

Distribution-Based Semantic Similarity of Nouns*

Igor A. Bolshakov and Alexander Gelbukh

Center for Computing Research (CIC), National Polytechnic Institute (IPN),
Av. Juan Dios Bátiz s/n, Col. Nueva Industrial Vallejo,
07738, Mexico City, Mexico
igor@cic.ipn.mx, gelbukh@gelbukh.com,
www.Gelbukh.com

Abstract. In our previous work we have proposed two methods for evaluating semantic similarity / dissimilarity of nouns based on their modifier sets registered in Oxford Collocation Dictionary for Student of English. In this paper we provide further details on the experimental support and discussion of these methods. Given two nouns, in the first method the similarity is measured by the relative size of the intersection of the sets of modifiers applicable to both of them. In the second method, the dissimilarity is measured by the difference between the mean values of cohesion between a noun and the two sets of modifiers: its own ones and those of the other noun in question. Here, the cohesion between words is measured via Web statistics for co-occurrences of words. The two proposed measures prove to be in approximately inverse dependency. Our experiments show that Web-based weighting (the second method) gives better results.

Keywords: Semantic relatedness, word space model, lexical resources, Web as corpus, natural language processing.

1 Introduction

Several works evaluate semantic similarity or dissimilarity between words, see [3, 11] and references therein. The majority of evaluations are based on semantic hierarchies of WordNet [4, 5]. In this class of methods, semantic dissimilarity between words is considered proportional to the number of steps separating corresponding nodes of the hierarchy. The nodes are synsets that include the words under evaluation, while the arcs are subset-to-superset links between the synsets. The greater the distance, the greater dissimilarity. This measure proved to be useful in many applications and tasks of computational linguistics, such as word sense disambiguation [9], information retrieval, etc.

Another possible way for estimation of semantic proximity of words consists in comparing the sets of other words frequently co-occurring in texts in close vicinity to the two words in question [6]. The more similar the recorded beforehand sets of standard neighbors of any two words of the same POS, the more semantically similar the

* Work done under partial support of Mexican Government (CONACyT, SNI, SIP-IPN, COTEPABE-IPN). Authors thank anonymous reviewers for valuable comments.

words [7]. As applied to nouns, the accompanying words are primordially their modifiers. In European languages, these are usually adjectives or participles; in English these are also nouns staying in preposition and used attributively.

We evaluate semantic similarity or dissimilarity of English nouns by two methods in this paper. Both of them are based on those standard modifier sets for few tens of commonly used English nouns that are registered for them in OCDSE that seems the most reliable source of English collocations so far [10]. The nouns were preferred with more numerous collections of modifiers recorded.

In the first method, the similarity $Sim(N_1, N_2)$ of the noun N_1 to the noun N_2 is measured by the ratio of the number of modifiers commonly applicable to the both nouns and the number of modifiers of N_2 .

In the second method, we weight the relatedness between the noun and its modifiers by the Web co-occurrence statistics. Namely, the dissimilarity $DSim(N_1, N_2)$ of N_1 from N_2 is measured by the residual of two mean values of specially introduced *Stable Connection Index*. *SCI* is exteriorly like Mutual Information of two words [8] and operates by raw statistics of Web pages containing these words considered separately and in their close co-occurrences. In contrast to Mutual Information, it does not require repetitive evaluation of the total amount of pages under search engine's control. One mean value covers *SCIs* of all 'noun \rightarrow its own modifier' pairs, another mean value covers *SCIs* of all ' $N_1 \rightarrow$ a modifier of N_2 ' pairs. English modifiers usually precede their nouns forming bigrams with them, thus facilitating reliable Web statistic evaluations. In other words, *Sim* is determined through coinciding modifiers of nouns, while *DSim* is determined through alien modifiers.

The main idea of the two methods discussed here was briefly presented in our previous work [2]. In this paper, we give more details on the experiments conducted to compare these two methods.

Namely, our experimental data show that though the *Sim* and *DSim* measures can be rather arbitrary in each specific case, on average they show an inverse monotonic interdependence. However, in our experiments *DSim* showed higher resolution. By higher resolution we mean that while many noun pairs have zero *Sim* values as measured according to the OCDSE, they differ significantly in their *DSim* values.

