

Word Sense Disambiguation based on Semantic Density

Rada Mihalcea and Dan I. Moldovan
Department of Computer Science and Engineering
Southern Methodist University
Dallas, Texas, 75275-0122
{rada,moldovan}@seas.smu.edu

Abstract

This paper presents a Word Sense Disambiguation method based on the idea of semantic density between words. The disambiguation is done in the context of WordNet. The Internet is used as a raw corpora to provide statistical information for word associations. A metric is introduced and used to measure the semantic density and to rank all possible combinations of the senses of two words. This method provides a precision of 58% in indicating the correct sense for both words at the same time. The precision increases as we consider more choices: 70% for top two ranked and 73% for top three ranked.

1 Introduction

Word Sense Disambiguation (WSD) is an open problem in Natural Language Processing. Its solution impacts other tasks such as discourse, reference resolution, coherence, inference and others. WSD methods can be broadly classified into three types:

1. WSD that make use of the information provided by machine readable dictionaries (Cowie et al.1992), (Miller et al.1994), (Agirre and Rigau, 1995), (Li et al.1995), (McRoy, 1992);
2. WSD that use information gathered from training on a corpus that has already been semantically disambiguated (supervised training methods) (Gale, Church et al., 1992), (Ng and Lee, 1996);
3. WSD that use information gathered from raw corpora (unsupervised training methods) (Yarowsky 1995) (Resnik 1997).

There are also hybrid methods that combine several sources of knowledge such as lexicon information, heuristics, collocations and others (McRoy, 1992) (Bruce and Wiebe, 1994) (Ng and Lee, 1996) (Rigau, Asterias et al., 1997).

Statistical methods produce high accuracy results for small number of preselected words. A lack of widely available semantically tagged corpora almost excludes supervised learning methods. On the other hand, the disambiguation using unsupervised methods has the disadvantage that the senses are not well defined. To our knowledge, none of the statistical methods disambiguate adjectives or adverbs so far.

One approach to WSD is to determine the *conceptual distance* between words, that is to measure the semantic closeness of the words within a semantic network. Essentially, it is the length of the shortest path connecting the concepts (Rada et al.1989), (Rigau, Asterias et al., 1997). By measuring the

conceptual distance between words, it is possible to determine the likelihood of word sense associations. For example, the method proposed in (Li et al.1995) tries to determine the possible sense of a noun associated with a verb using WordNet and a large text. Based on other occurrences of the verb or semantically related verbs in the text, the possible object is determined by measuring the semantic similarity between the noun objects.

Methods that do not need large corpora are usually based exclusively on MRD. A proposal in this sense has been made in (Agirre and Rigau, 1995); they measure the conceptual density between nouns, by using WordNet, but the method proposed in their paper cannot be applied to measuring a conceptual distance between a verb and a noun, as no direct links are provided in MRDs between the nouns and verbs hierarchies. A WordNet-based method for measuring the semantic similarity between nouns was also proposed in (Richardson et al., 1994). Their method consists of using hierarchical concept graphs constructed from WordNet data files, and a semantic similarity formula. Still, the method does not provide a link between different part-of-speech words.

2 Our approach

The approach described in this paper is based on the idea of *semantic density*. This can be measured by the number of common words that are within a semantic distance of two or more words. The closer the semantic relationship between two words the higher the semantic density between them. The way it is defined here, the semantic density works well in the case of uniform MRD. In reality there are gaps in the knowledge representations and the semantic density can provide only an estimation of the actual semantic relatedness between words.

We introduce the semantic density because it is

relatively easy to measure it on a MRD like WordNet. This is done by counting the number of concepts two words have in common. A metric is introduced in this sense which when applied to all possible combinations of the senses of two or more words it ranks them.

