to have a number of common co-occurring words such as *crawl*, *squash*, *small*, *woods*, and so on, in a large-enough text corpus.

The distributional hypothesis only mentions semantic similarity and not semantic relatedness. This, coupled with the fact that the difference between semantic relatedness and semantic similarity is somewhat nuanced and can be missed, meant that almost all work employing the distributional hypothesis was labeled as estimating semantic similarity. However, it should be noted that distributional measures can be used to estimate both semantic similarity and semantic relatedness. Even though Schütze and Pedersen (1997) and Landauer, Foltz, and Laham (1998), for example, use the term *similarity* and not *relatedness*, their LSA-based distance measures in fact estimate semantic relatedness and not semantic similarity. We propose more-specific distributional hypotheses that make clear how distributional measures can be used to estimate semantic similarity and how they can be used to measure semantic relatedness:

**Hypothesis of the distributionally close and semantically related:**

Two target words are distributionally close and semantically related if they have many common strongly co-occurring words.

(For example, *doctor–surgeon* and *doctor–scalpel*. See example co-occurring words in Table 2.)

**Hypothesis of the distributionally close and semantically similar:**

Two target words are distributionally close and semantically similar if they have many common strongly co-occurring words that each have the same syntactic relation with the two targets.

(For example, *doctor–surgeon*, but not *doctor–scalpel*. See syntactic relations with example co-occurring words in Table 2.)

**Table 2**

Example: Common syntactic relations of target words with co-occurring words.

	Co-occurring words		
	<i>cut</i> (v)	<i>hardworking</i> (adj)	<i>patient</i> (n)
<b>Semantically similar</b>			
<b>target pair</b>			
<i>doctor</i> (n)	subject-verb	noun-qualifier	subject-object
<i>surgeon</i> (n)	subject-verb	noun-qualifier	subject-object
<b>Semantically related</b>			
<b>target pair</b>			
<i>doctor</i> (n)	subject-verb	noun-qualifier	subject-object
<i>scalpel</i> (n)	prepositional object-verb	–	prepositional object-object

The idea is that both semantically similar and semantically related word pairs will have many common co-occurring words. However, words that are semantically similar belong to the same broad part of speech (noun, verb, etc.), but the same need not be true for words that are semantically related. Therefore, words that are semantically similar will tend to have the same syntactic relation, such as verb-object or subject-verb, with most common co-occurring words. Thus, the two words are considered semantically related simply if they have many common co-occurring words. But to be semantically similar as well, the words must have the same syntactic relation with co-occurring words. Consider the word pair *doctor-operate*. In a large enough body of text, the two words are likely to have the following common co-occurring words: *patient, scalpel, surgery, recuperate*, and so on. All these words will contribute to a high score of relatedness. However, they do not have the same syntactic relation with the two targets. (The word *doctor* is almost always used as a noun while *operate* is a verb.) Thus, as per the two revised distributional hypotheses, *doctor* and *operate* will correctly be identified as semantically related but not semantically similar. The word pair *doctor-nurse*, on the other hand, will be identified as both semantically related and semantically similar.

In order to clearly differentiate from the distance as calculated by a WordNet-based semantic measure (described earlier in Section 2.1), the distance calculated by a corpus-based distributional measure will be referred to as **distributional distance**.

### 3.2 Corpus-based measures of distributional distance

We now describe specific distributional measures that rely on the distributional hypotheses; depending on which specific hypothesis they use, they mimic either semantic similarity or semantic relatedness.

**3.2.1 Spatial Metrics: Cos, L<sub>1</sub>, L<sub>2</sub>.** Consider a multidimensional space in which the number of dimensions is equal to the size of the vocabulary. A word  $w$  can be represented by a point in this space such that the component of  $\vec{w}$  in a dimension (corresponding to word  $x$ , say) is equal to the strength of association (SoA) of  $w$  with  $x$  ( $SoA(w, x)$ ). Thus, the vectors corresponding to two words are *close* together, and thereby get a low distributional distance score, if they share many co-occurring words and the co-occurring

words have more or less the same strength of association with the two target words. The distance between two vectors can be calculated in different ways as described below.

*Cosine.* The **cosine** method (denoted by **Cos**) is one of the earliest and most widely used distributional measures. Given two words  $w_1$  and  $w_2$ , the cosine measure calculates the cosine of the angle between  $\vec{w}_1$  and  $\vec{w}_2$ . If a large number of words co-occur with both  $w_1$  and  $w_2$ , then  $\vec{w}_1$  and  $\vec{w}_2$  will have a small angle between them and the cosine will be large; signifying a large relatedness/similarity between them. The cosine measure gives scores in the range from 0 (unrelated) to 1 (synonymous). So the higher the value, the less distant the target word-pair is.

$$\text{Cos}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cup C(w_2)} (P(w|w_1) \times P(w|w_2))}{\sqrt{\sum_{w \in C(w_1)} P(w|w_1)^2} \times \sqrt{\sum_{w \in C(w_2)} P(w|w_2)^2}} \quad (6)$$

where  $C(w)$  is the set of words that co-occur (within a certain window) with the word  $w$  in a corpus. In this instantiation of the cosine measure, conditional probability of the co-occurring words given the target words is used as the strength of association.

The cosine was used, among others, by Schütze and Pedersen (1997) and Yoshida, Yukawa, and Kuwabara (2003), who suggest methods of automatically generating distributional thesauri from text corpora. Schütze and Pedersen (1997) use the Tipster category B corpus (Harman 1993) (450,000 unique terms) and the *Wall Street Journal* to create a large but sparse co-occurrence matrix of 3,000 medium-frequency words (frequency rank between 2,000 and 5,000). Latent semantic indexing (singular value decomposition) (Schütze and Pedersen 1997) is used to reduce the dimensionality of the matrix and get for each term a word vector of its 20 strongest co-occurrences. The cosine of a target word’s vector with each of the other word vectors is calculated and the words that give the highest scores comprise the thesaurus entry for the target word.

Yoshida, Yukawa, and Kuwabara (2003) believe that words that are closely related for one person may be distant for another. They use around 40,000 HTML documents to generate personalized thesauri for six different people. Documents used to create the thesaurus for a person are retrieved from the subject’s home page and a web crawler which accesses linked documents. The authors also suggest a root-mean-squared method to determine the similarity of two different thesaurus entries for the same word.

*Manhattan and Euclidean Distances.* Distance between two points (words) in vector space can also be calculated using the formulae for **Manhattan distance** a.k.a. the **L<sub>1</sub> norm** (denoted by **L<sub>1</sub>**) or **Euclidean distance** a.k.a. the **L<sub>2</sub> norm** (denoted by **L<sub>2</sub>**). In the Manhattan distance (7) (Dagan, Lee, and Pereira (1997), Dagan, Lee, and Pereira (1999), and Lee (1999)), the difference in strength of association of  $w_1$  and  $w_2$  with each word that they co-occur with is summed. The greater the difference, the greater is the distributional distance between the two words. Euclidean distance (8) (Lee 1999) employs the root mean square of the difference in association to get the final distributional distance. Both the **L<sub>1</sub>** and **L<sub>2</sub>** norms give scores in the range between 0 (zero distance or synonymous) and infinity (maximally distant or unrelated).

$$L_1(w_1, w_2) = \sum_{w \in C(w_1) \cup C(w_2)} |P(w|w_1) - P(w|w_2)| \quad (7)$$

$$L_2(w_1, w_2) = \sqrt{\sum_{w \in C(w_1) \cup C(w_2)} (P(w|w_1) - P(w|w_2))^2} \quad (8)$$

The above formulae use conditional probability of the co-occurring words given a target word as the strength of association.

Lee (1999) compared the ability of all three spatial metrics to determine the probability of an unseen (not found in training data) word pair. The measures in order of their performance (from better to worse) were:  $L_1$  norm, cosine, and  $L_2$  norm. Weeds (2003) determined the correlation of word pair ranking as per a handful of distributional measures with human rankings (Miller and Charles (1991) word pairs). She used verb-object pairs from the *British National Corpus* (BNC) and found the correlation of  $L_1$  norm with human rankings to be 0.39.

**3.2.2 Mutual information-based measures: Hindle, Lin.** Hindle (1990) was one of the first to factor the strength of association of co-occurring words into a distributional similarity measure.<sup>6</sup> Consider the nouns  $n_j$  and  $n_k$  that exist as objects of verb  $v_i$  in different instances within a text corpus. Hindle used the following formula to determine the distributional similarity of  $n_j$  and  $n_k$  solely from their occurrences as object of  $v_i$ :

$$Hin_{obj}(v_i, n_j, n_k) = \begin{cases} \min(I(v_i, n_j), I(v_i, n_k)), & \text{if } I(v_i, n_j) > 0 \text{ and } I(v_i, n_k) > 0 \\ |\max(I(v_i, n_j), I(v_i, n_k))|, & \text{if } I(v_i, n_j) < 0 \text{ and } I(v_i, n_k) < 0 \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

$I(n, v)$  stands for the pointwise mutual information (PMI) between the noun  $n$  and verb  $v$  (note that in case of negative PMI values, the maximum function captures the PMI which is lower in absolute value). The measure follows from the distributional hypothesis—the more similar the associations of co-occurring words with the two target words, the more semantically similar they are. Hindle used PMI<sup>7</sup> as the strength of association. Using the minimum of the two PMIs captures the similarity in the strength of association of  $v_i$  with each of the two nouns.

Hindle used an analogous formula to calculate distributional similarity ( $Hin_{subj}$ ) using the subject-verb relation. The overall distributional similarity between any two nouns is calculated by the formula:

$$Hin(n_1, n_2) = \sum_{i=0}^N (Hin_{obj}(v_i, n_1, n_2) + Hin_{subj}(v_i, n_1, n_2)) \quad (10)$$

The measure gives similarity scores from 0 (maximally dissimilar) to infinity (maximally similar or synonymous). Note that in Hindle's measure, the set of co-occurring words used is restricted to include only those words that have the same syntactic relation with

<sup>6</sup> See Grefenstette (1992) for an approach that does not incorporate strength of association of co-occurring words. He, like Hindle (1990), uses syntactic dependencies to characterize the set of contexts of a target word. The Jaccard coefficient is used to determine how similar the two sets of contexts are.

<sup>7</sup> In their respective papers, Hindle (1990) and Lin (1998b) refer to pointwise mutual information as mutual information.

both target words (either verb–object or verb–subject). This is therefore a measure that mimics semantic similarity and not semantic relatedness. A form of Hindle’s measure where all co-occurring words are used, making it a measure that mimics semantic relatedness, is shown below:

$$Hin_{rel}(w_1, w_2) = \sum_{w \in C(w)} \begin{cases} \min(I(w, w_1), I(w, w_2)), & \text{if } I(w, w_1) > 0 \text{ and } I(w, w_2) > 0 \\ |\max(I(w, w_1), I(w, w_2))|, & \text{if } I(w, w_1) < 0 \text{ and } I(w, w_2) < 0 \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

where  $C(w)$  is the set of words that co-occur with word  $w$ .

