



Archived at the Flinders Academic Commons:

<http://dspace.flinders.edu.au/dspace/>

This is a copy of a paper presented at the Third International WordNet Conference (GWC-06), South Jeju Island, Korea.

Please cite this as: Yang, D. & Powers, D.M., 2006. Verb similarity on the taxonomy of WordNet. In P. Sojka, K.-S. Choi, C. Fellbaum and P. Vossen (Eds.) Proceedings of the Third International WordNet Conference GWC 2006, South Jeju Island, Korea. Masaryk University, Brno, pp. 121-128.

Copyright © Masaryk University, 2005.

The abstracts for the proceedings are available at:

<http://nlp.fi.muni.cz/gwc2006/proc/index.html>

And the published version of the paper is available at:

<http://semanticweb.kaist.ac.kr/conference/gwc/pdf2006/2.pdf>

Paper reproduced here in accordance with the publishers' copyright policy.

# Verb Similarity on the Taxonomy of WordNet

Dongqiang Yang and David M. W. Powers  
School of Informatics and Engineering  
Flinders University of South Australia  
PO Box 2100, Adelaide 5001, South Australia  
Dongqiang.Yang@flinders.edu.au  
David.Powers@flinders.edu.au

## Abstract

In this paper, we introduce two kinds of word similarity algorithms, SHE and RHE, to investigate the capability of WordNet in measuring verb similarity. In the absence of a standard verb set we have proposed two new verb similarity evaluation data sets.

## 1 Introduction

Many researchers have explored the similarity of nouns using a variety of methods including methods based on WordNet. However, little attention has been paid to verbs (Resnik and Diab, 2000), there is no standard evaluation set, and it is not clear that the WordNet verb hierarchy is rich enough to support verb similarity assessment.

This paper seeks to extend work done on nouns by Yang and Powers (2005) to verbs and uses the noun performance as a benchmark level for the work on verbs. In this study we introduce a verb evaluation set with both tuning and evaluation partitions, we present and adapt a successful noun similarity method based on WordNet to the verb similarity task, and we present a hybrid technique that seeks to increase accuracy by cross mapping into the noun hierarchy and back.

Measuring word similarity can be classified into knowledge-rich and knowledge-poor methods (Grefenstette, 1993; Gasperin et al., 2001). Here the knowledge refers to acquiring lexicon oriented information from a pre-handcrafted thesaurus or from learning from a corpus. We introduce both approaches before presenting our own results using knowledge-rich methods.

### 1.1 Knowledge-poor methods

Knowledge-poor methods mainly depend on information or probability information derived from a corpus or the Internet (Turney, 2001) rather than a knowledge base. Such methods may be further categorized according to how co-occurrence frequency data is handled:

#### 1.1.1 Vector space

These approaches assume that semantically related words are more likely to co-occur in the corpus. A matrix is constructed in word-by-word or word-by-document order with a cell value such as term frequency ( $TF$ ) or  $TF \times IDF$  (inverse document frequency, but more accurately the information conveyed by the fact of occurrence in a document).

Word similarity is established by comparing distance measures such as the cosine coefficient or Euclidean distance (Schütze, 1992).

#### 1.1.2 Syntactic dependency

These approaches assume that the semantic relatedness of words leads to their use in similar grammatical structures. Judging word similarity is achieved by tagging parts-of-speech in the corpus, shallow parsing of sentences, specifying the relationship between chunks and comparing the syntactic components along with their dependency relations (Grefenstette, 1993).

### 1.2 Knowledge-rich approaches

Knowledge-rich methods require semantic networks or a semantically tagged corpus to define the concept of word in the relation with other concepts or to other words in the surrounding context. Most methods that calculate semantic distance using ontology or thesaurus knowledge, such as WordNet (Miller, 1995; Fellbaum, 1998) or Roget's thesaurus (Jarmasz and Szpakowicz, 2003) fall into this category. The popular methodologies for measuring semantic relatedness with the help of a thesaurus can be classified into two categories: one uses solely semantic links (i.e. edge-counting), the other combines corpus statistics with taxonomic distance.

#### 1.2.1 Edge-counting

The edge-counting or shortest path method derives from the geometric model in Cognitive Psychology, where the shorter distance entails the stronger association between stimuli and response. It can be traced back to Quillian's semantic memory model (Quillian, 1967; Collins and Quillian, 1969) where concept nodes are planted within the hierarchical network and the number of hops between the nodes specifies the similarity of the concepts. Generally the similarity of words in the thesaurus space can be described as:

$$Sim(i, j) = 2D - Dist(i, j) \quad (1)$$

where  $D$  is a constant (e.g. the maximum depth in the taxonomy of WordNet, viz. 16 if we presume all the hierarchies have a common node),  $Dist(i, j)$  is the number of links between two concept nodes  $i$  and  $j$ . In the edge-counting methods distance is typically assessed by counting the edges traversed from  $c1$  to  $c2$  via  $ncn$  (the nearest common node),  $Dist(c1, c2)$  – we will introduce a few popular

edge-counting models working in the semantic hierarchy (cf. Pedersen et al., 2003).

