

EZ.WordNet: principles for automatic generation of a coarse grained WordNet

Rada Mihalcea and Dan I. Moldovan

Department of Computer Science and Engineering
Southern Methodist University
Dallas, Texas, 75275-0122
{rada, moldovan}@enr.smu.edu

Abstract

In this paper, we propose several principles that enable the automatic transformation of WordNet into a coarser grained dictionary, without affecting its existing semantic relations. We derive a new version of WordNet leading to a reduction of 26% in the average polysemy of words, while introducing a small error rate of 2.1%, as measured on a sense tagged corpus.

Introduction

In the Natural Language Processing (NLP) community, WordNet is well known as an invaluable resource: more and more applications that require machine readable dictionaries or world knowledge encoded in semantic networks use WordNet. There are certain applications, such as text inference or knowledge processing (Harabagiu & Moldovan 1998), which require the availability of a large set of relations among concepts and for which the large number of concepts and semantic links found in WordNet is well suited. On the other hand, other applications such as semantic or conceptual indexing, and sometimes word sense disambiguation or machine translation, do not always need such a fine grain distinction between concepts. This actually constitutes one of the most often “critiques” brought to WordNet: the fact that sometime word senses are so close together that a distinction is hard to be made even for humans.

We are proposing here a methodology that enables the automatic generation of a coarser WordNet by either collapsing together synsets (sets of synonym words) very similar in meaning, or dropping synsets very rarely used. The results obtained are encouraging: we show that using the rules presented in this paper one can automatically generate a new version of WordNet, which we call *EZ.WordNet.1*, with an average polysemy reduced with 26%, while the level of ambiguity introduced is only 2.1% as measured on SemCor (a corpus tagged with WordNet senses). We also derive an alternative version *EZ.WordNet.2*, using the same rules but with different parameters, that brings a reduction of 39% in average polysemy, with an error rate of 5.6%.

A method for defining the relatedness among synsets would be useful for several applications:

1. As pointed out in (Resnik & Yarowsky 1999), when evaluating Word Sense Disambiguation (WSD) systems it is important to know which senses of a word are similar and which are not.
2. WSD for certain applications such as semantic indexing (Gonzalo *et al.* 1998) do not need to make such a fine distinction among senses. More than this, having similar meanings collapsed together would bring more possible candidate words for the task of query expansion (Moldovan & Mihalcea 2000).
3. By reducing the semantic space, the task of WSD algorithms is made easier, as the correct sense has to be chosen among a smaller set of possible senses.
4. Machine translation could benefit from such a measure of relatedness among word meanings, as closed senses tend to have the same translation.

The rules presented in this paper are derived from sound principles for ambiguity testing, enabling a significant reduction in the polysemy of the words in WordNet, without introducing too much ambiguity. Simplified versions of WordNet can be automatically generated from its current version based on the rules described in this paper and can be used in various applications.

WordNet and synsets similarity

WordNet was developed at Princeton University by a group led by George A. Miller (Fellbaum 1998). It covers the vast majority of nouns, verbs, adjectives and adverbs from the English language. Table 1 shows the number of words, respectively the number of synsets defined in WordNet 1.6 for each part of speech.

WordNet defines one or more senses for each word. Depending on the number of senses it has, a word can be (1) *monosemous*, if it has only one sense or (2) *polysemous*, if it has two or more senses. Table 1 also shows the number of polysemous and monosemous words and the average polysemy measured on WordNet. It results an average polysemy of 2.91 for polysemous words and 1.34 if monosemous words are considered as well.