2 Modifier Sets Selected for Evaluations

English nouns with all their recorded modifiers—both adjectives and nouns in attributive use—were taken from OCDSE. The nouns were picked up in rather arbitrary manner, without taking into account their mental similarity. Our only preferences were with the nouns with larger modifier sets.

For 32 nouns taken, total amount of modifiers (partially repeating) is 1964, and the mean modifiers group size equals to 61.4, varying from 39 (for *comment* and *disease*) to 119 (for *eyes*). The second and the third ranks determined by the set sizes are with *expression* (115) and *effect* (105). The nouns selected and sizes of their modifier sets are shown in Table 1.

We have limited the number of nouns to 32 units, since the total amount of accesses to the Web in experiments of the second method (cf. Section 5) grows approximately as a square of the number of words in question, so that, taking into

Table 1. Selected nouns and sizes of their modifier sets

S/N	Noun	MSet Size	S/N	Noun	MSet Size
1	<i>answer</i>	44	17	<i>effect</i>	105
2	<i>chance</i>	43	18	<i>enquiries</i>	45
3	<i>change</i>	71	19	<i>evidence</i>	66
4	<i>charge</i>	48	20	<i>example</i>	52
5	<i>comment</i>	39	21	<i>exercises</i>	80
6	<i>concept</i>	45	22	<i>expansion</i>	44
7	<i>conditions</i>	49	23	<i>experience</i>	53
8	<i>conversation</i>	52	24	<i>explanation</i>	59
9	<i>copy</i>	61	25	<i>expression</i>	115
10	<i>decision</i>	40	26	<i>eyes</i>	119
11	<i>demands</i>	98	27	<i>face</i>	96
12	<i>difference</i>	53	28	<i>facility</i>	89
13	<i>disease</i>	39	29	<i>fashion</i>	61
14	<i>distribution</i>	58	30	<i>feature</i>	51
15	<i>duty</i>	48	31	<i>flat</i>	48
16	<i>economy</i>	42	32	<i>flavor</i>	50

account limitations of Internet searchers and the general trend of all statistics to grow, we could afford several days to acquire all necessary statistics but not a month.

Some nouns (*conditions*, *demands*, *enquiries*, *exercises*, and *eyes*) were taken in plural, since they are used with the recorded modifier sets in plural more frequently than in singular.

3 Influence of Intersection of Modifier Sets

In our first method, the similarity $Sim(N_i, N_j)$ is defined through the intersection ratio of modifier sets $M(N_i)$ and $M(N_j)$ of the two nouns by the formula

$$Sim(N_i, N_j) \equiv \frac{|M(N_i) \cap M(N_j)|}{|M(N_i)|}, \quad (1)$$

where $|M(N_i)|$ means cardinal number of the set $M(N_i)$, \cap designates set intersection.

With such definition, the similarity measure is generally asymmetric: $Sim(N_i, N_j) \neq Sim(N_j, N_i)$, though both values are proportional to the number of commonly applicable modifiers. We can explain the asymmetry by means of the following extreme case. If $M(N_i) \subset M(N_j)$, each member of $M(N_i)$ has its own counterpart in $M(N_j)$, thus $Sim(N_i, N_j)$ reaches the maximum equal to 1 (just as when $M(N_i) = M(N_j)$), but some members of $M(N_j)$ have no counterparts in $M(N_i)$, so that $Sim(N_j, N_i) < 1$.

To better visualize the similarity, we put to Table 2 symmetric ratios

$$Sym(N_i, N_j) \equiv \frac{|M(N_i) \cap M(N_j)|}{\sqrt{|M(N_i)| |M(N_j)|}},$$

Quite dissimilar pairs (with zero *Sym* value) are quite numerous (76): {*change, copy*}, {*charge, decision*}, {*comment, decision*}, {*answer, disease*}, {*chance, disease*}, etc. The nouns for human features *eyes* and *face* proved to be very productive in modifiers (119 and 96 relatively) but very specific (their *Sym* measures are close to zero for majority of noun pairs).