Wu and Palmer (1994) proposed to measure the verbal concept similarity in the projected domain hierarchy when translating from English verbs to Chinese. According to Wu and Palmer, the relatedness of two words is the weighted sum of all their senses comparison,

$$Sim(v_i, v_j) = \sum_k w_k \times \frac{2 \times dep(ncn(c_{i,k}, c_{j,k}))}{dep(c_{i,k}) + dep(c_{j,k})} \quad (2)$$

where  $ncn(c_{i,k}, c_{j,k})$  is the nearest common node ( $ncn$ ) for the conceptual nodes  $c_{i,k}, c_{j,k}$  of verbs  $v_i$  and  $v_j$ ,  $dep$  is the depth of the node relative to the root,  $w_k$  is the weight of each pair of concepts in each domain. The sum of  $w_k$  is 1. This model is appropriate for measuring both verbs and nouns in the "IS-A" hierarchical concept net.

Leacock and Chodorow (1998) adapted the concept of information content (Resnik, 1995) to evaluate the relatedness of two words using the following model:

$$\begin{aligned} Sim(W_i, W_j) &= Max \left[ -\log \frac{Dist(c_i, c_j)}{2 \times D} \right] \quad (3) \\ &= Max \left[ \log 2D - \log Dist(c_i, c_j) \right] \end{aligned}$$

where  $Dist(c_i, c_j)$  is the shortest distance between concepts  $c_i$  and  $c_j$ . In addition, they defined the similarity of two words as the maximized value of all the pairwise similarities. Note that in Equation (3)

$$\begin{aligned} Dist(c_i, c_j) &= dep(c_i) + dep(c_j) - 2dep(ncn(c_i, c_j)) \quad (4) \\ Sim(W_i, W_j) &= Max \left[ \log \frac{2D}{Dist(c_i, c_j)} \right] \quad (5) \end{aligned}$$

Hence, the concept model is similar to Wu and Palmer's apart from the  $\log$  normalization.

### 1.2.2 Information Content

Resnik (1995) argues that the links in the hierarchy of WordNet representing a uniform distance in the edge-counting measurement can not account for the semantic variability of a single link. He defines information content of  $ncn$  to explain the similarity of two words through frequency statistics retrieved from a corpus, not through the distance of edge-counting. Here the frequency of  $ncn$  subsumes all the frequency data of subordinate concept nodes. The information content can be quantified as the negative of the log likelihood,  $-\log P(c)$ .

However, Resnik still employs the structure of a conceptual net and one drawback is that the  $ncn$  for all concept pairs that have the same parent node is the same.

Building on Resnik's work, Jiang and Conrath (1997) further assumed that a combination of information content and edge-counting will improve the correlation co-efficient (compared with human judgment). They also considered the link type, depth, conceptual density, and information content of concepts. Their simplified formula can be expressed as

follows:

$$Dist(c_i, c_j) = IC(c_i) + IC(c_j) - 2 \times IC(ncn(c_i, c_j)) \quad (6)$$

$$Sim(c_i, c_j) = -Dist(c_i, c_j) \quad (7)$$

Lin (1997) introduced another way of computing the similarity to disambiguate word sense,

$$Sim(c_i, c_j) = \frac{2 \times IC(ncn(c_i, c_j))}{IC(c_i) + IC(c_j)} \quad (8)$$

which is essentially another normalized form of Jing and Conrad's model.

## 2 Multiplicative Models

### 2.1 The noun model

Generally speaking, similarity models in the taxonomy of WordNet, proposed by Wu and Palmer, Leacock and Chodorow, Jiang and Conrath, and Lin, can be abstracted into one of the following forms:

$$Sim(c1, c2) = 2\gamma \div (\alpha + \beta) \quad (9)$$

$$Sim(c1, c2) = 2\gamma - (\alpha + \beta) \quad (10)$$

where  $\alpha, \beta, \gamma$ , respectively denote attributes of concepts  $c1, c2$ , and the  $ncn$  of  $c1, c2$  in the "IS-A" hierarchy. The attribute can be viewed as some function of the depth in the taxonomy or the information content extracted from the outer corpus.

Yang and Powers (2005) proposed a new model to measure semantic similarity in the taxonomy of WordNet, based on a variation of edge-counting. In contrast with the above methods they also take into account the part-whole (hol/meronym) relationships in WordNet and compare two searching algorithms, a bidirectional depth-limit search (**BDLS**) and unidirectional breadth-first search (**UBFS**).