Lin (1998b) suggests a different measure derived from his information-theoretic definition of similarity (Lin 1998a). Further, he uses a broad set of syntactic relations apart from just subject–verb and verb–object relations and shows that using multiple relations is beneficial even for Hindle’s measure. He first extracts triples of the form  $(x, r, y)$  from the partially parsed text, where the word  $x$  is related to  $y$  by the syntactic relation  $r$ . Lin defines the distributional similarity between two words,  $w_1$  and  $w_2$ , as follows:

$$Lin(w_1, w_2) = \frac{\sum_{(r,w) \in T(w_1) \cap T(w_2)} (I(w_1, r, w) + I(w_2, r, w))}{\sum_{(r,w') \in T(w_1)} I(w_1, r, w') + \sum_{(r,w'') \in T(w_2)} I(w_2, r, w'')} \quad (12)$$

where  $T(x)$  is the set of all word pairs  $(r, y)$  such that the pointwise mutual information  $I(x, r, y)$ , is positive. Note that this is different from Hindle (1990) where even the cases of negative PMI were considered. Church and Hanks (1990) showed that it is hard to accurately predict negative word association ratios with confidence, and so, co-occurrence pairs with negative PMI are ignored. The measure gives similarity scores from 0 (maximally dissimilar) to 1 (maximally similar).

Like Hindle’s measure, Lin’s is a measure of distributional *similarity*. However, it distinguishes itself from that of Hindle in two respects. First, Lin normalizes the similarity score between two words (numerator of (12)) by their cumulative strengths of association with the rest of the co-occurring words (denominator of (12)). This is a significant improvement as now high PMI of the target words with shared co-occurring words alone does not guarantee a high distributional similarity score. As an additional requirement, the target words must have low PMI with words they do not both co-occur with. Second, Hindle uses the minimum of the PMI between each of the target words and the shared co-occurring word, while Lin uses the sum. Taking the sum has the drawback of not penalizing for a mismatch in strength of co-occurrence, as long as  $w_1$  and  $w_2$  both co-occur with a word.

Hindle (1990) used a portion of the *Associated Press* news stories (6 million words) to classify the nouns into semantically related classes. Lin (1998b) used his measure to generate a distributional thesaurus from a 64-million-word corpus of the *Wall Street Journal*, *San Jose Mercury*, and *AP Newswire*. He also provides a framework for evaluating such automatically generated thesauri by comparing them with WordNet-based and Roget-based thesauri. He shows that the distributional thesaurus created with his measure is closer to the WordNet and Roget-based thesauri than that created using Hindle’s measure.

### 3.2.3 Relative Entropy-Based Measures: KLD, ASD, JSD.

*Kullback-Leibler divergence.* Given two probability mass functions  $p(x)$  and  $q(x)$ , their **relative entropy**  $D(p\|q)$  is:

$$D(p\|q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \quad \text{for } q(x) \neq 0 \quad (13)$$

Intuitively, if  $p(x)$  is the accurate probability mass function corresponding to a random variable  $X$ , then  $D(p\|q)$  is the information lost when approximating  $p(x)$  by  $q(x)$ . In other words,  $D(p\|q)$  is indicative of how different the two distributions are. Relative entropy is also called the **Kullback-Leibler divergence** or the **Kullback-Leibler distance** (denoted by **KLD**).

Pereira, Tishby, and Lee (1993) and Dagan, Lee, and Pereira (1994) point out that words have probabilistic distributions with respect to neighboring syntactically related words. For example, there exists a certain probabilistic distribution ( $d_1(P(v|n_1))$ ), say of a particular noun  $n_1$  being the object of any verb. This distribution can be estimated by corpus counts of parsed or chunked text. Let  $d_2(P(v|n_2))$  be the corresponding distribution for noun  $n_2$ . These distributions ( $d_1$  and  $d_2$ ) define the contexts of the two nouns ( $n_1$  and  $n_2$ , respectively). As per the distributional hypothesis, the more these contexts are similar, the more  $n_1$  and  $n_2$  are semantically similar. Thus the Kullback-Leibler distance between the two distributions is indicative of the semantic distance between the nouns  $n_1$  and  $n_2$ .

$$\begin{aligned} KLD(n_1, n_2) &= D(d_1\|d_2) \\ &= \sum_{v \in Vb} P(v|n_1) \log \frac{P(v|n_1)}{P(v|n_2)} \quad \text{for } P(v|n_2) \neq 0 \\ &= \sum_{v \in Vb'(n_1) \cap Vb'(n_2)} P(v|n_1) \log \frac{P(v|n_1)}{P(v|n_2)} \quad \text{for } P(v|n_2) \neq 0 \end{aligned} \quad (14)$$

where  $Vb$  is the set of all verbs and  $Vb'(x)$  is the set of verbs that have  $x$  as the object. Note again that the set of co-occurring words used is restricted to include only verbs that each have the same syntactic relation (verb-object) with both target nouns. This too is therefore a measure that mimics semantic similarity and not semantic relatedness.

It should be noted that the verb-object relationship is not inherent to the measure and that one or more of any other syntactic relations may be used. One may also estimate semantic relatedness by using all words co-occurring with the target words. Thus a more generic expression of the Kullback-Leibler divergence is as follows:

$$\begin{aligned} KLD(w_1, w_2) &= D(d_1\|d_2) \\ &= \sum_{w \in V} P(w|w_1) \log \frac{P(w|w_1)}{P(w|w_2)} \quad \text{for } P(w|w_2) \neq 0 \\ &= \sum_{w \in C(w_1) \cup C(w_2)} P(w|w_1) \log \frac{P(w|w_1)}{P(w|w_2)} \quad \text{for } P(w|w_2) \neq 0 \end{aligned} \quad (15)$$

where  $V$  is the vocabulary (all the words found in a corpus).  $C(w)$ , as mentioned earlier, is the set of words occurring (within a certain window) with word  $w$ .

It should be noted that the Kullback-Leibler distance is not symmetric; that is, the distance from  $w_1$  to  $w_2$  is not necessarily, and even not likely, the same as the distance from  $w_2$  to  $w_1$ . This asymmetry is counterintuitive to the general notion of semantic similarity of words, although Weeds (2003) has argued in favor of asymmetric measures.

Further, it is very likely that there are instances such that  $P(w_1|v)$  is greater than 0 for a particular verb  $v$ , while due to data sparseness or grammatical and semantic constraints, the training data has no sentence where  $v$  has the object  $w_2$ . This makes  $P(w_2|v)$  equal to 0 and the ratio of the two probabilities infinite. Kullback-Leibler divergence is not defined in such cases, but approximations may be made by considering smoothed values for the denominator.

Pereira, Tishby, and Lee (1993) used KLD to create clusters of nouns from verb-object pairs corresponding to the thousand most frequent nouns in the *Grolier's Encyclopedia*, June 1991 version (10 million words). Dagan, Lee, and Pereira (1994) used KLD to estimate the probabilities of bigrams that were not seen in a text corpus. They point out that a significant number of possible bigrams are not seen in any given text corpus. The probabilities of such bigrams may be determined by taking a weighted average of the probabilities of bigrams composed of distributionally similar words. Use of Kullback-Leibler distance as the semantic distance metric yielded a 20% improvement in perplexity on the *Wall Street Journal* and dictation corpora provided by ARPA's HLT program (Paul 1991).

It should be noted here that the use of distributionally similar words to estimate unseen bigram probabilities will likely lead to erroneous results in case of less-preferred and strongly-preferred collocations (word pairs). Inkpen and Hirst (2002) point out that even though words like *task* and *job* are semantically very similar, the collocations they form with other words may have varying degrees of usage. While *daunting task* is a strongly-preferred collocation, *daunting job* is rarely used. Thus using the probability of one bigram to estimate that of another will not be beneficial in such cases.

*$\alpha$ -skew divergence.* The  **$\alpha$ -skew divergence** (ASD) is a slight modification of the Kullback-Leibler divergence that obviates the need for smoothed probabilities. It has the following formula:

$$ASD(w_1, w_2) = \sum_{w \in C(w_1) \cup C(w_2)} P(w|w_1) \log \frac{P(w|w_1)}{\alpha P(w|w_2) + (1 - \alpha)P(w|w_1)} \quad (16)$$

where  $\alpha$  is a parameter that may be varied but is usually set to 0.99. Note that the denominator within the logarithm is never zero with a non-zero numerator. Also, the measure retains the asymmetric nature of the Kullback-Leibler divergence. Lee (2001) shows that  $\alpha$ -skew divergence performs better than Kullback-Leibler divergence in estimating word co-occurrence probabilities. Weeds (2003) achieves a correlation of 0.48 and 0.26 with human judgment on the Miller and Charles word pairs using  $ASD(w_1, w_2)$  and  $ASD(w_2, w_1)$ , respectively.

*Jensen-Shannon divergence.* A relative entropy-based measure that overcomes the problem of asymmetry in Kullback-Leibler divergence is the **Jensen-Shannon divergence** a.k.a. **total divergence to the average** a.k.a. **information radius**. It is denoted by **JSD** and has the following formula:

$$JSD(w_1, w_2) = D \left( d_1 \parallel \frac{1}{2}(d_1 + d_2) \right) + D \left( d_2 \parallel \frac{1}{2}(d_1 + d_2) \right) \quad (17)$$

$$= \sum_{w \in C(w_1) \cup C(w_2)} \left( P(w|w_1) \log \frac{P(w|w_1)}{\frac{1}{2}(P(w|w_1) + P(w|w_2))} + P(w|w_2) \log \frac{P(w|w_2)}{\frac{1}{2}(P(w|w_1) + P(w|w_2))} \right) \quad (18)$$

The Jensen-Shannon divergence is the sum of the Kullback-Leibler divergence between each of the individual co-occurrence distributions  $d_1$  and  $d_2$  of the target words with the average distribution ( $\frac{d_1+d_2}{2}$ ). Further, it can be shown that the Jensen-Shannon divergence avoids the problem of zero denominator. The Jensen-Shannon divergence is therefore always well defined and, like  $\alpha$ -skew divergence, obviates the need for smoothed estimates.

The Kullback-Leibler divergence,  $\alpha$ -skew divergence, and Jensen-Shannon divergence all give distributional distance scores from 0 (synonymous) to infinity (unrelated).

Landauer, Foltz, and Laham (1998) proposed **Latent semantic analysis (LSA)**, which can be used to determine distributional distance between words or between sets of words.<sup>8</sup> Unlike the various approaches described earlier where a word-word co-occurrence matrix is created, the first step of LSA involves the creation of a word-paragraph, word-document, or similar such word-passage matrix, where a *passage* is some grouping of words. A cell for word  $w$  and passage  $p$  is populated with the number of times  $w$  occurs in  $p$  or, for even better results, a function of this frequency that captures how much information the occurrence of the word in a text passage carries.

Next, the dimensionality of this matrix is reduced by applying **singular value decomposition (SVD)**, a standard matrix decomposition technique. This smaller set of dimensions represents abstract (unknown) concepts. Then the original word-passage matrix is recreated, but this time from the reduced dimensions. Landauer, Foltz, and Laham (1998) point out that this results in new matrix cell values that are different from what they were before. More specifically, words that are expected to occur more often in a passage than what the original cell values reflect are incremented. Then a standard vector distance measure, such as cosine, that captures the distance between distributions of the two target words is applied.