Another idea of this paper is to use the Internet as a raw corpora. Thus we have two sources of information: (1) the Internet for gathering statistics and (2) WordNet for measuring semantic density. As will be shown below, a ranking of words senses results from each of these two sources. The issue now is how to combine these two rankings in order to provide an overall ranking. One possibility is to use them in parallel and the other one is to use them serially. We have tried both and the serial approach provided better results. Thus, for a verb - noun pair, the WSD method consists of two Algorithms, the first one ranks the noun senses, of which we retain only the best two senses; and a second Algorithm takes the output produced by the first Algorithm and ranks the pairs of verb - noun senses. Extensions of this method to other pairs than verb - noun are discussed, and larger windows of more than two words are considered.

An essential aspect of the WSD method presented here is that we provide a ranking of possible associations between words instead of a binary yes/no decision for each possible sense combination. This allows for a controllable precision as other modules may be able to distinguish later the correct sense association from such a small pool.

WordNet is a fine grain MRD and this makes it more difficult to pinpoint the correct sense combination since there are many to choose from and many are semantically close. For applications such as machine translation, fine grain disambiguation works well but for information extraction and some other applications this is an overkill, and some senses may be lumped together.

A simple sentence or question can usually be briefly described by an action and an object; for example, the main idea from the sentence *He has to investigate all the reports* can be described by the action-object pair *investigate-report*. Even the phrase may be ambiguous by having a poor context, still the results of a search or interface based on such a sentence can be improved if the possible associations between the senses of the verb and the noun are determined.

In WordNet (Miller 1990), the gloss of a verb synset provides a noun-context for that verb, i.e. the possible nouns occurring in the context of that particular verb. The glosses are used here in the same way a corpus is used.

3 Ranking the possible senses of the noun

In order to improve the precision of determining the conceptual density between a verb and a noun, the senses of the noun should be ranked, such as to indicate with a reasonable accuracy the first possible senses that it might have.

The approach we considered for this task is the use of unsupervised statistical methods on large texts. The larger the collection of texts, the bigger is the probability to provide an accurate ranking of senses. As the biggest number of texts electronically stored - and thus favoring an automatic processing - is contained on the Web, we thought of using the Internet as a source of corpora for ranking the senses of the words.

This first step of our method takes into consideration verb-noun pairs $V - N$, and it creates pairs in which the verb remains constant, i.e. V , and the noun is replaced by the words in its similarity lists. Using WordNet, a similarity list is created for each sense of the noun, and it contains: the words from the noun synset and the words from the noun hypernym synset.

Algorithm 1

Input: untagged verb - noun pair

Output: ranking of noun senses

Procedure:

1. *Form a similarity list for each noun sense.*
Consider, for example, that the noun N has m senses. This means that N appears in m similarity lists,
 $(N^1, N^{1(1)}, N^{1(2)}, \dots, N^{1(k_1)})$
 $(N^2, N^{2(1)}, N^{2(2)}, \dots, N^{2(k_2)})$
 \dots
 $(N^m, N^{m(1)}, N^{m(2)}, \dots, N^{m(k_m)})$
 where N^1, N^2, \dots, N^m represent the different senses of N , and $N^{i(s)}$ represents the synonym number s of the sense N^i of the noun N as defined in WordNet.
2. *Form verb - noun pairs.* The pairs that may be formed are:
 $(V - N^1, V - N^{1(1)}, V - N^{1(2)}, \dots, V - N^{1(k_1)})$
 $(V - N^2, V - N^{2(1)}, V - N^{2(2)}, \dots, V - N^{2(k_2)})$
 \dots
 $(V - N^m, V - N^{m(1)}, V - N^{m(2)}, \dots, V - N^{m(k_m)})$
3. *Search the Internet and rank senses.* A search performed on the Internet for each of these groups will indicate a ranking over the possible senses of the noun N .

In our experiments we used (AltaVista) since it is one of the most powerful search engines currently available.