On the assumption that a single link in the taxonomy always stands for the same depth-independent distance and that the distance between two conceptual nodes is the least number of links,  $\lambda$ , from one node to another, they define the similarity of two concepts multiplicatively as,

$$Sim(c1, c2) = \alpha_t * \beta^\lambda \quad (11)$$

Partially inspired by Hirst and St. Onge's algorithm (1995) for the detection and correction of malapropisms with different weights for identical words, synonyms or antonyms, and hyper/hyponyms, Yang and Powers deal with the identity case where  $c1$  and  $c2$  are identical as  $\alpha_{id} = 1, \gamma = 0$ , the syn/antonym as an intermediate weight,  $\alpha_{sa} = 0.9, \gamma = 0$ , assigning the lowest weight (e.g.  $\alpha = \alpha_{hh} = \alpha_{hm} = 0.85, \beta = \beta_{hh} = \beta_{hm} = 0.7$ ) for the hyper/hyponym, hol/meronym where searching depth  $\gamma$  is more than one – these weights being the result of tuning noun similarity.

These models are evaluated against a benchmark set by human similarity judgment, and achieve a much improved result compared with other methods: the correlation with

average human judgment on a standard 28 noun pair dataset (Resnik, 1995) is 0.921, which is better than anything reported in the literature and also significantly better than average individual human judgments. As this set has been effectively used for algorithm selection and tuning, they also validate on an independent 37 noun pair test set (0.876) and present cross-validated results for the full 65 noun-pair superset (0.897) (Rubenstein and Goodenough, 1965). Note that their best performance on these data sets is achieved for the maximum score across distinct senses in relation to the common case of words that are polysemous.

## 2.2 A multistrategy verb model

To investigate the appropriateness of such a model for judging word similarity we have sought to adapt it to apply to verbs, which are another significant hierarchy in WordNet. Unlike the noun taxonomy, which is rich in complexity and links, the verbs are organized into a relatively shallow hierarchy according to their hyper/troponymy relations and WordNet does not represent holo/metonymy relations. The maximum distance between contentive verbs (excluding stopwords like ‘be’, ‘make’ and ‘do’) is around 4 nodes, which makes it more difficult to find relationships between verbs (Fellbaum, 1998). Based on the Yang and Powers noun model and approach, we designed and tuned a new algorithm to account for the similarity of verbs in the face of the sparseness and limitations of the WordNet verb hierarchy. To supplement the verb hierarchy, we also considered derivational mapping into the noun hierarchy, the use of definitions (glosses), and the effect of stemming. Thus we consider the following factors in constructing this model of verb similarity, where at this stage stemming refers only to the simple suffix removal functions provided with WordNet2.

1. Similarity on the verb taxonomy is evaluated in the same basic way as for the noun hierarchy, viz. equation (11) and (12), except that there is no correlate of the holo/meronym relationships (viz. no metonymy by which a part of an action/scene may be related to the whole). We thus need to set up and tune parameters for the syno/antonyms and hyper/troponyms in the same way as with the noun model.
2. Some verbs have the *noun* form as a *stem*, or vice versa, as they are *derivationally* related. Thus we can project to the noun hierarchy from the verb hierarchy to enrich the relationships among verbs, introducing  $\alpha_{der}$  as a discount factor or fusion weight.
3. The definition of a verb, its *gloss*, can give a hint to the relation with other verbs *when there are no apparent linkages in the verb and noun hierarchies*. Lesk (1986) proposed calculating the overlaps of target word and other words in the context in the definitions to select an appropriate sense. Pedersen et al. (2003) treat the definitions in WordNet as a million word corpus, and build a co-occurrence matrix to specify how many times two concepts turn up together in the gloss of WordNet. In this paper we assume verbs

in the definition of WordNet, which are not in the frequent word list like "make", "do", etc., bring about a strong semantic relation with its target word. This thus introduces weight  $\alpha_{gls}$ .

4. The *stemming* effect seen above can also connect related verbs in the verb hierarchy without considering their individual senses. Rather, it allows us to capture a wider class of relationship that relate to the etymology of the word and its root meaning, but should not represent as strong a relationship as those that are represented directly by links. This gives us weight  $\alpha_{stm}$ .

Comprehensively considering these new factors and the existing link type and depth factors that we need to tune for the WordNet verb taxonomy, and noting that Yang and Powers have already well tuned for noun similarity and we need no adjustments, the new model is:

$$Sim(c1, c2) = \alpha_{stm} \alpha_t \prod_{i=1}^{Dist(c1, c2)} \beta_t Dist(c1, c2) < \gamma, \quad (12)$$

$$Sim(c1, c2) = 0, Dist(c1, c2) \geq \gamma;$$

$$Sim_{max}(v1, v2) = Max_{(i, j)} [Sim(c_{1,i}, c_{2,j})] \quad , \quad (13)$$

- $c1, c2$  represent concept nodes

where  $0 \leq Sim(c1, c2) \leq 1$ ,

- $t = ht$  (hyper/troponym),  $sa$  (syn/antonym),  $der$  (derived nouns) or  $gls$  (definition),
- $\alpha_t$  is a link type factor applied to a sequence of links of type  $t$ . ( $0 < \alpha_t \leq 1$ ),
- $\alpha_{stm}$  is the stemming factor, if  $c1$  links to  $c2$  without stemming,  $\alpha_{stm} = 1$
- $\beta_t$  is the depth factor depending on the link type
- $\gamma$  is an arbitrary threshold on the distance, which will no more than five in the verb taxonomy
- $Dist(c1, c2)$  is the distance (the shortest path) between  $c1$  and  $c2$