Humans tend to use more frequently words with higher polysemy, and this makes the distinction of word senses a problem. This fact is proven by statistics derived from SemCor (Miller *et al.* 1993) showing that words used in common texts have an average polysemy of 6.55. Table 2 lists the total number of occurrences

Part of speech	Word forms	Number synsets	Total senses	Polysemous words	Monosemous words	Total senses polys.words	Avg. polysemy (- monos)	Avg. polysemy (+ monos)
Noun	94473	66024	116318	12562	81911	34407	2.73	1.23
Verb	10318	12126	22067	4565	5753	16314	3.57	2.13
Adj	20169	17914	29882	5372	14797	15085	2.80	1.48
Adverb	4545	3574	5678	748	3797	1881	2.51	1.25
TOTAL	129505	99638	173945	23247	106258	67687	2.91	1.34

Table 1: Statistics on WordNet: number of word forms, number of synsets, number of monosemous and polysemous words, average polysemy in WordNet 1.6

of word forms in SemCor, and the number of senses defined in WordNet for these words forms.

Part of speech	Total word occ.	Total senses defined in WordNet	Average polysemy
Noun	88398	382765	4.33
Verb	90080	944368	10.48
Adj	35770	157751	4.41
Adv	20595	55617	2.70
TOTAL	234843	1540501	6.55

Table 2: Statistics on SemCor: total number of word occurrences, total number of senses defined in WordNet and average polysemy

It results that in common texts the average polysemy of words with respect to WordNet is very high, and this is why WordNet is classified as a fine grained dictionary. **Similarity measures in WordNet.** A proof that similarity among synsets is useful to be known is constituted by the fact that there are already some similarity relations defined in WordNet. The problem is that these relations are defined only for nouns and verbs, and do not always succeed in indicating properly two synsets as similar in meaning.

For verbs, a *VERBGROUP* pointer is defined, pointing to synsets similar in meanings. For nouns, there are three types of synset relations defined in WordNet: (1) *sisters*, representing synsets with a word in common and with the same direct hypernym; (2) *cousins*, which are synsets with a word in common and with at least one pair of direct or indirect hypernyms related, based on a predefined list (there is also a list of exception synsets defined for this relation) and (3) *twins*, which are synsets with three or more words in common.

These relations are intended to provide the user of WordNet with a measure of similarity among the different senses of a word. There are no such relations defined for adjectives and adverbs.

Several problems are associated with these relations: (1) coverage of verb groups is incomplete (as mentioned in the WordNet manuals); (2) the measures for noun synsets similarity are sometime too strong, and sometime too loose; for example, these relations will wrongly group together $\{house\#3\}$ with its meaning of “a building in which something is sheltered or located” and $\{house\#5, theater, theatre\}$ meaning “a building where theatrical performances or motion-picture shows can be presented”, but they will fail to group as similar the synsets $\{decapitation\#1, beheading\}$ defined as “execution by cutting off the victim’s head” and $\{decapitation\#2, beheading\}$ meaning “killing by cutting off the head”.

The methods we are proposing are able to group similar synsets together based on set of rules derived from sound principles. These methods can also be applied to adjectives and adverbs.

Semantic principles: tests for ambiguity

When talking about word meanings and ambiguity, one of the problems is how to actually determine if a word form is ambiguous or not. (Cruse 1986) describes three principles used by humans in testing the ambiguity of words. We present them here and show how these principles can be translated into methods of measuring the ambiguity level among the synsets in WordNet.

PRINCIPLE I. *If there exists a synonym or one occurrence of a word form which is not a synonym of a second, syntactically identical occurrence of the same word in a different context, then that word form is ambiguous, and the two occurrences exemplify different senses.*

The example given as an illustration of this principle is the word “match”, which in the contexts “Gay struck the match”, respectively “The match was draw” have two different synonyms, namely “lucifer” (or “friction match”) and “contest”. It results that this word is ambiguous in the two given contexts.

This principle can be reformulated in the following way: if a word form, in two different contexts, has the same synonyms, then that word form has similar meanings in the given contexts.

Translated in WordNet terminology, the context of a word meaning is represented by its synset and its relations to other synsets. We can infer that if a word has two senses with the same synonyms, then it is hard to make the distinction among the two senses, and thus we can collapse the appropriate synsets together. The following rule is inferred:

Rule SP1.1 If S1 and S2 are two synsets containing at least two words, and if S1 and S2 contain the same words, then S1 and S2 can be collapsed together into one single synset S12.