LSA was used by Schütze and Pedersen (1997), Turney (2001) and Rapp (2003) to measure distributional distance, with encouraging results. However, there is no non-heuristic way to determine when the dimension reduction should stop. Further, the generic concepts represented by the reduced dimensions are not interpretable; that is, one cannot determine which concepts they represent in a given sense inventory. This means that LSA cannot directly be used for tasks such as unsupervised sense disambiguation or estimating semantic similarity of known concepts. LSA is computationally expensive as singular value decomposition, a key component for dimensionality reduction, requires computationally intensive matrix operations. This makes LSA less scalable to large amounts of text (Gorman and Curran 2006). Finally, it too, like other distributional word-distance measures conflates the many senses of a word (see Section 4.6.1 ahead for more discussion on sense conflation).

---

<sup>8</sup> Landauer, Foltz, and Laham (1998) describe it as a measure of *similarity*, but in fact it is a distributional measure that mimics semantic relatedness.

**3.2.5 Co-occurrence Retrieval Models.** The distributional measures suggested by Weeds (2003) are based on a notion of substitutability: the more appropriate it is to substitute word  $w_1$  in place of word  $w_2$  in a suitable natural language task, the more semantically similar they are. She uses **co-occurrence retrieval** (the retrieval of words that co-occur with a target word from text) to determine the degree to which one word is substitutable by another. The degree of substitutability of  $w_2$  with  $w_1$  is dependent on how the proportion of co-occurrences of  $w_1$  that are also co-occurrences of  $w_2$  and the proportion of co-occurrences of  $w_2$  that are also co-occurrences of  $w_1$ . Thus Weeds’s distributional measures have a precision component and a recall component (which may or may not incorporate the strength of co-occurrence association). The final score is a weighted sum of the precision, recall, and standard  $F$  measure (see equation (19)<sup>9</sup>). The weights determine the importance of precision and recall and are determined empirically. If precision and recall are equally important, then it results in a symmetric measure which gives identical scores for the distributional similarity of  $w_1$  with  $w_2$  and  $w_2$  with  $w_1$ . Otherwise, we get an asymmetric measure which assigns different scores to the two cases.

$$CRM(w_1, w_2) = \gamma \left[ \frac{2 \times P \times R}{P + R} \right] + (1 - \gamma) \left[ \beta[P] + (1 - \beta)[R] \right] \quad (19)$$

$\gamma$  and  $\beta$  are tuned parameters that lie between 0 and 1.

Both precision and recall can be considered as the product of a core formula (denoted by *core*) and a penalty function (denoted by *penalty*). Weeds03 provides six (three times two) distinct formulae for precision and recall, depending on the strength of co-occurrence (three alternatives) and whether or not a penalty is applied for differences in strength of association of common co-occurring words (two alternatives).

Depending on the strength of association, the CRMs are classified as **type-based**, **token-based**, and **mutual information-based**. The CRMs that use simple counts of the common co-occurrences and not the strength of associations as core precision and recall values are called type-based CRMs (denoted by the superscript *type*). The CRMs that use conditional probabilities of the shared co-occurring words with the target words are called token-based CRMs (denoted by the superscript *token*). The CRMs that use pointwise mutual information of the shared co-occurring words with target words are called mutual information-based CRMs (denoted by the superscript *mi*). The core precision and recall formulae for type, token, and mutual information-based CRMs are listed below:

$$\text{core}_P^{\text{type}}(w_1, w_2) = \frac{|C(w_1) \cap C(w_2)|}{|C(w_1)|} \quad (20)$$

$$\text{core}_R^{\text{type}}(w_1, w_2) = \frac{|C(w_1) \cap C(w_2)|}{|C(w_2)|} \quad (21)$$

$$\text{core}_P^{\text{token}}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} P(w|w_1) \quad (22)$$

---

<sup>9</sup>  $P$  is short for  $P(w_1, w_2)$ , while  $R$  is short for  $R(w_1, w_2)$ . The abbreviations are made due to space constraints and to improve readability.

$$\text{core}_R^{\text{token}}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} P(w|w_2) \quad (23)$$

$$\text{core}_P^{\text{mi}}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cap C(w_2)} I(w, w_1)}{\sum_{w \in C(w_1)} I(w, w_1)} \quad (24)$$

$$\text{core}_R^{\text{mi}}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cap C(w_2)} I(w, w_2)}{\sum_{w \in C(w_2)} I(w, w_2)} \quad (25)$$

where  $C(x)$  is the set of words that co-occur with  $x$ .

Depending on the penalty function, the CRMs are classified as **additive** and **difference-weighted**. The CRMs that do not penalize difference in strength of co-occurrence are called additive CRMs (denoted by the subscript *add*); those that do penalize are called difference-weighted CRMs (subscript *d<sub>w</sub>*). The penalty is a conditional probability-based function (26, 27) for the token- and type-based CRMs, and a mutual information-based function (28, 29) for the mutual information-based CRM.

$$\text{penalty}_P^{\text{type}} = \text{penalty}_P^{\text{token}} = \frac{\min(P(w|w_1), P(w|w_2))}{P(w|w_1)} \quad (26)$$

$$\text{penalty}_R^{\text{type}} = \text{penalty}_R^{\text{token}} = \frac{\min(P(w|w_1), P(w|w_2))}{P(w|w_2)} \quad (27)$$

$$\text{penalty}_P^{\text{mi}} = \frac{\min(I(w, w_1), I(w, w_2))}{I(w, w_1)} \quad (28)$$

$$\text{penalty}_R^{\text{mi}} = \frac{\min(I(w, w_1), I(w, w_2))}{I(w, w_2)} \quad (29)$$

The six pairs of precision and recall difference-weighted CRMs are thus as follows:

$$P_{\text{add}}^{\text{type}}(w_1, w_2) = \frac{|C(w_1) \cap C(w_2)|}{|C(w_1)|} \quad (30)$$

$$R_{\text{add}}^{\text{type}}(w_1, w_2) = \frac{|C(w_1) \cap C(w_2)|}{|C(w_2)|} \quad (31)$$

$$P_{\text{d<sub>w</sub>}^{\text{type}}}(w_1, w_2) = \frac{\sum_{|C(w_1) \cap C(w_2)|} \frac{\min(P(w|w_1), P(w|w_2))}{P(w|w_1)}}{|C(w_1)|} \quad (32)$$

$$R_{\text{d<sub>w</sub>}^{\text{type}}}(w_1, w_2) = \frac{\sum_{|C(w_1) \cap C(w_2)|} \frac{\min(P(w|w_1), P(w|w_2))}{P(w|w_2)}}{|C(w_2)|} \quad (33)$$

$$P_{\text{add}}^{\text{token}}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} P(w|w_1) \quad (34)$$

$$R_{\text{add}}^{\text{token}}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} P(w|w_2) \quad (35)$$

$$P_{dw}^{token}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} \min(P(w|w_1), P(w|w_2)) \quad (36)$$

$$R_{dw}^{token}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} \min(P(w|w_2), P(w|w_1)) \quad (37)$$

$$P_{add}^{mi}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cap C(w_2)} I(w, w_1)}{\sum_{w \in C(w_1)} I(w, w_1)} \quad (38)$$

$$R_{add}^{mi}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cap C(w_2)} I(w, w_2)}{\sum_{w \in C(w_2)} I(w, w_2)} \quad (39)$$

$$P_{dw}^{mi}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cap C(w_2)} \min(I(w, w_1), I(w, w_2))}{\sum_{w \in C(w_1)} I(w, w_1)} \quad (40)$$

$$R_{dw}^{mi}(w_1, w_2) = \frac{\sum_{w \in C(w_1) \cap C(w_2)} \min(I(w, w_1), I(w, w_2))}{\sum_{w \in C(w_2)} I(w, w_2)} \quad (41)$$

Note that in case of the difference-weighted token and mutual information-based precision and recall formulae, there is a cancellation of a pair of terms obtained from the core formulae and the penalty.

Asymmetry in substitutability is intuitive, as in many cases it may be acceptable to substitute a word, say *dog*, with another, say *animal*, but the reverse is not acceptable as often. Since Weeds uses substitutability as a measure of semantic similarity, she believes that distributional similarity between two words should reflect this property as well. Hence, like the Kullback-Leibler divergence, all her distributional similarity models are asymmetric.

Weeds (2003) extracted verb-object pairs of 2,000 nouns from the *British National Corpus* (BNC). The verbs related to the target words by the verb-object relation were used. Thus each of the co-occurring verbs is related to the target nouns by the same syntactic relation and therefore the measures mimic semantic similarity, not relatedness. Correlation with human judgment (Miller and Charles word pairs) showed that difference-weighted ( $r = 0.61$ ) and additive mutual information-based measures ( $r = 0.62$ ) performed far better than the other CRMs.

#### 4. The anatomy of a distributional measure

Even though there are numerous distributional measures, many of which may seem dramatically different from each other, all distributional measures perform two functions: (1) create **distributional profiles (DPs)**, and (2) calculate the distance between two DPs.

The distributional profile of a word is the strength of association between it and each of the lexical, syntactic, and/or semantic units that co-occur with it. Commonly used **measures of strength of association** are conditional probability (0 to 1) and pointwise mutual information ( $-\infty$  to  $\infty$ ). Commonly used units of co-occurrence with the target are other *words*, and so we speak of the **lexical distributional profile of a word (lexical DPW)**. The co-occurring words may be all those in a predetermined window around the target, or may be restricted to those that have a certain syntactic (*e.g.*, verb-object) or semantic (*e.g.*, agent-theme) relation with the target word. We will refer to

**Table 3**

Measures of DP distance, measures of strength of association, and standard combinations. Measures of strength of association that are traditionally used are marked in bold. The use of other measures of association remains to be explored.

Measures of DP distance	Measures of strength of association
$\alpha$ -skew divergence (ASD)	$\phi$ coefficient (Phi)
cosine (Cos)	<b>conditional probability (CP)</b>
Dice coefficient (Dice)	cosine (Cos)
Euclidean distance ( $L_2$ norm)	Dice coefficient (Dice)
Hindle's measure (Hin)	odds ratio (Odds)
Kullback-Leibler divergence (KLD)	<b>pointwise mutual information (PMI)</b>
Manhattan distance ( $L_1$ norm)	Yule's coefficient (Yule)
Jensen-Shannon divergence (JSD)	
Lin's measure (Lin)	

Standard combinations
$\alpha$ -skew divergence— $\phi$ coefficient (ASD-CP)
cosine—conditional probability (Cos-CP)
Dice coefficient—conditional probability (Dice-CP)
Euclidean distance—conditional probability ( $L_2$ norm-CP)
Hindle's measure—pointwise mutual information (Hin-PMI)
Kullback-Leibler divergence—conditional probability (KLD-CP)
Manhattan distance—conditional probability ( $L_1$ norm-CP)
Jensen-Shannon divergence—conditional probability (JSD-CP)
Lin's measure—pointwise mutual information (Lin-PMI)

the former kind of DPs as **relation-free**. Usually in the latter case, separate association values are calculated for each of the different relations between the target and the co-occurring units. We will refer to such DPs as **relation-constrained**. Typical relation-free DPs are those of Schütze and Pedersen (1997) and Yoshida, Yukawa, and Kuwabara (2003). Typical relation-constrained DPs are those of Lin (1998a) and Lee (2001). Below are contrived, but plausible, examples of each for the word *pulse*; the numbers are conditional probabilities:

**relation-free DP**

*pulse*: *beat* .28, *racing* .2, *grow* .13, *beans* .09, *heart* .04, ...