The most strongly related concepts are the identity case where  $c1$  and  $c2$  are identical,  $\alpha_{id} = 1$  and  $Dist(c1, c2) = 0$ . For the link type of syn/antonym, we again assign an intermediate weight (e.g.  $\alpha_{sa} = 0.9$ ,  $Dist(c1, c2) = 0$ ), and we again tune to assign the lowest weight (e.g.  $\alpha_{ht} = 0.85$ ) for hyper/troponymy. Note that any syn/antonym and identity links constitute entire paths and cannot be part of a multilink path.

Given the fact that most verbs are polysemous we will again assign the maximum value of the similarity among all the  $n_i$  senses  $c_{i,j}$  of any polysemous word  $v_i$ . To make clear the final model of verb similarity in the WordNet we present it succinctly but informally as the following algorithm. The bidirectional search is as described in the original Yang and Powers algorithm (2005), deciding first if it is a direct

identity or synonym path, or otherwise discounting it if it is a hyper/tropo path and calculating the additional distance required to connect them, except that if it is unsuccessful, it is redone with a further discount allowing a connection through any derivationally related stem, not just through specific senses.

The basic algorithm is as follows, where the noun similarity and maximum similarity steps are exactly as described by Yang and Powers:

```

for each sense c1 and c2 of v1 and v2 resp.
  if c1 and c2 are synonymous or antonymous
    assign sim_sa(c1,c2)=  $\alpha_{sa}$ ; Goto next loop
  elsif c1 and c2 are hyper- tropo- and/or antonym connected
    with depth d less than  $\gamma$ 
    sim(c1,c2) = sim_hta(c1,c2)=  $\alpha_{ht}$  *  $\beta_{ht}^d$ 
  if=0 & c1 and c2 are stem hyper/tropo/antonym connected
    with depth d less than  $\gamma$ 
    sim(c1,c2) = sim_stm(c1,c2)=  $\alpha_{stm}$  *  $\alpha_{ht}$  *  $\beta_{ht}^d$ 
  endif
endif
endfor
calculate the maximum similarity score,
sim_max(c1∈v1, c2∈v2)
if=0
  sim(v1,v2) = sim_max(c1∈v1, c2∈v2)
elsif v1 can find v2 in its definition or vice versa
  sim(v1,v2) = sim_gls(v1,v2)=  $\alpha_{gls}$ 
else
  if both v1 and v2 have derived noun form
    go into noun taxonomy and perform BLS search:
    sim(v1,v2) = sim_der(c1,c2)=  $\alpha_{der}$  * sim_noun(c1,c2)
  endif
endif

```

### 3 Evaluation

#### 3.1 Task

Unfortunately, there is no benchmark data set for verbs in the literature. We have thus had to make our own data set and offer it as a standard for testing verb similarity. We selected 20 verb synonym tests from the 80 TOEFL<sup>1</sup> (Test of English as a Foreign Language) questions used by (Landauer and Dumais, 1997), and 16 from a set of 50 ESL (English as a second language) questions (Tatsuki, 1998) – these are widely used to assess non-native eligibility for university entry or employment in English speaking countries and we judged them as representing different levels of difficulty for non-native speakers, but as all well within the competence of a native speaker or university graduate in an English speaking country. Each of these 36 multiple choice questions consists of a question or target word and four other words or phrases to choose from. We managed to select examples with words rather than phrases, and then used each target word together with one of the four choices to construct a pair of verbs in the questionnaire, giving a total of 144 pairs of verbs.

We randomly arranged these word pairs and randomly reversed the order of target verb and choice verb. Six colleagues (2 academic staff and 4 postgraduate students) voluntarily rated these pairs for similarity. Four of them are native speakers of Australian English; the other two are near-native speakers who have used English as a primary language and a main communication tool (at high school, at

<sup>1</sup>Test of English as a Foreign Language (TOEFL), Educational Testing Service, Princeton, New Jersey, <http://www.ets.org/>

university and in everyday life) for over ten years. We gave them the following instructions:

*Indicate how strongly these words are related in meaning using integers from 0 to 4. The following are given as examples of kinds of descriptions that might apply to each number, but you must give your own judgement and if you think something falls in between two of these categories you must push it up or down (no halves or decimals).*

- 0: not at all related
- 1: vaguely related
- 2: indirectly related
- 3: strongly related
- 4: inseparably related

The word pairs were sorted in descending order of average score, and divided up to achieve a balanced set with 26 words in each category (eliminating some words with averages below 2 to eliminate an expected imbalance due to the questions being designed to have exactly one best answer and being biased to include more dissimilar words). We then randomly assigned 13 words from each category to one of two data sets, data1 and data2. The average correlation among these six subjects was  $r = 0.866$ .