Verb	Noun	Sense of noun in SemCor	Hits provided by AltaVista for V-N where N has the sense no.:								Result
			#1	#2	#3	#4	#5	#6	#7	#8	
rescind	action	6	9	49	0	0	2	27	1	0	1
set-aside	resolution	7	5	0	75	7	0	0	17	2	2
reject	amendment	1	48								1
allow	legislator	1	172								1
ask	person	1	15628	0	1912	0					1
endorse	support	8	101	162	31	13	361	3	0	134	4
expend	fund	1	846	123							1
provide	increase	2	110189	5268	4429	1543	0				2
defeat	person	1	340	0	428						2
wait	term	2	0	27321	1	0	762				1
receive	vote	1	1271	0	0	406	62				1
revise	law	2	224	2829	648	640	37	397	0		1
expect	resignation	3	12	0	554						1
comment-on	topic	2	1801	5517							1
hold	meeting	1	205	128	8	1164	20	69	2227		3
remedy	problem	2	107	345	266						1
place	burden	4	2327	2031	12842	3271					2
award	fee	1	1	1284							2
award	compensation	1	22	126							2
protect	court	1	2574	3120	360	540	916	722	433		2

Table 1: A sample of the result we obtained in ranking the noun senses using the Internet

Using the operators provided by AltaVista, the verb-noun groups derived above can be expressed in two query-forms:

- (a) ("V* N¹*" OR "V* Nⁱ⁽¹⁾*" OR "V* Nⁱ⁽²⁾*" OR ... OR "V* N^{i(k)}*")
 (b) ((V* NEAR N¹*) OR (V* NEAR Nⁱ⁽¹⁾*) OR (V* NEAR Nⁱ⁽²⁾*) OR ... OR (V* NEAR N^{i(k)}*))

where the asterisk (*) is used as a wildcard indicating that we want to find all words containing a match for the specified pattern of letters.

Using one of these queries, we can get the number of hits for each sense i of the noun and this provides a ranking of the m senses of the noun as they relate with the verb V .

We tested this method for 80 verb-noun pairs extracted from SemCor 1.5 of the Brown corpus.¹

Using query form (a) as an input to the search engine, we obtained an accuracy of 83% in providing a ranking over the noun senses, such as the sense indicated in SemCor was one of the first two senses in this classification. In Table 1, we present a sample of the results we obtained. The column *Result* in this table presents the ranking over the noun senses: a 1 in this column means that the sense indicated in SemCor was also indicated by our method; 2 means that the sense indicated in SemCor was in top two of the sense ranking provided by our method; similarly, 3 or 4 indicates that the sense of the noun, as specified in SemCor, was in the top three, respectively four, of this sense ranking.

We used also the query form (b), but the results we obtained have been proved to be similar; using the operator *NEAR*, a bigger number of hits is reported, but the sense ranking remains the same.

It is interesting to observe that even we are creating queries starting with a verb-noun pair, it is

¹These verb-noun pairs have been extracted from the file br-a0...

not guaranteed that the search on the web will identify only words linked by such a lexical relation. We based our idea on the fact that: (1) the noun directly following a verb is highly probable to be an object of the verb (as in the expression "Verb* Noun*") and (2) for our method, we are actually interested in determining possible senses of a verb and a noun that can share a common context.

4 Determining the conceptual density between verbs and nouns

A measure of the relatedness between words can be a knowledge source for several decisions in the NLP applications. The conceptual density between verbs and nouns seems difficult to determine, without large corpora or a without a machine-readable dictionary having semantic links between verbs and nouns. Such semantic links can be traced however if we consider the glosses for the verbs, which are providing a possible context of a verb.