**relation-constrained DP**

*pulse*: *langlebeat*, subject-*verbrangle* .34, *langleracing*, noun-qualifying *adjectiverangle* .22, *langlegrow*, subject-*verbrangle* .14, ...

Since the DPs represent the contexts of the two target words, the distance between the DPs is the distributional distance and, as per the distributional hypothesis, a proxy for semantic distance. A **measure of DP distance**, such as cosine, calculates the distance between two distributional profiles. While any of the measures of DP distance may be used with any of the measures of strength of association, in practice only certain com-

binations are used (see Table 3) and certain other combinations may not be meaningful (for example, Kullback-Leibler divergence with  $\phi$  coefficient). Observe from Table 3 that all standard-combination distributional measures (or at least those that are described in this paper) use either conditional probability or PMI as the measure of association.

#### 4.1 Simple co-occurrences versus syntactically related Words

Harris (1968), one of the early proponents of the distributional hypothesis, used syntactically related words to represent the context of a word. However, the strength of association of any word appearing in the context of the target words may be used to determine their distributional similarity. Dagan, Lee, and Pereira (1997), Lee (1999), and Weeds (2003) represent the context of a noun with verbs whose object it is (single syntactic relation), Hindle (1990) represents the context of a noun with verbs with which it shares the verb-object or subject-verb relation, while Lin (1998b) uses words related to a noun by any of the many pre-decided syntactic relations to determine distributional similarity. Schütze and Pedersen (1997) and Yoshida, Yukawa, and Kuwabara (2003) use all co-occurring words in a pre-decided window size. Although Lin (1998b) shows that the use of multiple syntactic relations is more beneficial as compared to just one, McCarthy et al. (2007) show that results obtained using just word co-occurrences produced almost as good results as those obtained using syntactically related words. Further, use of syntactically related words entails the requirement of chunking or parsing the data.

#### 4.2 Compositionality

The various measures of distributional similarity may be divided into two kinds as per their composition. In certain measures, each co-occurring word contributes to some *finite calculable* distributional distance between the target words. The final score of distributional distance is the sum of these contributions. We will call such measures **compositional measures**. The relative entropy-based measures,  $L_1$  norm and  $L_2$  norm, fall in this category. On the other hand, the cosine measure, along with Hindle's and Lin's mutual information-based measures, belong to the category of what we call **non-compositional** measures. Each co-occurring word shared by both target words contributes a score to the numerator and the denominator of the measures' formula. Words that co-occur with just one of the two target words contribute scores only to the denominator. The ratio is calculated once all co-occurring words are considered. Thus the distributional distance contributed by individual co-occurrences is not calculable and the final semantic distance cannot be broken down into compositional distances contributed by each of the co-occurrences. It is not clear as to which of the two kinds of measures (compositional or non-compositional) resembles human judgment more closely and how much they differ in their ranking of word pairs.

*Primary Compositional Measures.* The compositional measures of distributional similarity (or relatedness) capture the contribution to distance between the target words ( $w_1$  and  $w_2$ ) due to a co-occurring word by three primary mathematical manipulations of the co-occurrence distributions ( $d_1$  and  $d_2$ ): the **difference**, denoted by *Dif* (as in  $L_1$  norm), **division**, denoted by *Div* (as in the relative entropy-based measures), and **product**, denoted by *Pdt* (as in the conditional probability or mutual information-based cosine method). We will call the three types of compositional measures **primary compositional measures (PCM)**. Their form is depicted below:

$$Dif = \sum_{w \in C(w_1) \cup C(w_2)} |P(w|w_1) - P(w|w_2)| \quad (42)$$

$$Div = \sum_{w \in C(w_1) \cup C(w_2)} \left| \log \frac{P(w|w_1)}{P(w|w_2)} \right| \quad (43)$$

$$Pdt = \sum_{w \in C(w_1) \cup C(w_2)} \frac{P(w|w_1) \times P(w|w_2)}{\text{Scaling Factor}} \quad (44)$$

Observe that by taking absolute values in expressions (42) and (43), the variation in the distributions for different co-occurring words has an additive affect rather than one of cancellation. This corresponds to our distributional hypothesis — the more the disparity in distributions, the more is the semantic distance between the target words. The product form (44) also achieves this and is based on this theorem: The product of any two numbers will always be less than or equal to the square of their average. In other words, the more two numbers are close to each other in value, the higher is the ratio of their product to a suitable scaling factor (for example, the square of their average). Note that the difference and division measures give higher values when there is large disparity between the strength of association of co-occurring words with the target words. They are therefore measures of distributional distance and not distributional similarity. The product method gives higher values when the strengths of association are closer, and is a measure of distributional relatedness.

Although all three methods seem intuitive, each produces different distributional similarity values and more importantly, given a set of word pairs, each is likely to rank them differently. For example, consider the division and difference expressions applied to word pairs  $(w_1, w_2)$  and  $(w_3, w_4)$ . For simplicity, let there be just one word  $w'$  in the context of all the words. Given:

$$P(w'|w_1) = 0.91$$

$$P(w'|w_2) = 0.80$$

$$P(w'|w_3) = 0.60$$

$$P(w'|w_4) = 0.50$$

The distributional distance between word pairs as per the difference PCM:

$$Dif(w_1, w_2) = |0.91 - 0.8| = 0.11$$

$$Dif(w_3, w_4) = |0.6 - 0.5| = 0.1$$

The distributional distance between word pairs as per the division PCM:

$$Div(w_1, w_2) = \left| \log \frac{0.91}{0.8} \right| = 0.056$$

$$Div(w_3, w_4) = \left| \log \frac{0.6}{0.5} \right| = 0.079$$

Observe that for the same set of co-occurrence probabilities, the difference-based measure ranks the  $(w_3, w_4)$  pair more distributionally similar (lower distributional distance), while the division-based measure gives lower distributional similarity values for word pairs having large co-occurrence probabilities. This behavior is not intuitive and it remains to be seen, by experimentation, as to which of the three, difference, division or product, yields distributional similarity measures closest to human notions of semantic similarity.

The  $L_1$  norm is a basic implementation of the difference method. A simple product-based measure of distributional similarity is as proposed below:

$$Pdt^{Avg}(w_1, w_2) = \sum_{w \in C(w_1) \cup C(w_2)} \frac{P(w|w_1) \times P(w|w_2)}{\left(\frac{1}{2}(P(w|w_1) + P(w|w_2))\right)^2} \quad (45)$$

The scaling factor used is the square of the average probability. It can be proved that if the sum of two variables is equal to a constant ( $k$ , say), their values must be equal to  $k/2$  in order to get the largest product. Now, let  $k$  be equal to the sum of  $P(w|w_1)/(P(w|w_1) + P(w|w_2))$  and  $P(w|w_2)/(P(w|w_1) + P(w|w_2))$ . This sum will always be equal to 1 and hence the product ( $Z$ ) will be largest only when the two numbers are equal i.e.  $P(w|w_1)$  is equal to  $P(w|w_2)$ . In other words, the farther  $P(w|w_1)$  and  $P(w|w_2)$  are from their average, the smaller is the product  $Z$ . Therefore, the measure gives high scores for low disparity in strengths of co-occurrence and low scores otherwise. The incorporation of  $1/2$  in the scaling factor results in a measure that ranges between 0 and 1.

The relative entropy-based methods use a weighted division method. Observe that both Kullback-Leibler divergence (formula repeated below for convenience — equation (46)) and Jensen-Shannon divergence do not take absolute values of the division of co-occurrence probabilities. This will mean that if  $P(w|w_1) > P(w|w_2)$ , the logarithm of their ratio will be positive and if  $P(w|w_1) < P(w|w_2)$ , the logarithm will be a negative number. Therefore, there will be a cancellation of contributions to distributional distance by words that have higher co-occurrence probability with respect to  $w_1$  and words that have a higher co-occurrence probability with respect to  $w_2$ . Observe however that the weight  $P(w|w_1)$  multiplying the logarithm means that in general the positive logarithm values receive higher weight than the negative ones, resulting in a net positive score. Therefore, with no absolute value of the logarithm, as in the KLD, the weight plays a crucial role. A modified Kullback-Leibler divergence ( $D^{Abs}$ ) which incorporates the absolute value is suggested in equation (47).<sup>10</sup>

$$KLD(w_1, w_2) = D(d_1||d_2) = \sum_{w \in C(w_1) \cup C(w_2)} P(w|w_1) \log \frac{P(w|w_1)}{P(w|w_2)} \quad (46)$$

<sup>10</sup> It should be noted that any changes to the formula for Kullback-Leibler divergence means that the resulting measure is no longer Kullback-Leibler divergence; these measure are denoted by KLD (and a suitable subscript and/or superscript simply to indicate that they are derived from the Kullback-Leibler divergence.

$$KLD^{Abs}(w_1, w_2) = D^{Abs}(d_1 || d_2) = \sum_{w \in C(w_1) \cup C(w_2)} P(w|w_1) \left| \log \frac{P(w|w_1)}{P(w|w_2)} \right| \quad (47)$$

The updated Jensen-Shannon divergence measure will remain the same as in equation (17), except that it is a manipulation of  $D^{Abs}$  and not the original Kullback-Leibler divergence (relative entropy).

$$JSD^{Abs}(w_1, w_2) = D^{Abs}(d_1 || \frac{1}{2}(d_1 + d_2)) + D^{Abs}(d_2 || \frac{1}{2}(d_1 + d_2)) \quad (48)$$

Note that once the absolute value of the logarithm is taken, it no longer makes much sense to use an asymmetric weight ( $P(w|w_1)$ ) as in the KLD or as necessary to use a weight at all. Equation (49) shows a simple division-based measure. It is an unweighted form of  $KLD^{Abs}(w_1, w_2)$  and so we will call it  $KLD_{Unw}^{Abs}$ .

$$KLD_{Unw}^{Abs}(w_1, w_2) = Div(w_1, w_2) = \sum_{w \in C(w_1) \cup C(w_2)} \left| \log \frac{P(w|w_1)}{P(w|w_2)} \right| \quad (49)$$

Experimental evaluation of these suggested modifications of Kullback-Leibler divergence will be interesting.

*Weighting the PCMs.* The performance of the primary compositional measures may be improved by adding suitable weights to the distributional distance contributed by each co-occurrence. The idea is that some co-occurrences may be better indicators of semantic distance than others. Usually, a formulation of the strength of association of the co-occurring word with the target words is used as weight, the hypothesis being that a strong co-occurrence is likely to be a strong indicator of semantic closeness.