We next optimized the verb model for each data set through calculating the correlation with average human scores, using a greedy approach to optimizing the parameters (choosing the mid-value when there was no significant difference). Here we show how we regulated the verb model on data1.

To distinguish the different effect of each factor we proposed, we assumed the contribution of the verb hierarchy similarity, derived noun hierarchy similarity and gloss similarity are independent. Thus we first sought the optimal parameterization for the verb hierarchy, and then to set  $\alpha_{der}$  and  $\alpha_{gls}$  considered how helpful the derived noun similarity was and then how helpful the gloss similarity was.

#### 3.2 Tuning

There were three parameters we needed to adjust in relation to the application of the Yang and Powers algorithm to the verb hierarchy, the path type factor  $\alpha_t$ , the link type factor  $\beta$  and the depth factor  $\gamma$  (optional, noting that this last factor was originally and primarily conceived to minimize CPU time, but may also serve as a threshold to stop relationships that are too strained being discovered). Then in order to factor in the alternative source of information we needed to set the stem similarity weighting  $\alpha_{stm}$ , the derived noun similarity weighting  $\alpha_{der}$ , and the gloss similarity weighting  $\alpha_{gls}$ . In this case the three values are fallback weights: given the algorithm for the verb hierarchy hasn't given us a non-zero value, we retry, ignoring sense and inflectional variations of verbs (discounted using  $\alpha_{stm}$ ), and if it is still non-zero, we use the noun version algorithm to seek a value for derivationally related nouns (discounted by  $\alpha_{der}$ ), or failing that we try to find a connection via the glosses ( $\alpha_{gls}$ ).