Algorithm 2

Input: untagged verb - noun pair and a ranking of noun senses (as determined by Algorithm 1)

Output: sense tagged verb - noun pair

Procedure:

- Given a verb-noun pair $V - N$, determine all the possible senses for the verb and the noun, by using WordNet. Let us denote them by $\langle v_1, v_2, \dots, v_k \rangle$ and $\langle n_1, n_2, \dots, n_l \rangle$ respectively.
- Using the method described in section 3, the senses of the noun are ranked. Only the first two possible senses indicated by this step will be considered.
- For each possible pair $v_i - n_j$, the conceptual density is computed as follows:

- (a) extract all the glosses from the sub-hierarchy including v_i (the rationale of the method used to determine these sub-hierarchies is explained below)
 - (b) Determine the nouns from these glosses. These constitute the noun-context of the verb. All these nouns are stored together with the level of the associated verb within the sub-hierarchy of v_i .
 - (c) Determine the nouns from the sub-hierarchy including n_i .
 - (d) Determine the number C_{ij} of common concepts between the nouns obtained at (b) and the nouns obtained at (c).
4. The most suitable combinations between the senses of the verb and the noun $v_i - n_j$ are the ones that provide the biggest values for C_{ij} .

In order to determine the sub-hierarchies that should be used for v_i and n_j , we used statistics provided by SemCor, a sense tagged version of the Brown corpus (Francis and Kucera, 1967) (Miller, Leacock et al., 1993), containing 250,000 words. Each word (noun, verb, adjective, adverb) is included in a synset within a hierarchy. The tops of these hierarchies denominate the class of the word. The sense in SemCor for a word W is indicated by the class C of the word W , and the sense of the word within the class C . For example, the SemCor entry:

```
<wf cmd=done pos=NN lemma=investigation vnsn=1
lexsn=1:09:00::>investigation</wf>
```

indicates:

word: investigation
part of speech: common noun
sense in WordNet: 1

A statistic measure performed on SemCor, indicates the following probabilities for the sense of a word within a class:

Parts of speech	Number of words in SemCor	within a class, the probability to have sense number:					
		0	1	2	3	4	5
noun	47,799	85%	10%	3%	1%	-	-
verb	27,637	60%	14%	5%	12%	3%	2%

Table 2: The probabilities for the sense of a word within a class

As shown in Table 2, the class of the noun indicates with a probability of 85% a correct sense 1 within that class.

Thus, for this algorithm, we consider for a noun the hierarchy including the noun (if the class of the noun n_i is C , then the method considers all the nouns from the class C).

This does not work for the verbs, as the probability to indicate a correct sense knowing the class is much smaller (only 60%). For this reason, and based on the experiments we computed, the sub-hierarchy

including a verb v_i is determined as follows: (i) consider the hypernym h_i of the verb v_i and (ii) consider the hierarchy having h_i as top.

It is necessary to consider a bigger hierarchy than just the one provided by synonyms and direct hyponyms, since providing accuracy in a metric computation needs large corpora. As we replaced the corpora with the glosses, better results are achieved if more glosses are considered. Still, we do not have to enlarge too much the context, in order not to miss the correct answers.

Conceptual Density Metric

For determining the conceptual density between a noun n_i and a verb v_j , the algorithm considers:

- the list of nouns sv_k associated with the glosses of the verbs within the hierarchy determined by h_j : (sv_k, w_k) , where:
 - h_j is the hypernym of v_j
 - w_k is the level in this hierarchy
- the list of nouns sn_l within the class of n_i : (sn_l)

The common words between these two lists (sv_k, w_k) and (sn_l) will produce a list of common concepts with the associated weights $cd_{ij} < w_k >$. The conceptual density between n_i and v_j is given by the formula:

$$(1) \quad C_{ij} = \frac{\sum_k^{|cd_{ij}|} w_k}{\log(desc_i)}$$

where:

- $|cd_{ij}|$ is the number of common concepts between the hierarchies of n_i and v_j
- w_k are the weights associated with the nouns from the noun-context of the verb v_j
- $desc_i$ is the total number of words within the hierarchy of noun n_i

As the nouns with a big hierarchy tend to indicate a big value for $|cd_{ij}|$, the weighted sum of common concepts has to be normalized in respect with the dimension of the noun hierarchy. This is estimated as the logarithm of the total number of descendants in the hierarchy (i.e. $\log(desc_i)$).