4 Words Cohesion in Internet

Any words W_1 and W_2 may be considered forming a stable combination if their co-occurrence number $N(W_1, W_2)$ in a text corpus divided by S (the total number of words in the corpus) is greater than the product of relative frequencies $N(W_1)/S$ and $N(W_2)/S$ of the words considered apart. Using logarithms, we have a measure of word cohesion known as log-likelihood ratio or Mutual Information [8]:

$$MI(W_1, W_2) \equiv \log \frac{S \cdot N(W_1, W_2)}{N(W_1) \cdot N(W_2)}.$$

MI has important feature of scalability: if the values of all its building blocks S , $N(W_1)$, $N(W_2)$, and $N(W_1, W_2)$ are multiplied by the same factor, *MI* preserves its value.

Any Web search engine automatically delivers statistics on a queried word or a word combination measured in numbers of relevant Web pages, and no direct information on word occurrences or co-occurrences is available. We can re-conceptualize *MI* with all $N()$ as numbers of relevant pages and S as the page total managed by the engine. However, now $N()/S$ are not the empirical probabilities of corresponding events: the words that occur at the same a page are indistinguishable in the raw statistics, being counted only once, and the same page is counted repeatedly for each word included. We only hope that the ratios $N()/S$ are monotonically connected with the corresponding empirical probabilities for the events under consideration.

In such a situation a different word cohesion measure was construed from the same building blocks [1]. It conserves the feature of scalability, gives very close to *MI* results for statistical description of rather large sets of word combinations, but at the same time is simpler to be reached, since does not require repeated evaluation of the whole number of pages under the searcher's control. The new cohesion measure was named Stable Connection Index:

$$SCI(W_1, W_2) \equiv 16 + \log_2 \frac{N(W_1, W_2)}{\sqrt{N(W_1) \cdot N(W_2)}}. \quad (2)$$

The additive constant 16 and the logarithmic base 2 were chosen rather arbitrary. The constant 16 does not affect the comparisons discussed in this paper and is included purely for sake of tradition (since this is how the notion of *SCI* has been introduced previously); the reader can safely ignore it.

Since our experiments with Internet searchers need minimally several days to perform, some additional words on Web searchers are worthwhile here.

The statistics of searcher have two sources of variation in time. The first one is monotonic growing because of steady enlargement of searcher's DB. In our experience, for huge searcher's BDs and the queried words forming stable combinations, the

raw statistics $N(W_1)$, $N(W_2)$, $N(W_1, W_2)$ grow approximately with the same speed, so that *SCI* keeps its value with the precision to the second decimal digit, even if the statistics are got in different moments along the experimental day.

The second, fluctuating source of instability of Internet statistics is selection by the searcher of a specific processor and a specific path through searcher's DB—for each specific query. With respect to this, the searchers are rather different. For example, Google, after giving several very close statistics for a repeating query, can play a trick, suddenly giving twice fewer amount (with the same set of initial snippets!), thus shifting *SCI* significantly. Since we did not suffer of such troubles so far on behalf of AltaVista, we preferred it for our experiments.

5 Dissimilarity Based on Mean Cohesion Values

Consider first the mean cohesion values

$$\frac{1}{|M(N_i)|} \sum_{x \in M(N_i)} SCI(N_i, x)$$

between the noun N_i and all modifiers in its own modifier set $M(N_i)$. One can see in Table 3 that all mean *SCI* values are positive and mainly rather big (4 to 8), except for *enquiries*. As to the latter, we may suppose that occurrence statistics of British National Corpus—the base for selection of collocations in OCDSE—differ radically from Internet statistics, probably because OCDSE is oriented to the British variant of the English language, while Internet is mostly composed of texts written in American English or in international sort of English. Hence the collocations *intellectual / joint / open / critical / sociological... enquiries*, being rather rare in whole Internet, were inserted to OCDSE by purely British reasons. This is not unique case of British vs. USA language discrepancies. We had rejected orthographic differences like *flavour* vs. *flavor*, but we did not feel free to sift out such OCDSE collocations as *coastal flat* 'property by the sea,' which proved to be rare in Internet as a whole.