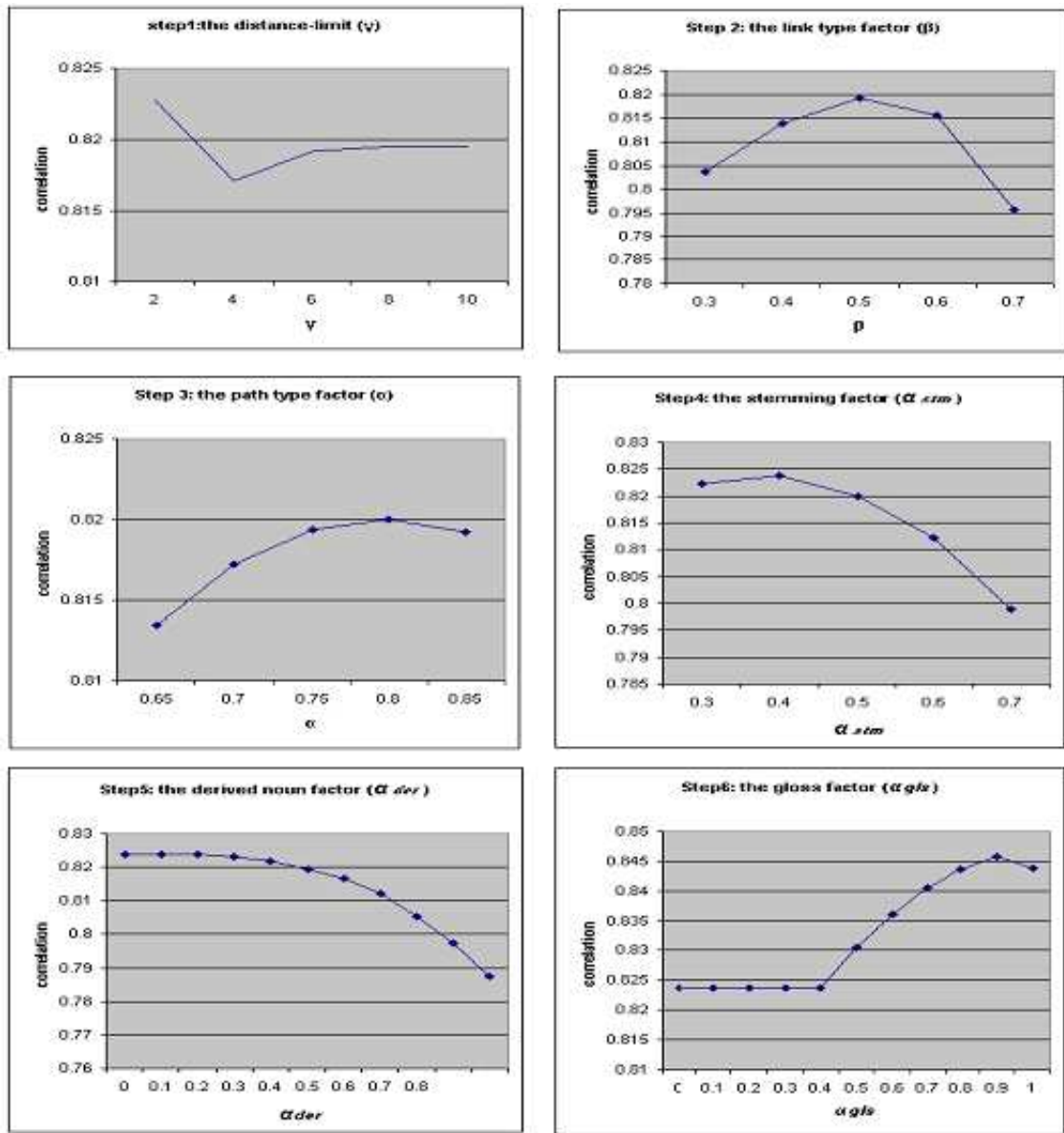


Figure 1: The tuning process on the RHE

### 3.3 Step 1: the distance limit ( $\gamma$ )

Once the values of  $\alpha$ ,  $\alpha_{stem}$  and  $\beta$  had been assigned initially, i.e. respectively 0.85, 0.5 and 0.5, we varied the distance-limit  $\gamma$  (for the combined path length), enlarging the search distance of each node from 1 to 5 (essentially the maximum distance is no more than 5 in WordNet), viz. the total distance of two nodes in the BDLS ranged from 2 to 10, to investigate if by expanding the distance-limit, the model could produce a judgment that is more accurate. We can see in Figure 1( $\gamma$ ) that there is a drop in the correlation when we increase the searching scope from 1 level to 2 level, after that the curve approached level. Our purpose in the paper is to investigate the function of verb hierarchy, so we use  $\gamma = 6$

for a rich hierarchy exploration (RHE) but also use  $\gamma = 2$  as a reference point for shallow hierarchy exploration (SHE). In the following part we just illustrate how to calibrate the model using the RHE variant.

### 3.4 Step 2: the link type factor ( $\beta$ )

We tested  $\beta$  over the range 0.3 to 0.7, tuning by increments of 0.1, to see if it affected the correlation with human judgment. Note that each link in the taxonomy is of uniform distance if we give  $\beta = 1$ . In fact, we see from Figure 1( $\beta$ ) that the performance of the system begins to deteriorate significantly for  $\beta$  bigger than 0.6 with the maximum at 0.5.

### 3.5 Step 3: the path type factor ( $\alpha$ )

We varied the value of  $\alpha$ , by increments of 0.05 from 0.5 to 0.95. The optimal value for  $\alpha$  is around 0.8 but there is very little sensitivity to its precise value as seen in in Figure 1( $\alpha$ ).

### 3.6 Step 4: the stemming factor ( $\alpha_{stm}$ )

After the optimal value, 0.4, Figure 1( $\alpha_{stm}$ ) shows that the correlation begins to drop quickly but prior to that there is little change.

### 3.7 Step 5: the derived noun factor ( $\alpha_{der}$ )

Similarly, there is little difference as  $\alpha_{der}$  increase from 0 to 0.5, but after that the correlation deteriorated slowly – see Figure 1( $\alpha_{der}$ ). We chose 0.4 as a compromise value, as with the shallower verb hierarchy we did expect to see smaller values, but a larger value will maximize utilization of the information in the network.

### 3.8 Step 6: the gloss factor ( $\alpha_{gls}$ )

There is an initial jump at 0.4, rising to a clear optimum at 0.9, as seen in Figure 1( $\alpha_{gls}$ ).

## 4 Results

Table 1: The final result on the each 65 data sets and the total dataset. (r\_t: the correlation on the tuning set, r\_e: the correlation on the evaluation set, where data1 is the evaluation set for data2, and vice versa.)

		$\gamma$	$\beta$	$\alpha$	$\alpha_{stm}$	$\alpha_{der}$	$\alpha_{gls}$	r t	r e
R	Data1(65)	2	0.5	0.8	0.4	0.1	0.9	0.846	0.775
H	Data2(65)	2	0.2	0.85	0.7	0.8	0.5	0.864	0.823
E	Total (130)	2	0.5	0.8	0.5	0.75	0.6	0.808	
S	Data1 (65)	0	0.6	0.75	0.4	0.7	0.9	0.838	0.824
H	Data2 (65)	0	0.4	0.8	0.6	0.7	0.5	0.846	0.835
E	Total (130)	0	0.5	0.8	0.5	0.75	0.6	0.833	

After we had tuned the verb model on each data set we found the selected values did not correspond very well with each other, reducing the score for the 2-fold cross validation. This was not unexpected due to the relative flatness (lack of significant difference) for much of the curves, which forced an arbitrary selection within a range. Unfortunately the tuning is a time intensive process, so we have not yet been able to perform a higher order cross validation. Owing to the sensitivity of each data set as measured by the correlation, r, to tuning on the other, we adopted a compromise tuning based on both subsets for future comparison against human performance, noting that apart from the Yang and Powers paper where identical results were achieved for each mode of the cross-validation, results for work on noun similarity do *not* do tuning and validation on separate subsets of the data. Table 1 shows the final parameters and correlations with the average human scores for both **RHE** and **SHE**. There is little difference on the final verb model due to the choice of **RHE** or **SHE**.