Weighting the primary compositional measures results in some of the existing measures. For example, as pointed out earlier, the Kullback-Leibler divergence is a weighted form of the division measure (not considering the absolute value). Here, the conditional probability of a co-occurring word with respect to the first word ( $P(w|w_1)$ ) is used as the weight. Since the absolute value of the logarithm is not taken and because the weight ( $P(w|w_1)$ ) is dependent on the first word and not the other, Kullback-Leibler divergence is asymmetric. Below is a symmetric weight function:

$$weight_{AvgWt}(w_1, w_2) = \frac{1}{2} (P(w|w_1) + P(w|w_2)) \quad (50)$$

$$(51)$$

$L_2$  norm is a weighted version of the  $L_1$  norm, the weight being  $P(w|w_1) - P(w|w_2)$ . A simple product measure with weights is shown below:

$$Pdt_{AvgWt}^{Avg} = \sum_{w \in C(w_1) \cup C(w_2)} \frac{1}{2} (P(w|w_1) + P(w|w_2)) \frac{P(w|w_1) \times P(w|w_2)}{(\frac{1}{2}(P(w|w_1) + P(w|w_2)))^2}$$

$$= \sum_{w \in C(w_1) \cup C(w_2)} \frac{P(w|w_1) \times P(w|w_2)}{\frac{1}{2}(P(w|w_1) + P(w|w_2))} \quad (52)$$

A possibly better weight function (which is also symmetric) hinges on the following hypothesis: The stronger the association of a co-occurring word with a target word, the better the indicator of semantic properties of the target word it is. Equation (53) shows the corresponding weight function:

$$weight_{MaxWt}(w_1, w_2) = \frac{\max(P(w|w_1), P(w|w_2))}{\sum_{w' \in C(w_1) \cup C(w_2)} \max(P(w'|w_1), P(w'|w_2))} \quad (53)$$

The co-occurring word is likely to have different strengths of associations with the two target words. Taking the maximum of the two as the weight (Dagan, Marcus, and Markovitch (1995)) will mean that more weight is given to a co-occurring word if it has high strength of association with any of the two target words. As Dagan, Marcus, and Markovitch (1995) point out, there is strong evidence for dissimilarity if the strength of association with the other target word is much lower than the maximum, and strong evidence of similarity if the strength of association with both target words is more or less the same.

### 4.3 Predictors of Semantic Relatedness

Given a pair of target words, the vocabulary may be divided into three sets: (1) the set of words that co-occur with both target words (common); (2) words that co-occur with exactly one of the two target words (exclusive); (3) words that do not co-occur with either of the two target words. Hindle (1990) uses evidence only from words that co-occur with both target words to determine the distributional similarity. All the other measures discussed in this paper so far also use words that co-occur with just one target word.

One can argue that the more there are common co-occurrences between two words, the more they are related. For example, *drink* and *sip* may be considered related as they have a number of common co-occurrences such as *water*, *tea* and so on. Similarly, *drink* and *chess* can be deemed unrelated as words that co-occur with one do not with the other. For example, *water* and *tea* do not usually co-occur with *chess*, while *castle* and *move* are not found close to *drink*. Measures that use all co-occurrences (common and exclusive) tap into this intuitive notion. However, certain strong exclusive co-occurrences can adversely effect the measure. Consider the classic *strong coffee vs powerful coffee* example (Halliday (1966)). The words *strong* and *powerful* are semantically very related. However, the word *coffee* is likely to co-occur with *strong* but not with *powerful*. Further, *strong* and *coffee* can be expected to have a large value of association as given by a suitable measure, say PMI. This large PMI value, if used in the distributional relatedness formula, can greatly reduce the final value. Thus it is not clear if the benefit of using all co-occurrences is outweighed by the drawback pointed out.

A further advantage of using only common co-occurrences is that the Kullback-Leibler divergence can now be used without the need of smoothed probabilities.

$$KLD_{Com}(w_1, w_2) = \sum_{w \in C(w_1) \cap C(w_2)} P(w|w_1) \log \frac{P(w|w_1)}{P(w|w_2)} \quad (54)$$

Observe that we are taking the intersection of the set of co-occurring words instead of union as in the original formula (15).

#### 4.4 Capitalizing on Asymmetry

Given a hypernym-hyponym pair (*automobile-car*, say) asymmetric distributional measures such as the Kullback-Leibler divergence,  $\alpha$  skew divergence, and the CRMs generate different values as the distributional distance of  $w_1$  with  $w_2$  as compared to that of  $w_2$  with  $w_1$ . Usually, if  $w_1$  is a more generic concept than  $w_2$ , the measures find  $w_1$  to be more distributionally similar to  $w_2$  than the other way round (see (Mirkin, Dagan, and Geffet 2007) for work on lexical entailment using the Kullback-Leibler divergence). Weeds (2003) argues that this behavior is intuitive as it is more often acceptable to substitute a generic concept in place of a specific one than vice versa, and substitutability is a indicator of semantic similarity.

On the other hand, in most cases the notion of asymmetric semantic similarity is counterintuitive, and possibly detrimental. In many natural language tasks, one needs the distance between two words and there is no order information. Further, in case two words share a hypernym-hyponym relation, they are likely to be highly semantically similar. Thus given two words, it may make sense to always choose the higher of the two distributional similarity values suggested by an asymmetric measure as the final distributional similarity between the two. This way an asymmetric measure ( $Sim_{Asym}$ ) can easily be converted into a symmetric one ( $Sim_{Max}$ ), while still capitalizing on the asymmetry to generate more suitable distributional distance values for hypernym-hyponym word pairs. Equation (55) states the formula for the proposed conversion. A specific implementation of the KL divergence formula is given in equation (56):

$$Sim_{Max}(w_1, w_2) = \max(Sim_{Asym}(w_1, w_2), Sim_{Asym}(w_2, w_1)) \quad (55)$$

$$KLD_{Max}(w_1, w_2) = \max(KLD(w_1, w_2), KLD(w_2, w_1)) \quad (56)$$

Another way to convert an asymmetric measure of distributional distance into a symmetric one is by taking the average (formula 57) of the two possible similarity values. A specific implementation on the KL divergence formula is given in equations (58) through (61):

$$Sim_{Avg}(w_1, w_2) = \frac{1}{2}(Sim_{Asym}(w_1, w_2) + Sim_{Asym}(w_2, w_1)) \quad (57)$$

$$KLD_{Avg}(w_1, w_2) = \frac{1}{2}(KLD(w_1, w_2) + KLD(w_2, w_1)) \quad (58)$$

$$= \frac{1}{2} \sum_{w \in C(w_1) \cup C(w_2)} \left( P(w|w_1) \log \frac{P(w|w_1)}{P(w|w_2)} + P(w|w_2) \log \frac{P(w|w_2)}{P(w|w_1)} \right) \quad (59)$$

$$= \frac{1}{2} \sum_{w \in C(w_1) \cup C(w_2)} \left( P(w|w_1) \log \frac{P(w|w_1)}{P(w|w_2)} - P(w|w_2) \log \frac{P(w|w_1)}{P(w|w_2)} \right) \quad (60)$$

$$= \frac{1}{2} \sum_{w \in C(w_1) \cup C(w_2)} (P(w|w_1) - P(w|w_2)) \log \frac{P(w|w_1)}{P(w|w_2)} \quad (61)$$

#### 4.5 Summarizing the distributional measures

Tables 4 and 5 in the appendix summarize the properties of various distributional measures discussed in this paper.

#### 4.6 Challenges

**4.6.1 Conflation of word senses.** The distributional hypothesis (Firth 1957) states that words that occur in similar contexts tend to be semantically close. But when words have more than one sense, it is not at all clear what semantic distance between them actually means. Further, a word in each of its senses is likely to co-occur with different sets of words. For example, *bank* in the FINANCIAL INSTITUTION sense is likely to co-occur with *interest*, *money*, *accounts*, and so on, whereas the RIVER BANK sense might have words such as *river*, *erosion*, and *silt* around it. Since words that occur together in text tend to refer to senses that are closest in meaning to one another, in most natural language applications, what is needed is the distance between the closest senses of the two target words. However, because distributional measures calculate distance from occurrences of the target word in all its occurrences and hence all its senses, they fail to get the desired result. Also note that the dimensionality reduction inherent to latent semantic analysis, a special kind of distributional measure, has the effect of making the predominant senses of the words more dominant while de-emphasizing the other senses. Therefore, an LSA-based approach will also conflate information from the different senses, and even more emphasis will be placed on the predominant senses. Given the semantically close target nouns *play* and *actor*, for example, a distributional measure will give a score that is some sort of a dominance-based average of the distances between their senses. The noun *play* has the predominant sense of CHILDREN'S RECREATION (and not DRAMA), so a distributional measure will tend to give the target pair a large (and thus erroneous) distance score. WordNet-based measures do not suffer from this problem as they give distance between concepts, not words.

**4.6.2 Lack of explicitly-encoded world knowledge and data sparseness.** It is becoming increasingly clear that more-accurate results can be achieved in a large number of natural language tasks, including the estimation of semantic distance, by combining corpus statistics with a knowledge source, such as a dictionary, published thesaurus, or WordNet. This is because such knowledge sources capture semantic information about concepts and, to some extent, world knowledge. For example, WordNet, as discussed earlier, has an extensive is-a hierarchy. If it lists one concept, say GERMAN SHEPHERD as a hyponym of another, say DOG, then we can be sure that the two are semantically close. On the other hand, distributional measures do not have access to such explicitly encoded information. Further, unless the corpus used by a distributional measure has sufficient instances of GERMAN SHEPHERD and DOG, it will be unable to deem them semantically close. Since Zipf's law seems to hold even for the largest of corpora, there will always be words that occur too few times to accurately determine their distributional distance from others.

**4.6.3 Limitations shared with WordNet-based measures.** In addition to the limitations described above, which are unique to the knowledge-lean distributional measures, like the knowledge-rich measures they also suffer from problems of requiring the calculation of large distance matrices (as described in Section 2.3.4 earlier) and the reluctance to cross the language barrier (Section 2.3.5).

## 5. A hybrid approach: Distributional measures of concept-distance

So far we have looked at knowledge source-based approaches that exploit the structure of a resource such as WordNet, and studied in detail corpus-based distributional measures that make use of co-occurrence statistics. A **hybrid approach to semantic distance** will reconcile and combine the information about concepts, explicitly encoded in a linguistic resource, and information about words, implicitly encoded in text via co-occurrence.

Mohammad and Hirst (2006b) and Mohammad et al. (2007) have proposed one such hybrid, yet distinctly distributional, approach that combines corpus statistics with a published thesaurus to give the semantic distance between concepts (rather than words).

### 5.1 The distributional hypothesis for concepts

As pointed out in Section 4.6.1, words when used in different senses tend to keep different “company” (co-occurring words). For example, consider the contrived but plausible distributional profile of *star*:

*star*: *space* 0.21, *movie* 0.16, *famous* 0.15, *light* 0.12, *constellation* 0.11, *heat* 0.08, *rich* 0.07, *hydrogen* 0.07, ...

Observe that it has words that co-occur both with *star*’s CELESTIAL BODY sense and *star*’s CELEBRITY sense. Thus, it is clear that different senses of a word may have very different distributional profiles. Using a single DP for the word will mean the union of those profiles. While this might be useful for certain applications, Mohammad and Hirst (2006b) argue that in a number of tasks (including estimating semantic distance), acquiring different DPs for the different senses is not only more intuitive, but also, as they show through experiments, more useful. They show that the closer the distributional profiles of two concepts, the smaller is their semantic distance. Below are example distributional profiles of two senses of *star*:

CELESTIAL BODY: *space* 0.36, *light* 0.27, *constellation* 0.11, *hydrogen* 0.07, ...  
 CELEBRITY: *famous* 0.24, *movie* 0.14, *rich* 0.14, *fan* 0.10, ...

The values are the strength of association (usually pointwise mutual information or conditional probability) of the target concept with co-occurring words.