When calculating the *SCI* value of 'noun \rightarrow modifier of a different noun' pairs that mainly are not normal collocations, we frequently observe the cases with zero co-occurrence number in Internet. Then formula (2) gives *SCI* value equal to $-\infty$. To avoid the singularity, we take the value -16 for such cases, i.e. the maximally possible positive value, but with the opposite sign.

We define the dissimilarity measure as

$$DSim(N_i, N_j) = \frac{1}{|M(N_i)|} \sum_{x \in M(N_i)} (SCI(N_i, x) - SCI(N_j, x)) \quad (3)$$

i.e., as the mean difference between the *SCI* value of the modifiers of N_i with N_i and N_j , respectively. Note that in this formula the noun in question is compared with the set of its own modifiers defined by the dictionary and with the set of the modifiers of the other noun. Two things can be observed as to this definition.

Table 3. The mean SCI values of nouns with their own modifiers

S/N	Noun	Mean SCI	S/N	Noun	Mean SCI
1	<i>answer</i>	6.3	17	<i>effect</i>	6.7
2	<i>chance</i>	4.9	18	<i>enquiries</i>	1.4
3	<i>change</i>	6.5	19	<i>evidence</i>	8.0
4	<i>charge</i>	5.6	20	<i>example</i>	6.1
5	<i>comment</i>	4.4	21	<i>exercises</i>	4.0
6	<i>concept</i>	5.9	22	<i>expansion</i>	6.4
7	<i>conditions</i>	6.5	23	<i>experience</i>	7.7
8	<i>conversation</i>	6.0	24	<i>explanation</i>	6.1
9	<i>copy</i>	5.4	25	<i>expression</i>	4.9
10	<i>decision</i>	7.2	26	<i>eyes</i>	6.0
11	<i>demands</i>	4.1	27	<i>face</i>	5.7
12	<i>difference</i>	6.2	28	<i>facility</i>	4.5
13	<i>disease</i>	8.3	29	<i>fashion</i>	5.1
14	<i>distribution</i>	6.7	30	<i>feature</i>	5.9
15	<i>duty</i>	5.6	31	<i>flat</i>	4.3
16	<i>economy</i>	6.7	32	<i>flavor</i>	6.1

First, the formula is not symmetric. As it was discussed above, we consider the relations between different nouns more as inclusion than as distance: *cat* is a perfect *animal*, i.e., in our terminology we would say that *cat* is no different from *animal*, while *animal* by no means is a perfect *cat*.

Another observation about this definition is more theoretical. It seems to be contradicting: while we use the objective reality, the Web (as corpus) to measure the relatedness between a noun and a modifier, we seemingly arbitrary restrict the set of participating modifiers to be considered by those found in a dictionary, which were subjectively selected by a lexicographer. What is more, this seemingly leads to the necessity to use in our method a specialized large lexical resource, which does not exist in all languages, and it is not clear how the results obtained with different such resources would coincide.

Though we did not conduct any corresponding experiments, we believe that the formula above can be modified to use the whole set of words of the language (occurring in a large corpus or in the Web). The formula is then to be modified to take into account the cohesion between each word and the noun in question; those words that have low value of such cohesion would be weighted out. However, this would be a bit impractical. So we here use an approximation to such a totally unsupervised approach. Our approximation takes advantage of an already existing resource to roughly indicate which words are expected to correlate with the given noun.

Note that in this sense the second method can be thought of as a weighted variant of the first one.

Table 4 shows the pairs with the smallest and the greatest dissimilarity measure in our small dataset. One can notice the pairs with the smallest dissimilarity, such as {*enquiries*, *explanation*}, do have similar or related meaning, while those with greater dissimilarity, such as {*disease*, *enquiries*}, look totally unrelated.