## 5 Discussion

The Yang and Powers noun similarity study advocated the Wilcoxon Signed Rank Test as a principled non-parametric

Table 2: Significance test on both RHE and SHE, r\_a: the correlation with average human,  $\sigma$ : standard deviation,  $\mu$ : mean, sig: significance

	r a	$\sigma/\mu$	RHE		SHE	
			z-score	sig	z-score	sig
subject1	0.88	0.292	-3.25	0.001	-2.113	0.035
subject2	0.733	0.45	0	1	-0.802	0.423
subject3	0.878	0.488	-3.07	0.002	-3.421	<0.001
subject4	0.926	0.485	-3.52	<0.001	-1.14	0.254
subject5	0.913	0.397	-4.47	<0.001	-3.396	<0.001
subject6	0.868	0.402	-1.89	0.059	-1.61	0.107
RHE	0.808	0.308	0	1	-1.484	0.138
SHE	0.833	0.561	-1.484	0.138	0	1

modification to the two-sample t test for comparing their results against human judgment. We in the same way performed this test (at 95 percent level) for the present verb similarity study, achieving the results listed in the Table 2. The choice of **RHE** versus **SHE** makes no significant difference in the ability of judging verb similarity, and they are only significantly better than one subject (a non-native speaker). However, three other subjects fail to do significantly better than **SHE** (shallow), whilst just one just misses out on being significantly better than **RHE** (rich), although all their judgments retain a high correlation with the average human. Thus while there is no significant difference between the rich and shallow variants themselves with respect to the group, the richer variant doesn't keep step with individual human subjects as well as the shallower variant, implying that the additional levels of the verb hierarchy are less useful in modeling human behavior than the gloss derived noun fallbacks we have introduced.

## 6 Conclusions and Future Work

The maximum depth in the verb model is much less than the  $\gamma$  determined for the noun model. Moreover the link type factor  $\beta$  in the verb model also more quickly reduces the similarity of a node with distance in the hierarchy. So too does the path type factor discount relationships multiple link paths more severely. All of these facts confirm that the verb hierarchy is very shallow (in WordNet if not in humans), and means that the verb hierarchy is of limited help in assessing the similarity of verbs.

Thus the Yang and Powers noun similarity model does not adapt so directly or so well to verbs in the WordNet hierarchy. This is clearly connected to our previous observation that the verb taxonomy is shallower, but another significant factor is that the verb hierarchy does not include a second part-whole analog to the holo/meronym links of the noun hierarchy.

Such relationships do exist and correspond to the concept of metonymy, where there is a relationship between a word that describes a complex action or scene and one that describes a more specific aspect of that activity. For example, one of the poorly handled pairs in our data set is 'market' versus 'sell'. When we compare the noun sense of 'market' with 'sell' or 'sale' we do much better. Similarly if we

could recognize that marketing is a complex activity which involves price setting, product packaging, advertising, and selling, as metonymously related activities, we could again do better. The first improvement can be made by connecting the two hierarchies into one and using a single bidirectional search to evaluate similarity of any noun or verb against any other noun or verb – this is straightforward and is planned as part of our refinement of these techniques. The second improvement is not so straightforward, as it would seem to require manual augmentation of WordNet with the additional hierarchy, although of course there is always the possibility that WordNet-like hierarchies and variations could be self-organized based on corpus data and this we are also exploring.

The fallback into the use of glosses, stems, or noun similarity, does improve the situation but this increases the set of parameters to nine – three for the noun similarity, three for the basic verb similarity, and three for the three fallback options. However, this increase in the number of parameters does not seem to make the system brittle, as the tuning curves have fairly flat peaks and the tuning effects are relatively minor compared with the improvement due to the fallback mechanisms.

We note that the fallback model is a very primitive data fusion technique and thus also propose to investigate other fusion models.

#### Acknowledgements

Our thanks go to helpful comments and pointers to other relevant work from the anonymous referees.