We also took into consideration other metrics, like:

- (2) The number of common concepts between the noun and verb hierarchies, without considering the weights.
- (3) A weighted summation of the common concepts between the noun and verb hierarchies, as indicated in (1), but without a normalization in rapport with the noun hierarchy.

We considered also the metrics indicated in (Agirre and Rigau, 1995). But after running the program on several examples, the formula indicated in (1) provided the best results.

A possible improvement to the metric (1) is to consider the weights for the levels in the noun hierarchy, in addition to the levels in the verb hierarchy.

5 An example

Consider as example of a verb-noun pair the phrase *revise law*. The verb *revise* has two possible senses in WordNet 1.5:

Sense 1

revise, make revisions in

gloss: (revise a thesis, for example)

⇒ rewrite, write differently, alter by writing

gloss: ("The student rewrote his thesis")

Sense 2

re tool, revise

⇒ reorganize, shake up, organize an

The noun *law* has 7 possible senses

Sense 1

law, jurisprudence

gloss: (the collection of rules imposed by authority; "civilization presupposes respect for the law")

⇒ collection, aggregation, accumulation, assemblage

gloss: (several things grouped together)

Sense 2

law

gloss: (one of a set of rules governing a particular activity or a legal document setting forth such a rule; "there is a law against kidnapping")

⇒ rule, prescript

gloss: (prescribed guide for conduct or action)

⇒ legal document, legal instrument, official document, instrument

Sense 3

law, natural law

gloss: (a rule or body of rules of conduct inherent in human nature and essential to or binding upon human society)

⇒ concept, conception

gloss: (an abstract or general idea inferred or derived from specific instances)

Sense 4

law, law of nature

gloss: (a generalization based on recurring facts or events (in science or mathematics etc): "the laws of thermodynamics")

⇒ concept, conception

gloss: (an abstract or general idea inferred or derived from specific instances)

Sense 5

jurisprudence, law, legal philosophy

gloss: (the branch of philosophy concerned with the law)

⇒ philosophy

gloss: (the rational investigation of questions about existence and knowledge and ethics)

Sense 6

police, police force, constabulary, law

gloss: (the force of policemen and officers; "the law came looking for him")

⇒ force, personnel

gloss: (group of people willing to obey orders)

Sense 7

law, practice of law

gloss: (the learned profession that is mastered by graduate study in a law school and that is responsible for the judicial system; "he studied law at Yale")

⇒ learned profession

gloss: (one of the three professions traditionally believed to require advanced learning and high principles)

We searched on Internet, using AltaVista, for all possible pairs V-N that may be created using *revise* and the words from the similarity lists of *law*. Over the seven possible senses for this noun, the first step of our method indicated the following ranking (we indicate the number of hits between parenthesis): *law*#2(2829), *law*#3(648), *law*#4(640), *law*#6(397), *law*#1(224), *law*#5(37), *law*#7(0). Thus, only the sense #2 and #3 of the noun *law* are eligible to be used for the next algorithm.

For each of the two senses of the verb, we determined the noun-context, including the nouns from the glosses in the sub-hierarchy of the verb, and the associated weights.

For each of the two possible senses of the noun, we determined the nouns from the class of each sense.

In Table 3, we present: (1) the values obtained for the combinations of different senses, i.e. the number of common concepts between the verb and noun hierarchies - $|cd_{ij}|$ (columns 2-3); (2) the summations of the weights associated with each noun within the noun-context of the verb v_j (columns 4-5); (3) the total number of nouns within the hierarchy of each sense n_i , i.e. $desc_i$ (columns 6-7); (4) the conceptual density C_{ij} for each pair $n_i - v_j$, derived using the formula presented above (columns 8-9).