Example: Sense #1 and #3 for *paper* are:
S1 = {newspaper, paper} (a daily or weekly publication on folded sheets)
S3 = {newspaper, paper} (a newspaper as a physical object)
These two synsets can be collapsed together, as they are very similar in meaning.

There are cases when the synsets contain only one single word, and thus we cannot directly apply this principle of measuring ambiguity among synsets. A good approximation of the *synonymy* relation is represented by the *hypernymy* links found in WordNet: if a particular sense of a word has no synonyms, than its meaning can be judged by looking at its hypernym. Thus, another rule that can be inferred from principle I is:

Rule SP1.2

If S1 and S2 are two synsets with the same hypernym, and if S1 and S2 contain the same words, than S1 and S2 can be collapsed together into one single synset S12.

Example: Consider senses #1 and #2 for the verb *eat*:

S1 = {*eat*} (*take in solid food*)
 \Rightarrow {*consume, ingest, take in, take, have*}
 S2 = {*eat*} (*eat a meal*)
 \Rightarrow {*consume, ingest, take in, take, have*}

These two senses are fine distinctions of the possible meanings of *eat*, and the two synsets S1 and S2 can be collapsed together.

By relaxing this first principle in its requirement of having *all* the synonyms of two different word meanings identical in order to fuse the appropriate synsets, to a requirement of having at least K identical synonyms, we can infer Rule SP1.3. By taking K=3, we obtain the *twins* similarity measure from WordNet.

Rule SP1.3 If S1 and S2 are two synsets with at least K words in common, then S1 and S2 can be collapsed together into one single synset S12.

Example Sense #1 and #3 for noun *teaching* are:
 S1 = {*teaching, instruction, pedagogy*} (*the profession of a teacher*)
 S3 = {*education, instruction, teaching, pedagogy, educational activity*} (*activities that impart knowledge*)

PRINCIPLE II. *If there exists a word or expression standing in a relation of oppositeness to one occurrence of a word form, which does not stand in the same relation to a second, syntactically identical occurrence of the same word form in a different context, then that word form is ambiguous, and the two occurrences exemplify different senses.*

The example given for this ambiguity test is the word “*light*”, which is determined to be ambiguous in the sentences “*The room was painted in light colours*”, respectively “*Arthur has a rather light teaching load*”, as it has two different antonyms: “*dark*” and “*heavy*”.

The reformulation of this principle is: if a word form, in two different contexts, has the same antonyms, then that word form has similar meanings in the given contexts. We can translate this in WordNet terms, and measure the grade of similarity among two synsets with the following rule:

Rule SP2 If S1 and S2 are two synsets representing two senses of a given word, and if S1 and S2 have the same antonym, then S1 and S2 can be collapsed together into one single synset S12.

Example: Senses #1 and #2 for the verb *dress* are:
 S1 = {*dress, get dressed*} (*put on clothes*)
 antonym \Rightarrow {*undress, discase, uncase, unclothe, strip, strip down, disrobe*}
 S2 = {*dress, clothe, enclothe, garb, raiment, tog, garment, habitate, fit out, apparel*} (*provide with clothes or put clothes on*)
 antonym \Rightarrow {*undress, discase, uncase, unclothe, strip, strip down, disrobe*}

Finally, the last test of ambiguity refers to paronymic relations¹ related to a particular word form.

PRINCIPLE III. *If there exists a word which stands in a paronymic relation to one occurrence of a word form, but does not stand in the same relation to a second, syntactically identical occurrence of the same word form in a different context, then that word form is ambiguous, and the two occurrences exemplify different senses.*

This principle is illustrated using the noun “*race*”. In the sentence “*The race was won by Arkle*”, this noun has the verb “*to race*” related to it, while the same noun in the sentence “*They are a war-like race*” has two other related words, namely “*racial*” and “*racist*”. By having different paronymically related words, the noun “*race*” is determined to have different senses in these two examples.

Reformulated, this principle becomes: if a word form, in two different contexts, is paronymically related to the same words, than the word form has similar meanings in the given contexts.