Table 4. Most and least similar noun pairs in our sample

Least dissimilar noun pairs				Most dissimilar noun pairs			
Noun ₁	Noun ₂	DSim	Sim	Noun ₁	Noun ₂	DSim	Sim
<i>enquiries</i>	<i>explanation</i>	0.3	0.156	<i>disease</i>	<i>enquiries</i>	18.5	0.000
<i>enquiries</i>	<i>distribution</i>	0.5	0.022	<i>eyes</i>	<i>enquiries</i>	15.8	0.017
<i>enquiries</i>	<i>comment</i>	0.6	0.111	<i>effect</i>	<i>enquiries</i>	14.8	0.029
<i>enquiries</i>	<i>conversation</i>	0.6	0.089	<i>face</i>	<i>enquiries</i>	14.7	0.010
<i>enquiries</i>	<i>change</i>	0.9	0.044	<i>experience</i>	<i>enquiries</i>	14.4	0.000
<i>difference</i>	<i>change</i>	1.1	0.321	<i>disease</i>	<i>economy</i>	14.2	0.000
<i>enquiries</i>	<i>fashion</i>	1.1	0.022	<i>disease</i>	<i>chance</i>	14.0	0.000
<i>enquiries</i>	<i>charge</i>	1.2	0.067	<i>flavor</i>	<i>enquiries</i>	14.0	0.020

In fact, the very small *DSim* measure can indicate that the words are nearly synonyms or nearly antonyms, but this results from a different our research.

6 Comparison and Discussion

Comparing the *Sim* and *DSim* values for the 16 pairs in Table 4, one can see that the pairs with maximal *Sim* values usually have minimal *DSim* values and vice versa, i.e. an inverse monotonic dependency exists between the two measures. More representative comparison is given in Figure 1 that gives correlations between *Sim* and *DSim* on the plane.

A statistically proved inverse monotonic dependency is quite clear from Figure 1. One can also comprehend that *DSim* has higher resolution for semantically most different nouns. Indeed, the numerous pairs with zero *Sim* values have quite diverse *DSim* values, from 14.0 for {*disease*, *flat*} to 4.2 for {*flat*, *answer*}. Hence the use of *DSim* measure seems preferable.

7 Conclusions and Future Work

Two methods of numerical evaluation of semantic similarity of any nouns is proposed. The evaluations are based on comparison of standard modifiers of the nouns registered in OCDSE. The first method evaluates similarity by the portion of common modifiers of the nouns, while the second one evaluates dissimilarity by the change of the mean cohesion of a given modifier set with its own noun and an alien one.

Cohesion measurements are based on raw Web statistics of occurrences and co-occurrences of supposedly cohesive words. It is shown that dissimilarity measured through the Web has higher resolution and thus may have greater reliability.

Both methods do not depend on language and can be easily tested on the resources of other languages. Currently we are conducting experiments with Spanish and Russian, which are morphologically-rich languages. For English, it is worthwhile to repeat evaluations for a greater number of nouns and for different source of modifiers sets, e.g. for a large corpus of American origin. Finally, we believe that this method can be applied to words of parts of speech other than nouns, though one should be much more careful with, say, verbs, where the co-occurrence patterns are much more lexicalized and less semantic than those of nouns.