We have seen that creating distributional profiles of words involves simple word-word co-occurrence counts. The creation of DPCs, on the other hand, requires: (1) a concept inventory that list all the concepts and words that refer to them, and (2) counts of how often a concept co-occurs with a word in text. These two aspects will be discussed in the next two sub-sections; however once created, any of the many distributional measures can be used to estimate the distance between the DPs of two target concepts (just as in the case of traditional word-distance measures, where distributional measures are used to estimate the distance between the DPs of two target words). For example, here is how they adapt the formula for cosine (described earlier in Section 3.2.1) to

estimate distributional distance between two concepts:

$$\text{Cos}_{cp}(c_1, c_2) = \frac{\sum_{w \in C(c_1) \cup C(c_2)} (P(w|c_1) \times P(w|c_2))}{\sqrt{\sum_{w \in C(c_1)} P(w|c_1)^2} \times \sqrt{\sum_{w \in C(c_2)} P(w|c_2)^2}} \quad (62)$$

$C(x)$  is now the set of words that co-occur with *concept*  $x$  within a pre-determined window. The conditional probabilities in the formula are taken from the distributional profiles of concepts.

**5.1.1 A suitable knowledge source and concept inventory.** Mohammad and Hirst (2006b) use the categories in the *Macquarie Thesaurus*, 812 in all, as very coarse-grained word senses or concepts, in contrast to approaches that use WordNet or other similarly fine-grained sense inventories.<sup>11</sup> Their approach to determining word–concept co-occurrence counts (described in the next sub-section) requires a set of possibly ambiguous words that together unambiguously represent each concept—for which a thesaurus is a natural choice. The use of categories in a thesaurus as concepts means that this approach requires a concept–concept distance matrix of size only  $812 \times 812$ —much smaller than (less than 0.01% of) the matrix required by the WordNet-based and distributional measures.

**5.1.2 Estimating distributional profiles of concepts.** A **word–category co-occurrence matrix (WCCM)** is created having word types  $w$  as one dimension and thesaurus categories  $c$  as another.

	$c_1$	$c_2$	...	$c_j$	...
$w_1$	$m_{11}$	$m_{12}$	...	$m_{1j}$	...
$w_2$	$m_{21}$	$m_{22}$	...	$m_{2j}$	...
⋮	⋮	⋮	⋱	⋮	⋮
$w_i$	$m_{i1}$	$m_{i2}$	...	$m_{ij}$	...
⋮	⋮	⋮	...	⋮	⋱

The matrix is populated with co-occurrence counts from a large corpus. A particular cell  $m_{ij}$ , corresponding to word  $w_i$  and category or concept  $c_j$ , is populated with the number of times  $w_i$  co-occurs (they use a window of  $\pm 5$  words) with any word that has  $c_j$  as one of its senses (i.e.,  $w_i$  co-occurs with any word listed under concept  $c_j$  in the thesaurus). This matrix, created after a first pass of the corpus, is called the **base word–category co-occurrence matrix (base WCCM)**. A contingency table for any particular word  $w$  and category  $c$  can be easily generated from the WCCM by collapsing cells for all other words and categories into one and summing up their frequencies.

	$c$	$\neg c$
$w$	$n_{wc}$	$n_{w\neg c}$
$\neg w$	$n_{\neg w c}$	$n_{\neg w \neg c}$

<sup>11</sup> It has been suggested for some time now that WordNet is much too fine-grained for certain natural language applications (Agirre and Lopez de Lacalle Lekuona (2003) and citations therein).

A suitable statistic, such as pointwise mutual information or conditional probability, will then yield the strength of association between the word and the category.

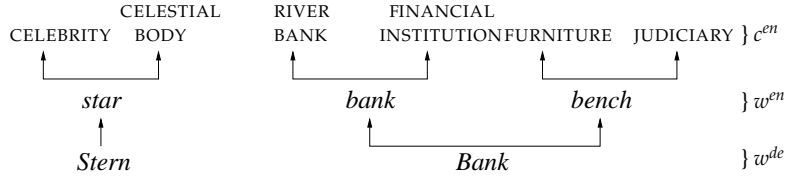
As the base WCCM is created from unannotated text, it will be noisy but nonetheless capture strong word–category co-occurrence associations reasonably accurately. This is because the errors in determining the true category that a word co-occurs with will be distributed thinly across a number of other categories. (For more discussion of the general principle see Resnik (1998).) A second pass of the corpus is made and the base WCCM is used to disambiguate the words in it. A new **bootstrapped WCCM** is created such that each cell  $m_{ij}$ , corresponding to word  $w_i$  and concept  $c_j$ , is populated with the number of times  $w_i$  co-occurs with any word *used in sense*  $c_j$ . Mohammad and Hirst (2006a) showed that this WCCM, created after simple sense disambiguation, better captures word–concept co-occurrence values, and hence strengths of association values, than the base WCCM.

**5.1.3 Mimicking semantic relatedness and semantic similarity.** The distributional profiles created by the above methodology are relation-free. This is because (1) all co-occurring words (not just those that are related to the target by certain syntactic relations) are used, and (2) the WCCM, as described above, does not maintain separate counts for the different syntactic relations between the target and co-occurring words. Thus, distributional measures that use these WCCMs will estimate semantic *relatedness* between concepts. Distributional measures that mimic semantic *similarity*, which require relation-constrained DPCs, can easily be created from WCCMs that have rows for each word–syntactic relation pair (rather than just for words).

**5.1.4 Performance.** Mohammad and Hirst (2006b) evaluate this approach on two tasks: ranking word pairs in order of their semantic distance and correcting real-word spelling errors. On both tasks, distributional concept-distance measures markedly outperformed distributional word-distance measures. The WordNet-based measures performed better in the word-pair ranking task, but in the spelling correction task three of the four distributional measure outperformed all WordNet-based measures except the Jiang–Conrath measure. It should be noted, however, that these experiments evaluated only semantic similarity of noun–noun pairs—for all other part-of-speech combinations and semantic relatedness estimates, the WordNet-based measures are markedly less accurate.

## 5.2 Multilinguality

Some of the best algorithms for semantic distance cannot be applied to most languages because of a lack of high-quality linguistic resources. Mohammad et al. (2007) showed how text in one language  $L_1$  can be combined with a knowledge source in another  $L_2$  using a bilingual lexicon  $L_1$ – $L_2$  and a bootstrapping and concept-disambiguation algorithm to create **cross-lingual distributional profiles of concepts**. These cross-lingual DPCs model co-occurrence distributions of concepts, as per a knowledge source in one language, with words from another language, and obtain semantic distance in a resource-poor language using a knowledge source from a resource-rich one. Cross-lingual semantic distance and cross-lingual DPCs are also useful in tasks that inherently involve two or more languages, such as machine translation, multilingual multidocument tasks of clustering, coreference resolution, and information retrieval. We summarize their approach here using German as  $L_1$  and English as  $L_2$ ; however, the algorithm is language-pair independent.



**Figure 1**  
The cross-lingual candidate senses of German words *Stern* and *Bank*.

**5.2.1 Cross-lingual senses, cross-lingual distributional profiles, and cross-lingual distributional distance.** Given a German word  $w^{de}$  in context, Mohammad et al. (2007) use a German–English bilingual lexicon to determine its different possible English translations. Each English translation  $w^{en}$  may have one or more possible coarse senses, as listed in an English thesaurus. These English thesaurus concepts ( $c^{en}$ ) will be referred to as the **cross-lingual candidate senses** of the German word  $w^{de}$ . Figure 1 depicts examples.<sup>12</sup>

As in the monolingual estimation of distributional concept-distance, the distance between two concepts is calculated by first determining their DPs. However, a concept is now glossed by near-synonymous words in an *English* thesaurus, whereas its profile is made up of the strengths of association with co-occurring *German* words. Here are constructed example cross-lingual distributional profiles of the two cross-lingual candidate senses of the German word *Stern*:<sup>13</sup>

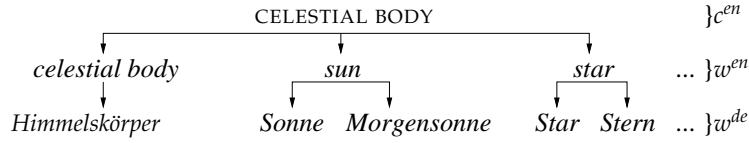
CELESTIAL BODY (*celestial body, sun, . . .*): *Raum* 0.36, *Licht* 0.27, *Konstellation* 0.11, . . .  
 CELEBRITY (*celebrity, hero, . . .*): *berühmt* 0.24, *Film* 0.14, *reich* 0.14, . . .

The cross-lingual DPCs are created from a cross-lingual word–category co-occurrence matrix without the use of any word-aligned parallel corpora or sense-annotated data (as described in the next subsection). Just as in the case of monolingual distributional concept-distance measures, distributional measures can be used to estimate the distance between the cross-lingual DPs of two target concepts. For example, the cosine formula can be adapted to estimate cross-lingual distributional distance between two concepts as shown below:

$$\text{Cos}(c_1^{en}, c_2^{en}) = \frac{\sum_{w^{de} \in C(c_1^{en}) \cup C(c_2^{en})} (P(w^{de}|c_1^{en}) \times P(w^{de}|c_2^{en}))}{\sqrt{\sum_{w^{de} \in C(c_1^{en})} P(w^{de}|c_1^{en})^2} \times \sqrt{\sum_{w^{de} \in C(c_2^{en})} P(w^{de}|c_2^{en})^2}} \quad (63)$$

12 They are called “candidate senses” because some of the senses of  $w^{en}$  might not really be senses of  $w^{de}$ . For example, CELESTIAL BODY and CELEBRITY are both senses of the English word *star*, but the German word *Stern* can only mean CELESTIAL BODY and not CELEBRITY. Similarly, the German *Bank* can mean FINANCIAL INSTITUTION or FURNITURE, but not RIVER BANK or JUDICIARY. An automated system has no straightforward method of teasing out the actual cross-lingual senses of  $w^{de}$  from those that are an artifact of the translation step.

13 Vocabulary of German words needed to understand this discussion: *Bank*: 1. financial institution, 2. bench (furniture); *berühmt*: famous; *Film*: movie (motion picture); *Himmelskörper*: heavenly body; *Konstellation*: constellation; *Licht*: light; *Morgensonne*: morning sun; *Raum*: space; *reich*: rich; *Sonne*: sun; *Star*: star (celebrity); *Stern*: star (celestial body)



**Figure 2**  
Words having CELESTIAL BODY as one of their cross-lingual candidate senses.

$C(x)$  is now the set of German words that co-occur with English concept  $x$  within a pre-determined window. The conditional probabilities in the formula are taken from the cross-lingual DPCs.

**5.2.2 Creating cross-lingual word–category co-occurrence matrix.** A German–English cross-lingual word–category co-occurrence matrix has German word types  $w^{de}$  as one dimension and English thesaurus concepts  $c^{en}$  as another.