	$ cd_{ij} $		weights		$desc_i$		C_{ij}	
	2	3	4	5	6	7	8	9
	n_2	n_3	n_2	n_3	n_2	n_3	n_2	n_3
v_1	5	4	2.06	2	975	1265	0.30	0.28
v_2	0	0	0	0	975	1265	0	0

Table 3: Values used in computing the conceptual density and the conceptual density C_{ij}

In this table:

- v_i indicates the sense number i of verb *revise*
- n_i indicates the sense number i of noun *law*

The biggest value for conceptual density is given by $v_1 - n_2$:

$$revise\#1/2 - law\#2/5 \quad C_{11} = 0.30$$

This combination of verb-noun senses² appears in SemCor, file br-a01.

6 Tests against SemCor

We tested this method by using verb-noun pairs from SemCor. A randomly selected sample from the entire table with 80 pairs is presented in Table 4.

For each pair verb-noun, we indicate the sense of the verb (column B), the sense of the noun (column C), as they result from SemCor; the total number of possible senses for both the verb (column D)

²The notation $\#i/n$ means sense i out of n possible.

these cases, the *NEAR* operator should be used for the first step of this algorithm).

2. The number of words considered at a time can be increased, from two to three, four or even more words.

7 Conclusion

In this paper, we have presented a method for WSD that is based on measuring the conceptual density between words using WordNet. The metric proposed may be further improved by considering the weights for verbs as well as for nouns. The senses of the words are ranked, and an user may select the first choice or the first few choices, depending upon the application. We have also proposed to use the Internet as a source of statistics on a raw corpora.

The method extends well to considering more than two words at a time, and also for all parts of speech covered by WordNet.

It is difficult to compare the precision obtained by this method with other methods, since we consider here the collective meaning of two or more words, while most of other methods consider one word at a time. However, an estimation can be done by extracting the square root of the accuracy for a pair of verb-noun words; and that is 76.15% for the first choice, 83.66% for the first two choices and 85.44% for the first three choices. Since the disambiguation precision for nouns is usually higher than for verbs, those numbers provide only an average.

References

- Digital Equipment Corporation. AltaVista Home Page. URL:<http://www.altavista.digital.com>.
- E. Agirre and G. Rigau, A Proposal for Word Sense Disambiguation using Conceptual Distance, Proceedings of the 1st International Conference on Recent Advances in Natural Language Processing, Velingrad, 1995
- R. Bruce and J. Wiebe, Word Sense Disambiguation using Decomposable Models, Proceedings of the 32nd Annual Meeting of the Association for Computational Linguistics (ACL-94), 139-146, LasCruces, NM, June 1994.
- J. Cowie, L. Guthrie and J. Guthrie, Lexical disambiguation using simulated annealing. *Proceedings of the Fifth International Conference on Computational Linguistics COLING-92*, 157-161, 1992.
- S. Francis and H. Kucera, *Computational Analysis of present-day American English*, Providence, RI: Brown University Press, 1967
- W. Gale, K. Church and D. Yarowsky, One Sense per Discourse, Proceedings of the DARPA Speech and Natural Language Workshop, Harriman, New York, 1992.
- X. Li, S. Szpakowicz and S. Matwin. A WordNet-based algorithm for word semantic sense disambiguation. Proceedings of the 14th International Joint Conference on Artificial Intelligence IJCAI-95, Montreal, Canada, 1995.
- S. McRoy, Using multiple Knowledge Sources for Word Sense Disambiguation, *Computational Linguistics*, 18(1):1-30, 1992.
- G.A. Miller, WordNet: An on-line lexical database. *International Journal of Lexicography*, 3(4):235-312, 1990.
- G.A. Miller, C. Leacock, T. Randee and R. Bunker, A Semantic Concordance. Proceedings of the 3rd DARPA Workshop on Human Language Technology, 303-308, Plainsboro, New Jersey, 1993
- G.A. Miller, M. Chodorow, S. Landes, C. Leacock and R.G. Thomas, Using a semantic concordance for sense identification. *Proceedings of the ARPA Human Language Technology Workshop*, 240-243, 1994.
- H.T. Ng and H.B. Lee, Integrating Multiple Knowledge Sources to Disambiguate Word Sense: An Exemplar-Based Approach, Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics (ACL-96), Santa Cruz, 1996.
- R. Rada, H. Mili, E. Bickell and M. Blettner. Development and Application of a Metric on Semantic Nets. *IEEE Transactions on Systems, Man and Cybernetics*, vol. 19, pp 17-30, Jan/Feb 1989.
- P. Resnik, Selectional Preference and Sense Disambiguation, Proceedings of ACL Siglex Workshop on Tagging Text with Lexical Semantics, *Why, What and How?*, Washington, April 4-5, 1997.
- P. Resnik and D. Yarowsky, A Perspective on Word Sense Disambiguation Methods and Their Evaluation, Proceedings of ACL Siglex Workshop on Tagging Text with Lexical Semantics, *Why, What and How?*, Washington, April 4-5, 1997.
- R. Richardson, A.F. Smeaton and J. Murphy, Using WordNet as a Knowledge Base for Measuring Semantic Similarity between Words, Technical Report, Working paper CA-1294, School of Computer Applications, Dublin City University, Dublin, Ireland, 1994.
- G. Rigau, J. Atserias and E. Agirre. Combining Unsupervised Lexical Knowledge Methods for Word Sense Disambiguation. *Computational Linguistics /9704007*, 1997.
- D. Yarowsky, Unsupervised Word Sense Disambiguation rivaling Supervised Methods. Proceedings of the 33rd Association of Computational Linguistics, 1995.