The translation into WordNet terminology cannot be done for all parts of speech, as there are no relations defined in WordNet among verbs and nouns, that would help us determine paronymic relations of the form *act - actor*. Still, this principle can be used for adjectives and adverbs, where a *pertainymy* relation is defined. The following rule is derived:

Rule SP3 If S1 and S2 are two synsets representing two senses of a given word, and if S1 and S2 have the same pertainym, then S1 and S2 can be collapsed together into one single synset S12.

Example: Senses #1 and #5 for the adverb *lightly* are:
 S1 = {*lightly*} (*without good reason*)
 pertainym \Rightarrow {*light#5*} (*psychologically light*)
 S5 = {*lightly*} (*with indifference or without dejection*)
 pertainym \Rightarrow {*light#5*} (*psychologically light*)

Probabilistic principles

Besides the principles presented in the previous section, the polysemy of WordNet can be reduced based on the frequency of senses and the probability of having particular synsets used in a text. By dropping synsets with very low probability of occurrence, we can reduce the number of senses a word might have.

We need (1) a distribution of sense frequencies for the different parts of speech and (2) a method of deriving

¹Relations involving identity of the root, but different syntactic categories (e.g. *act - actor*)

Part of speech	Probability of having sense number								
	1	2	3	4	5	6	7	8	>8
Noun	78.52%	12.73%	4.40%	2.07%	0.98%	0.52%	0.34%	0.17%	<0.1%
Verb	61.01%	19.22%	7.89%	4.12%	2.64%	1.47%	0.98%	0.65%	<0.5%
Adj	80.98%	12.35%	3.96%	1.41%	0.51%	0.25%	0.16%	0.16%	<0.05%
Adv	83.84%	11.24%	3.67%	0.61%	0.42%	0.15%	0.03%	0.009%	<0.009%

Table 3: Statistics on SemCor: distribution of senses for nouns, verbs, adjectives and adverbs

Part of speech	Rule applied														Total reduced word senses
	SP0		SP1.1		SP1.2		SP1.3		SP2		SP3		PP1		
	(s)	(w)	(s)	(w)	(s)	(w)	(s)	(w)	(s)	(w)	(s)	(w)	(s)	(w)	
<i>EZ.WordNet.1 (K=3 MaxP=2)</i>															
Noun	—	—	349	743	216	216	100	328	2	3	—	—	1074	1142	2432
Verb	244	252	117	242	131	131	71	226	7	7	—	—	889	969	1827
Adj	—	—	115	244	—	—	29	90	8	8	12	12	931	1018	1372
Adv	—	—	23	52	—	—	1	3	6	6	96	100	84	85	246
<i>EZ.WordNet.2 (K=2 MaxP=5)</i>															
Noun	—	—	349	743	216	216	973	2015	2	3	—	—	3159	3344	6321
Verb	244	252	117	242	131	131	677	1257	7	7	—	—	1615	1767	3656
Adj	—	—	115	244	—	—	356	714	8	8	12	12	1628	1795	2773
Adv	—	—	23	52	—	—	47	92	6	6	96	100	158	165	418

Table 4: Reductions in the number of synsets ((s) columns) and in the number of word senses ((w) columns) obtained for each rule.

the probability of a synset occurring in a text, starting with the probabilities of its component words.

To determine the distribution of word sense frequencies, we used again SemCor, as this is the only corpus available in which words are sense tagged using WordNet. Table 3 show these sense distributions for nouns, verbs, adjectives and adverbs.

Let us denote a synset with $S = \{W_{i_1}, W_{i_1} \dots, W_{i_n}\}$, meaning that the synset is composed by words having senses $i_1, i_2 \dots, i_n$. If we denote with P_{i_k} the probability of occurrence of a word having sense i_k , then the probability of occurrence P_S for the synset S is equal with the summation of the probabilities of occurrence for the component words, i.e. $P_S = \sum_{k=1}^{k=n} P_{i_k}$.