#### References

- Collins, A. M. and M. R. Quillian (1969). "Retrieval time from semantic memory." *Journal of Verbal Learning and Verbal Behavior* **8**: 240–47.
- Fellbaum, C. (1998). *An Electronic Lexical Database*. London, England, The MIT Press.
- Gasperin, C., et al. (2001). Using Syntactic Contexts for Measuring Word Similarity. Workshop on Semantic Knowledge Acquisition & Categorisation (ESSLLI 2001). Helsinki.
- Grefenstette, G. (1993). Evaluation Techniques for Automatic Semantic Extraction: Comparing Syntactic and Window Based Approaches. Workshop on acquisition of lexical knowledge from text Columbus.
- Hirst, G. and D. St. Onge (1995). Lexical chains as representations of context for the detection and correction of malapropisms. WordNet. C. Fellbaum. Cambridge, MA, The Mit Press.
- Jarmasz, M. and S. Szpakowicz (2003). Roget's Thesaurus and Semantic similarity.
- Jiang, J. and D. Conrath (1997). Semantic similarity based on corpus statistics and lexical taxonomy. The International Conference on Research in Computational Linguistics, Taiwan.
- Landauer, T. K. and S. T. Dumais (1997). "A solution to Plato's problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge." *Psychological Review* **104**: 211-240.
- Leacock, C. and M. Chodorow (1998). Combining local context and WordNet similarity for word sense identification. WordNet: An electronic lexical database. C. Fellbaum, MIT Press: 265-283.
- Lesk, M. (1986). Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. The 5th Annual International Conference on Systems Documentation, ACM Press.
- Lin, D. (1997). Using syntactic dependency as a local context to resolve word sense ambiguity. The 35th Annual Meeting of the Association for Computational Linguistics, Madrid.
- Miller, G. (1995). "A lexical database for English." *Communications of the ACM* **38,11**: 39-41.
- Pedersen, T., et al. (2003). Maximizing Semantic Relatedness to Perform Word Sense Disambiguation.
- Quillian, M. R. (1967). "Word concepts: A theory and simulation of some basic semantic capabilities." *Behavioral Science* **12**: 410-30.
- Resnik, P. (1995). Using information content to evaluate semantic similarity in a taxonomy. Proceedings of IJCAI-95.
- Resnik, P. and M. Diab (2000). Measuring verb similarity. The 22nd Annual Meeting of the Cognitive Science Society (CogSci 2000).
- Rubenstein, H. and J. B. Goodenough (1965). "Contextual correlates of synonymy." *Communications of the ACM* **8(10)**: 627–633.
- Schütze, H. (1992). Dimensions of meaning. The 1992 ACM/IEEE Conference on Supercomputing, Minneapolis, Minnesota, United States.
- Tatsuki, D. (1998). Basic 2000 Words - Synonym Match 1. In: Interactive JavaScript Quizzes for ESL Students. <http://www.aitech.ac.jp/~iteslj/quizzes/js/dt/mc-2000-01syn.html>.
- Turney, P. D. (2001). Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. The Twelfth European Conference on Machine Learning (ECML2001), Freiburg, Germany.
- Wu, Z. and M. Palmer (1994). Verb semantics and lexical selection. The 32nd. Annual Meeting of the Association for Computational Linguistics.
- Yang, D. and D. M. W. Powers (2005). Measuring Semantic Similarity in the Taxonomy of WordNet. The Twenty-Eighth Australasian Computer Science Conference (ACSC 2005), Newcastle, Australia, ACS.



## 7 Appendix: the 130 pairs of verbs

brag	boast	hail	acclaim	refer	explain	request	levy	anger	approve
concoct	devise	dissipate	disperse	finance	build	arrange	study	approve	boast
divide	split	approve	support	expect	deserve	relieve	hinder	research	distribute
build	construct	impose	levy	terminate	postpone	move	swell	request	concoct
end	terminate	hasten	accelerate	yell	boast	weave	print	boast	yield
accentuate	highlight	rap	tap	swell	curl	swear	think	furnish	impress
demonstrate	show	lean	rest	rotate	situate	forget	resolve	refine	sustain
solve	figure out	make	earn	seize	request	supervise	concoct	acknowledge	distribute
consume	eat	show	publish	approve	scorn	situate	isolate	clean	concoct
position	situate	sell	market	supply	consume	explain	boast	lean	grate
swear	vow	weave	intertwine	clip	twist	ache	spin	postpone	show
furnish	supply	refer	direct	divide	figure out	evaluate	terminate	hail	judge
merit	deserve	distribute	commercialize	advise	furnish	recognize	succeed	remember	hail
submit	yield	twist	intertwine	complain	boast	dilute	market	scrape	lean
seize	take	drain	tap	want	deserve	hasten	permit	sweat	spin
spin	twirl	depict	recognize	twist	fasten	scorn	yield	highlight	restore
enlarge	swell	build	organize	swing	crash	swear	describe	seize	refer
swing	sway	hail	address	make	trade	arrange	explain	levy	believe
circulate	distribute	call	refer	hinder	yield	discard	arrange	alter	highlight
recognize	acknowledge	swing	bounce	build	propose	list	figure out	refer	carry
resolve	settle	yield	seize	express	figure out	stamp	weave	empty	situate
prolong	sustain	split	crush	resolve	examine	market	sweeten	flush	spin
tap	knock	challenge	yield	bruise	split	boil	tap	shake	swell
block	hinder	hinder	assist	swing	break	sustain	lower	imitate	highlight
arrange	plan	welcome	recognize	catch	consume	resolve	publicize	correlate	levy
twist	curl	need	deserve	swear	explain	dissipate	isolate	refer	lean