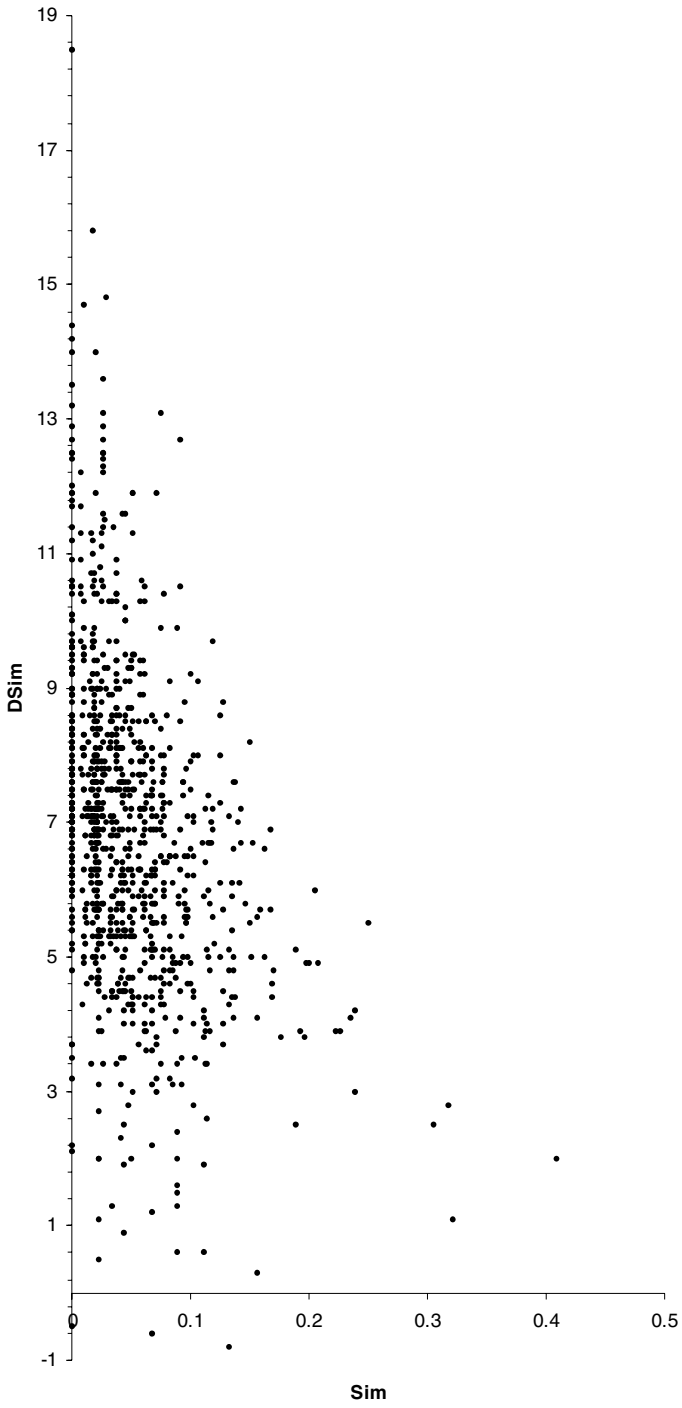


Fig. 1. Correlations between *Sim* and *DSim*

References

1. Bolshakov, I.A., Bolshakova, E.I.: Measurements of Lexico-Syntactic Cohesion by means of Internet. In: Gelbukh, A., de Albornoz, Á., Terashima-Marín, H. (eds.) MICAI 2005. LNCS (LNAI), vol. 3789, pp. 790–799. Springer, Heidelberg (2005)
2. Bolshakov, I.A., Gelbukh, A.: Two Methods of Evaluation of Semantic Similarity of Nouns Based on Their Modifier Sets. In: LNCS, vol. 4592, Springer, Heidelberg (2007)
3. Cilibrasi, R.L., Vitányi, P.M.B.: The Google Similarity Distance. *IEEE Transactions on Knowledge and Data Engineering* 19(3), 370–383 (2007), www.cwi.nl/paulv/papers/tkde06.pdf
4. Fellbaum, C. (ed.): *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge (1998)
5. Hirst, G., Budanitsky, A.: Correcting Real-Word Spelling Errors by Restoring Lexical Cohesion. *Natural Language Engineering* 11(1), 87–111 (2005)
6. Keller, F., Lapata, M.: Using the Web to Obtain Frequencies for Unseen Bigram. *Computational linguistics* 29(3), 459–484 (2003)
7. Lin, D.: Automatic retrieval and clustering of similar words. In: *COLING-ACL 1998, Canada* (1998)
8. Manning, C.D., Schütze, H.: *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge (1999)
9. McCarthy, D., Koeling, R., Weeds, J., Carroll, J.: Finding Predominant Word Senses in Untagged Text. In: *Proc. 42nd Annual Meeting of the ACL, Barcelona, Spain* (2004)
10. *Oxford Collocations Dictionary for Students of English*. Oxford University Press (2003)
11. Patwardhan, S., Banerjee, S., Pedersen, T.: Using Measures of Semantic Relatedness for Word Sense Disambiguation. In: Gelbukh, A. (ed.) *CICLing 2003*. LNCS, vol. 2588, Springer, Heidelberg (2003)