	$c_1^{en}$	$c_2^{en}$	...	$c_j^{en}$	...
$w_1^{de}$	$m_{11}$	$m_{12}$	...	$m_{1j}$	...
$w_2^{de}$	$m_{21}$	$m_{22}$	...	$m_{2j}$	...
⋮	⋮	⋮	⋱	⋮	⋮
$w_i^{de}$	$m_{i1}$	$m_{i2}$	...	$m_{ij}$	...
⋮	⋮	⋮	...	⋮	⋱

The matrix is populated with co-occurrence counts from a large German corpus. A particular cell  $m_{ij}$ , corresponding to word  $w_i^{de}$  and concept  $c_j^{en}$ , is populated with the number of times the German word  $w_i^{de}$  co-occurs (say a window of  $\pm 5$  words) with any German word having  $c_j^{en}$  as one of its *cross-lingual candidate senses*. For example, the *Raum*–CELESTIAL BODY cell will have the sum of the number of times *Raum* co-occurs with *Himmelskörper*, *Sonne*, *Morgensonne*, *Star*, *Stern*, and so on (see Figure 2). This matrix, created after a first pass of the corpus, is called the **cross-lingual base WCCM**. A contingency table for any particular German word  $w^{de}$  and English category  $c^{en}$  can be easily generated from the WCCM by collapsing cells for all other words and categories into one and summing up their frequencies. A suitable statistic, such as PMI or conditional probability, will yield the strength of association between the German word and the English category. Then a new bootstrapped cross-lingual WCCM is created, just as in the monolingual case.

**5.2.3 Performance.** Mohammad et al. (2007) evaluated the cross-lingual measures of semantic distance on two tasks: (1) estimating semantic distance between words and ranking the word pairs according to semantic distance, and (2) solving *Reader’s Digest* ‘Word Power’ problems. They compared these results with those obtained by conventional state-of-the-art monolingual approaches with and without a knowledge source in the target language  $L_1$  (GermaNet). The cross-lingual approach obtained much better results than monolingual approaches that do not use a knowledge source. Further, in both tasks, the cross-lingual measures performed as well if not slightly better than the GermaNet-based monolingual approaches, as well. This shows that the cross-lingual

approach is able to keep losses due to the translation step at a minimum, while allowing the use of a superior knowledge source in another language to get better results.

### 5.3 Challenges

Distributional measures of concept-distance have many desirable features of both knowledge-rich approaches and strictly corpus-based approaches—they have the high accuracies of knowledge-rich approaches, they can measure both semantic relatedness and semantic similarity, and they have a strong corpus-reliance making them domain adaptable. Further, with the cross-lingual approach, a lack of high-quality knowledge source in the target language is no longer a problem. However, certain issues remain.

**5.3.1 Integrating domain-specific terminology.** The reliance on a knowledge source means that the approach cannot measure the distance between words not listed in the thesaurus. This is especially a problem for domain-specific jargon, which may not find place in a general purpose knowledge source. Automatic ways of integrating domain-specific terminology into a general purpose knowledge source will be valuable to this end.

**5.3.2 Choosing the right concept-granularity.** So far Mohammad and Hirst (2006b) and Mohammad et al. (2007) have reported results using the categories of the thesaurus as very coarse word senses. This level of granularity has worked well for the tasks they experimented with; however, a relatively finer sense inventory may be more suitable for other tasks. Words within a thesaurus category are grouped into paragraphs; and using them (instead of categories) and determining when this finer-grained sense-inventory is more suitable for use remains to be explored.

**5.3.3 Identifying lexical semantic relations.** Word pairs can be semantically close because of any of the classical lexical semantic relations, such as hypernymy, near-synonymy, antonymy, troponymy, and meronymy, or the innumerable non-classical relations. The various distributional approaches discussed in this paper determine semantic distance without explicitly identifying the nature of the relationship. Already, there is some work on determining lexical entailment (Mirkin, Dagan, and Geffet 2007) and determining near-synonymy (Lin et al. 2003). Identifying antonymy (or more generally, contrasting word-pairs) is especially useful in many natural language tasks, even if it is simply to discard them from a list of distributionally close terms. Also, it will be interesting for measures of semantic distance to characterize the nature of any non-classical relationship shared by two words—not only determining if two terms are close but also specifying (in some intuitive way) the aspect of meaning they share.

## 6. Conclusion

A large number of important natural language problems, including machine translation, information retrieval, and word sense disambiguation, can be viewed in part as semantic distance problems. Numerous measures of semantic distance exist—those that use a knowledge source and those that rely on corpora. Yet, their use in real-world applications has been limited. In this paper, we investigated how automatic measures can be brought more in line with human notions of semantic distance, how they can be made applicable to a large number of natural language tasks, and how they can be used even for languages deficient in high-quality linguistic resources.

Even though corpus-based distributional measures of distance have traditionally performed poorly when compared to WordNet-based measures, we have shown that (1) there are a number of reasons that make distributional measures uniquely attractive, and (2) that their potential is yet to be fully realized. Distributional measures can be easily applied to most languages (they can make do even with just raw text) and they can be used to mimic both semantic similarity and semantic relatedness. With this in mind, the paper presented a detailed study of many important distributional measures, analyzed their limitations, and explained why their performance has been relatively poor so far. Understanding these limitations is crucial in the development of new and better approaches, whether they have a distributional base or otherwise. We concluded the paper with the discussion of a hybrid, yet distinctly distributional approach, that presents one way to more accurately measure distributional distance without compromising too much on essential properties such as the applicability to resource-poor languages.

## **7. Appendix**

Tables 4 and 5 details the properties of various distributional measures discussed in this paper.

**Table 4**  
Distributional measures and their properties.

distributional measure	measure type	compositional	PCM	formula	symmetric	strength of association
<i>ASD</i>	distance	✓	division	$\sum_{w \in C(w_1) \cup C(w_2)} P(w w_1) \log \frac{P(w w_1)}{\alpha P(w w_2) + (1-\alpha)P(w w_1)}$	X	CP
cos	closeness	X	n.a. <sup>0</sup>	$\frac{\sum_{w \in C(w_1) \cup C(w_2)} (P(w w_1) \times P(w w_2))}{\sqrt{\sum_{w \in C(w_1)} P(w w_1)^2} \times \sqrt{\sum_{w \in C(w_2)} P(w w_2)^2}}$	✓	CP
<i>CRMs</i>	closeness	X	n.a.	$\gamma \left[ \frac{2 \times P \times R}{P + R} \right] + (1 - \gamma) \left[ \beta [P] + (1 - \beta) [R] \right]$	X	both
<i>Dice<sup>CP</sup></i>	closeness	X	n.a.	$\frac{2 \times \sum_{w \in C(w_1) \cup C(w_2)} \min(P(w w_1), P(w w_2))}{\sum_{w \in C(w_1)} P(w w_1) + \sum_{w \in C(w_2)} P(w w_2)}$	✓	CP
<i>Dif</i> or $L_1$	distance	✓	difference	$\sum_{w \in C(w_1) \cup C(w_2)}  P(w w_1) - P(w w_2) $	✓	CP
<i>Div</i>	distance	✓	division	$\sum_{w \in C(w_1) \cup C(w_2)} \left  \log \frac{P(w w_1)}{P(w w_2)} \right $	✓	CP
<i>Hindle</i>	closeness	X	n.a.	$\sum_{w \in C(w)} \begin{cases} \min(I(w, w_1), I(w, w_2)), & \text{if } I(w, w_1) > 0 \text{ and } I(w, w_2) > 0 \\   \max(I(w, w_1), I(w, w_2))  , & \text{if } I(w, w_1) < 0 \text{ and } I(w, w_2) < 0 \\ 0, & \text{otherwise} \end{cases}$	✓	PMI
<i>Jaccard<sup>CP</sup></i>	closeness	X	n.a.	$\frac{\sum_{w \in C(w_1) \cup C(w_2)} \min(P(w w_1), P(w w_2))}{\sum_{w \in C(w_1) \cap C(w_2)} \max(P(w w_1), P(w w_2))}$	✓	CP
<i>JSD</i>	distance	✓	division	$\sum_{w \in C(w_1) \cup C(w_2)} \left( P(w w_1) \log \frac{P(w w_1)}{\frac{1}{2}(P(w w_1) + P(w w_2))} + P(w w_2) \log \frac{P(w w_2)}{\frac{1}{2}(P(w w_1) + P(w w_2))} \right)$	✓	CP

**Table 5**  
Distributional measures and their properties.

distributional measure	measure type	compositional	PCM	formula	symmetric	strength of association
$KLD$	distance	✓	div.	$\sum_{w \in C(w_1) \cup C(w_2)} P(w w_1) \log \frac{P(w w_1)}{P(w w_2)}$	X	CP
$KLD_{Com}$	distance	✓	div.	$\sum_{w \in C(w_1) \cap C(w_2)} P(w w_1) \log \frac{P(w w_1)}{P(w w_2)}$	X	CP
$KLD^{Abs}$	distance	✓	div.	$\sum_{w \in C(w_1) \cup C(w_2)} P(w w_1) \left  \log \frac{P(w w_1)}{P(w w_2)} \right $	X	CP
$KLD_{Avg}$	distance	✓	div.	$\frac{1}{2} \sum_{w \in C(w_1) \cup C(w_2)} (P(w w_1) - P(w w_2)) \log \frac{P(w w_1)}{P(w w_2)}$	✓	CP
$KLD_{Max}$	distance	✓	div.	$\max(KLD(w_1, w_2), KLD(w_2, w_1))$	✓	CP
$L_2$	distance	✓	difference	$\sqrt{\sum_{w \in C(w_1) \cup C(w_2)} (P(w w_1) - P(w w_2))^2}$	✓	CP
$Lin$	closeness	X	n.a.	$\frac{\sum_{(r,w) \in T(w_1) \cap T(w_2)} (I(w_1, r, w) + I(w_2, r, w))}{\sum_{(r,w') \in T(w_1)} I(w_1, r, w') + \sum_{(r,w'') \in T(w_2)} I(w_2, r, w'')}$	✓	PMI
$Pdt^{Avg}$	closeness	✓	pdt.	$\sum_{w \in C(w_1) \cup C(w_2)} \frac{P(w w_1) \times P(w w_2)}{(\frac{1}{2}(P(w w_1) + P(w w_2)))^2}$	✓	CP
$Pdt_{AvgWt}^{Avg}$	closeness	✓	pdt.	$\sum_{w \in C(w_1) \cup C(w_2)} \frac{P(w w_1) \times P(w w_2)}{\frac{1}{2}(P(w w_1) + P(w w_2))}$	✓	CP

### Acknowledgments

We thank Suzanne Stevenson, Gerald Penn, and the Computational Linguistics group at the University of Toronto for helpful discussions. This research is financially supported by the Natural Sciences and Engineering Research Council of Canada, the University of Toronto.