In order to reduce the granularity of WordNet without introducing too much ambiguity, we use this formula together with probabilities derived from SemCor, and drop those synsets with a probability of occurrence P_S smaller than a given threshold. The following rule is derived:

Rule PP1 If S is a synset $S = \{W_{i_1}, W_{i_1} \dots, W_{i_n}\}$ with the probability of occurrence $P_S = \sum_{k=1}^{k=n} P_{i_k} < Max_P$ than S can be considered as a very rarely occurring synset and it can be dropped.

Example: The noun synset $S = \{draft\#11, draught\#5, drawing\#6\}$ (the act of moving a load by drawing or pulling) has the probability of occurrence $P_S = P_{11} + P_5 + P_6 = 0.1\% + 0.98\% + 0.52\% = 1.6\%$. For Max_P set to 2.0, this synset can be dropped.

Note that this way of computing the probability of a synset does not make reference to the component words themselves, but to their senses, and thus we do not have to deal with the problem of data sparseness that would result from the limited size of the corpus.

Applying the principles on WordNet

We applied these semantic and probabilistic principles on WordNet and generated two new versions, called *EZ.WordNet.1* and *EZ.WordNet.2*.

The semantic principles resulted in collapsed synsets, while the probabilistic principles determined which synsets can be dropped. By applying these rules, we obtain a reduction in the number of synsets, and implicitly a reduction in the number of word senses.

There are two variables used by the reduction rules, namely the K minimum number of common synonyms among two synsets, as required by *Rule SP1.3*, and Max_P , which is the maximum probability threshold for the *Rule PP1*. Depending on the values chosen for these parameters, one can obtain sense inventories closer to the original WordNet, but with a smaller reduction in polysemy, or versions of WordNet with a higher reduction in polysemy but with more synsets modified respect to WordNet.

We chose two sets of values for these variables, and consequently we obtained two versions of WordNet:

- *EZ.WordNet.1*, for $K = 3$ and $Max_P = 2$.
- *EZ.WordNet.2*, for $K = 2$ and $Max_P = 5$.

Table 4 shows, for each of these versions, the reduction obtained in number of synsets, respectively in the number of senses, for each of the semantic and probabilistic rules. *Rule SP0* is applicable only for verbs and it corresponds to the *VERBROUP* pointers already defined in WordNet. Combining the information derived from this table with the statistics shown in Table 1, we can calculate the average polysemy for the two new versions of WordNet. Table 5 shows the total senses for the polysemous words, as computed in *EZ.WordNet.1* and *EZ.WordNet.2*, as well as the average polysemy computed for these sense inventories.

A much more important and interesting result would be to measure the reduction in the number of senses for

Part of speech	Polysemous words	Monosemous words	<i>EZ.WordNet.1</i>			<i>EZ.WordNet.2</i>		
			Total senses polys.words	Avg.polys. (- monos)	Avg.polys. (+ monos)	Total senses polys.words	Avg.polys. (- monos)	Avg.polys. (+ monos)
Noun	12562	81911	31975	2.54	1.20	28086	2.23	1.16
Verb	4565	5753	14487	3.17	1.96	12658	2.77	1.78
Adj	5372	14797	13713	2.55	1.41	12312	2.29	1.34
Adv	748	3797	1635	2.18	1.19	1463	1.95	1.15
TOTAL	23247	106258	61810	2.65	1.29	54429	2.34	1.24

Table 5: Statistics on EZ.WordNet.1 and EZ.WordNet.2: number of senses for polysemous words, average polysemy for polysemous words only and for all words

Part of speech	Total word occurrences	<i>WordNet</i>		<i>EZ.WordNet.1</i>				<i>EZ.WordNet.2</i>			
		Total senses	Avg. polys.	Total senses	Avg. polys.	Missing senses	Error rate	Total senses	Avg. polys.	Missing senses	Error rate
Noun	88398	382765	4.33	316796	3.58	1461	1.65%	256178	2.89	4668	5.28%
Verb	90080	944368	10.48	668712	7.42	2879	3.1%	542629	6.02	6310	7.0%
Adj	35770	157751	4.41	119044	3.32	545	1.52%	101907	2.84	1366	3.81%
Adv	20595	55617	2.70	45928	2.23	200	0.97%	39732	1.92	818	3.97%
TOTAL	234843	1540501	6.55	1150480	4.89	5085	2.16%	940446	4.0	13162	5.6%

Table 6: Average polysemy and error rate obtained on SemCor for *EZ.WordNet.1* and *EZ.WordNet.2*. We also replicate, for comparison purposes, the total number of senses and average polysemy in WordNet, as shown in Table 2

the words commonly used by people, i.e. to determine the reduction in the polysemy of the words in SemCor. We can also make use of this corpus to determine the ambiguity introduced by the new sense inventories with respect to the original WordNet sense tagging.

We compute two measures on SemCor: (1) *the average polysemy* determined as the total number of senses for the words in SemCor, with respect to *EZ.WordNet.1* or 2, divided by the total number of words, as it results from Table 2; and (2) *error rate*, defined as the total number of words from SemCor that are not defined anymore in the new WordNet versions, divided by the total number of words in SemCor. Table 6 shows these figures computed for *EZ.WordNet.1* and *EZ.WordNet.2*.

Interpretation of results. Looking at both Table 1 and 5, it can be seen that the rules proposed in this paper enable an average reduction in the number of senses for polysemous words of 9% for *EZ.WordNet.1*, respectively 20% for *EZ.WordNet.2*, with respect to the original WordNet.

There is a difference between the polysemy of the words in a dictionary and the polysemy of the words actually used by humans (the words in the *active* vocabulary). This difference is clearly shown by Table 2, where one can see that the average polysemy of the words in a common text like SemCor is much higher than the average polysemy of the words in a dictionary.

Hence, a more representative result is obtained by comparing the average polysemies obtained with different sense inventories on a corpus, such as SemCor. We can also determine the error rate introduced by the usage of the new sense inventories.

From Table 6, it results a reduction of 26% in polysemy, with an error rate of 2.16%. The second version has a larger error rate, namely 5.6%, but it also brings a larger reduction of 39% in polysemy. The error rates of 2.1% and 5.6% are acceptable as it is considered that the accuracy obtained by humans in sense tagging is not larger than 92-94%. Depending on the application, one

of these versions or newly compiled versions of WordNet can be used.

Conclusions

One of the main problems associated with WordNet is its fine granularity. In this paper, we present a methodology for reducing the average polysemy of the words defined in WordNet. We derive a set of semantic and probabilistic rules that are used to either collapse synsets very similar in meaning, or drop synsets that are very rarely used. In this way, we obtain a new version of WordNet leading to a reduction of 26% in the average polysemy of words, while introducing a small error rate of 2.1%, as measured on SemCor. An alternative version is also derived, that brings a polysemy reduction of 39% with an error rate of 5.6%. These results are encouraging, as a coarse grained WordNet is known to be beneficial for a large range of applications.

References

- Cruse, D. 1986. *Lexical Semantics*. Cambridge Univ. Press.
- Fellbaum, C. 1998. *WordNet, An Electronic Lexical Database*. The MIT Press.
- Gonzalo, J.; Verdejo, F.; Chugur, I.; and Cigarran, J. 1998. Indexing with WordNet synsets can improve text retrieval. In *Proceedings of COLING-ACL '98 Workshop on Usage of WordNet in Natural Language Processing Systems*.
- Harabagiu, S., and Moldovan, D. 1998. *Knowledge Processing on an Extended WordNet*. The MIT Press. 289-405.
- Miller, G.; Leacock, C.; Randee, T.; and Bunker, R. 1993. A semantic concordance. In *Proceedings of the 3rd DARPA Workshop on Human Language Technology*, 303-308.
- Moldovan, D., and Mihalcea, R. 2000. Using WordNet and lexical operators to improve Internet searches. *IEEE Internet Computing* 4(1):34-43.
- Resnik, P., and Yarowsky, D. 1999. Distinguishing systems and distinguishing senses: new evaluation methods for word sense disambiguation. *Natural Language Engineering* 5(2):113-134.