### References

- Agirre, Eneko and Oier Lopez de Lacalle Lekuona. 2003. Clustering WordNet word senses. In *Proceedings of the 1st International Conference on Recent Advances in Natural Language Processing (RANLP-2003)*, Borovets, Bulgaria.
- Banerjee, Satanjeev and Ted Pedersen. 2003. Extended gloss overlaps as a measure of semantic relatedness. In *Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence (IJCAI-03)*, pages 805–810, Acapulco, Mexico.
- Budanitsky, Alexander and Graeme Hirst. 2006. Evaluating WordNet-based measures of semantic distance. *Computational Linguistics*, 32(1):13–47.
- Church, Kenneth W. and Patrick Hanks. 1990. Word association norms, mutual information and lexicography. *Computational Linguistics*, 16(1):22–29.
- Cruse, D. Allen. 1986. *Lexical semantics*. Cambridge University Press, Cambridge, UK.
- Curran, James R. 2004. *From Distributional to Semantic Similarity*. Ph.D. thesis, School of Informatics, University of Edinburgh, Edinburgh, UK.
- Dagan, Ido, Lillian Lee, and Fernando Pereira. 1994. Similarity-based estimation of word cooccurrence probabilities. In *Proceedings of the 32nd Annual Meeting of the Association of Computational Linguistics (ACL-1994)*, pages 272–278, Las Cruces, New Mexico.
- Dagan, Ido, Lillian Lee, and Fernando Pereira. 1997. Similarity-based methods for word sense disambiguation. In *Proceedings of the 35th Annual Meeting of the Association of Computational Linguistics (ACL-1997)*, pages 56–63, Madrid, Spain.
- Dagan, Ido, Lillian Lee, and Fernando Pereira. 1999. Similarity-based models of cooccurrence probabilities. *Machine Learning*, 34(1–3):43–69.
- Dagan, Ido, Shaul Marcus, and Shaul Markovitch. 1995. Contextual word similarity and estimation from sparse data. *Computer Speech and Language*, 9:123–152.
- Firth, John R. 1957. A synopsis of linguistic theory 1930–55. In *Studies in Linguistic Analysis (special volume of the Philological Society)*, pages 1–32, Oxford, England. The Philological Society.
- Gorman, James and James R. Curran. 2006. Scaling distributional similarity to large corpora. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 361–368, Sydney, Australia.
- Grefenstette, Gregory. 1992. Sextant: Exploring unexplored contexts for semantic extraction from syntactic analysis. In *Proceedings of the 30th Annual Meeting of the Association for Computational Linguistics (ACL-1992)*, pages 324–326, Newark, Delaware.
- Gurevych, Iryna. 2005. Using the structure of a conceptual network in computing semantic relatedness. In *Proceedings of the 2nd International Joint Conference on Natural Language Processing (IJCNLP-2005)*, pages 767–778, Jeju Island, Republic of Korea.
- Halliday, Michael A. K. 1966. Lexis as a linguistic level. In J.C. Catford C.E. Bazell, M.A.K Halliday, and R.H. Robins, editors, *In memory of J.R. Firth*, pages 148–162, London, UK. Longmans Linguistics Library.
- Harman, Donna. 1993. Overview of the first text retrieval conference. In *Proceedings of the 16th Annual International Association for Computing Machinery Special Interest Group on Information Retrieval (ACM-SIGIR) conference on Research and Development in information retrieval*, pages 36–47, Pittsburgh, Pennsylvania.
- Harris, Zellig. 1968. *Mathematical Structures of Language*. Interscience Publishers, New York, NY.
- Hindle, Donald. 1990. Noun classification from predicate-argument structures. In *Proceedings of the 28th Annual Meeting of the Association of Computational Linguistics (ACL-1990)*, pages 268–275, Pittsburgh, Pennsylvania.
- Hirst, Graeme and Alexander Budanitsky. 2005. Correcting real-word spelling errors by restoring lexical cohesion. *Natural Language Engineering*, 11(1):87–111.
- Hirst, Graeme and David St-Onge. 1998. Lexical chains as representations of context for the detection and correction of malapropisms. In Christiane Fellbaum, editor, *WordNet: An Electronic Lexical Database*. The MIT Press, Cambridge, MA,

- chapter 13, pages 305–332.
- Inkpen, Diana and Graeme Hirst. 2002. Acquiring collocations for lexical choice between near-synonyms. In *SIGLEX Workshop on Unsupervised Lexical Acquisition, 40th meeting of the Association for Computational Linguistics*, Philadelphia, PA.
- Jarmasz, Mario and Stan Szpakowicz. 2003. Roget's Thesaurus and semantic similarity. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP-2003)*, pages 212–219, Borovets, Bulgaria.
- Jiang, Jay J. and David W. Conrath. 1997. Semantic similarity based on corpus statistics and lexical taxonomy. In *Proceedings of International Conference on Research on Computational Linguistics (ROCLING X)*, Taipei, Taiwan.
- Landauer, Thomas K., Peter W. Foltz, and Darrell Laham. 1998. Introduction to latent semantic analysis. *Discourse Processes*, 25(2–3):259–284.
- Leacock, Claudia and Martin Chodorow. 1998. Combining local context and WordNet similarity for word sense identification. In Christiane Fellbaum, editor, *WordNet: An Electronic Lexical Database*. The MIT Press, Cambridge, MA, chapter 11, pages 265–283.
- Lee, Lillian. 1999. Measures of distributional similarity. In *Proceedings of the 37th conference on Association for Computational Linguistics (ACL-1999)*, pages 25–32, College Park, Maryland.
- Lee, Lillian. 2001. On the effectiveness of the skew divergence for statistical language analysis. In *Proceedings of the Eighth International Workshop on Artificial Intelligence and Statistics (AISTATS-2001)*, pages 65–72, Key West, Florida.
- Lesk, Michael. 1986. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *Proceedings of the 5th Annual International Conference on Systems Documentation*, pages 24–26, Toronto, Canada.
- Lin, Dekang. 1997. Using syntactic dependency as local context to resolve word sense ambiguity. In *Proceedings of the 8th Conference of the European Chapter of the Association for Computational Linguistics (ACL, EACL-1997)*, pages 64–71, Madrid, Spain.
- Lin, Dekang. 1998a. Automatic retrieval and clustering of similar words. In *Proceedings of the 17th International Conference on Computational Linguistics (COLING-1998)*, pages 768–773, Montreal, Canada.
- Lin, Dekang. 1998b. Automatic retrieval and clustering of similar words. In *Proceedings of the 17th International Conference on Computational Linguistics (COLING-1998)*, pages 768–773, Montreal, Canada.
- Lin, Dekang, Shaojun Zhao, Lijuan Qin, and Ming Zhou. 2003. Identifying synonyms among distributionally similar words. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence (IJCAI-95)*, pages 1492–1493, Acapulco, Mexico.
- McCarthy, Diana, Rob Koeling, Julie Weeds, and John Carroll. 2007. Unsupervised acquisition of predominant word senses. *Computational linguistics*, 33(4).
- Miller, George A. and Walter G. Charles. 1991. Contextual correlates of semantic similarity. *Language and Cognitive Processes*, 6(1):1–28.
- Mirkin, Shachar, Ido Dagan, and Maayan Geffet. 2007. Integrating pattern-based and distributional similarity methods for lexical entailment acquisition. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics Main Conference Poster Sessions*, pages 579–586, Sydney, Australia.
- Mohammad, Saif, Iryna Gurevych, Graeme Hirst, and Torsten Zesch. 2007. Cross-lingual distributional profiles of concepts for measuring semantic distance. In *Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP/CoNLL-2007)*, pages 571–580, Prague, Czech Republic.
- Mohammad, Saif and Graeme Hirst. 2006a. Determining word sense dominance using a thesaurus. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 121–128, Trento, Italy.
- Mohammad, Saif and Graeme Hirst. 2006b. Distributional measures of concept-distance: A task-oriented evaluation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP-2006)*, pages 35–43, Sydney, Australia.
- Morris, Jane and Graeme Hirst. 2004. Non-classical lexical semantic relations. In *Proceedings of the Workshop on Computational Lexical Semantics, Human Language Technology Conference of the North*

- American Chapter of the Association for Computational Linguistics*, pages 46–51, Boston, Massachusetts.
- Patwardhan, Siddharth, Satanjeev Banerjee, and Ted Pedersen. 2003. Using measures of semantic relatedness for word sense disambiguation. In *Proceedings of the Fourth International Conference on Intelligent Text Processing and Computational Linguistics (CICLING-03)*, pages 17–21, Mexico City, Mexico.
- Patwardhan, Siddharth and Ted Pedersen. 2006. Using WordNet based context vectors to estimate the semantic relatedness of concepts. In *Proceedings of the European Chapter of the Association for Computational Linguistics Workshop Making Sense of Sense—Bringing Computational Linguistics and Psycholinguistics Together*, pages 1–8, Trento, Italy.
- Paul, Douglas B. 1991. Experience with a stack decoder-based HMM CSR and back-off n-gram language models. In *Proceedings of the Speech and Natural Language Workshop*, pages 284–288, Palo Alto, California.
- Pereira, Fernando, Naftali Tishby, and Lillian Lee. 1993. Distributional clustering of English words. In *Proceedings of the 31st Annual Meeting of the Association of Computational Linguistics (ACL-1993)*, pages 183–190, Columbus, Ohio.
- Rada, Roy, Hafedh Mili, Ellen Bicknell, and Maria Blettner. 1989. Development and application of a metric on semantic nets. *IEEE Transactions on Systems, Man, and Cybernetics*, 19(1):17–30.
- Rapp, Reinhard. 2003. Word sense discovery based on sense descriptor dissimilarity. In *Proceedings of the Machine Translation Summit IX*, pages 315–322, New Orleans, Louisiana.
- Resnik, Philip. 1995. Using information content to evaluate semantic similarity. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI-95)*, pages 448–453, Montreal, Canada.
- Resnik, Philip. 1998. WordNet and class-based probabilities. In Christiane Fellbaum, editor, *WordNet: An Electronic Lexical Database*. The MIT Press, Cambridge, Massachusetts, pages 239–263.
- Resnik, Philip. 1999. Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language. *Communications of the Association of Computing Machinery*, 11:95–130.
- Resnik, Philip and Mona Diab. 2000. Measuring verb similarity. In *Proceedings of the 22nd Annual Meeting of the Cognitive Science Society (CogSci 2000)*, pages 399–404, Philadelphia, Pennsylvania.
- Rubenstein, Herbert and John B. Goodenough. 1965. Contextual correlates of synonymy. *Communications of the Association of Computing Machinery*, 8(10):627–633.
- Schütze, Hinrich and Jan O. Pedersen. 1997. A cooccurrence-based thesaurus and two applications to information retrieval. *Information Processing and Management*, 33(3):307–318.
- Turney, Peter. 2001. Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. In *Proceedings of the Twelfth European Conference on Machine Learning (ECML-2001)*, pages 491–502, Freiburg, Germany.
- Weeds, Julie, David Weir, and Diana McCarthy. 2004. Characterising measures of lexical distributional similarity. In *Proceedings of the 20th International Conference on Computational Linguistics (COLING-04)*, pages 1015–1021, Geneva, Switzerland.
- Weeds, Julie E. 2003. *Measures and Applications of Lexical Distributional Similarity*. Ph.D. thesis, Department of Informatics, University of Sussex, Brighton, UK.
- Yang, Dongqiang and David Powers. 2005. Measuring semantic similarity in the taxonomy of WordNet. In *Proceedings of the Twenty-eighth Australasian Computer Science Conference (ACSC-2005)*, pages 315–322, Newcastle, Australia.
- Yang, Dongqiang and David Powers. 2006. Verb similarity on the taxonomy of WordNet. In *Proceedings of the Third International WordNet Conference (GWC-2006)*, pages 121–128, Jeju Island, Republic of Korea.
- Yoshida, Sen, Takashi Yukawa, and Kazuhiro Kuwabara. 2003. Constructing and examining personalized cooccurrence-based thesauri on web pages. In *Proceedings of the 12th International World Wide Web Conference*, pages 20–24, Budapest, Hungary.
- Zesch, Torsten, Iryna Gurevych, and Max Mühlhäuser. 2007. Comparing Wikipedia and German WordNet by evaluating semantic relatedness on multiple datasets. In *Proceedings of Human Language Technologies: The Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL/HLT-2007)*, pages 205–208, Rochester, New York.