Constructing Bayesian Networks from WordNet for Word-Sense Disambiguation: Representational and Processing Issues *

Janyce Wiebe and Tom O'Hara
Department of Computer Science and
Computing Research Laboratory
New Mexico State University
Las Cruces, NM 88003
wiebe, tomohara@cs.nmsu.edu

Rebecca Bruce
Department of Computer Science
University of North Carolina at Asheville
Asheville, NC 28804-3299
bruce@cs.unca.edu

Abstract

This paper describes a probabilistic model that is formed from the integration of an analytical and empirical component. The analytical component is a Bayesian network derived from WordNet, and the empirical component is composed of compatible probabilistic models formulated from tagged training data. The components are integrated in a formal, uniform framework based on the semantics of causal dependence. The paper explores various representational issues that must be addressed when formulating a Bayesian network representation of lexical information such as that expressed in WordNet. These issues are essential to the design of such a network and they have not been previously explored. We describe two choices for the representation of lexical items and two choices for the representation of lexical relations. The effect of each combination of choices on evidence propagation in the network is discussed.

1 Introduction

There is a long tradition in AI of resolving interdependent lexical ambiguities through spreading activation, from Quillian's (1968) seminal work on semantic networks, through Hirst's work (1988) on Polaroid words, to more recent work by Voorhees (1993) and Veronis and Ide (1990) on large-scale disambiguation. This research investigates a probabilistic realization of spreading activation to resolve interdependent word-sense ambiguities. The core idea is to exploit belief propagation in Bayesian networks: Words are mapped to nodes, lexical relations are mapped to edges, and evidence is propagated from word senses to other related word senses.

The lexical relations are derived from an existing knowledge source, because this information cannot be automatically extracted from training data with existing techniques. The knowledge source we use is the WordNet *is-a* hierarchy, i.e., the *hypernym/hyponym* taxonomy (Miller, 1990). Although this hierarchy was developed for other purposes, it

has been frequently applied to word-sense disambiguation (Resnik, 1995; Sussna, 1993). In this work, we investigate various approaches to constructing a Bayesian network representation of the *is-a* hierarchy for use in word-sense disambiguation. As this work continues, other relations such as part/whole and entailment relations will also be included in the network.

Another contribution of our work is a novel proposal for integrating symbolic and statistical information for the purpose of performing NLP tasks. Statistical approaches to word-sense disambiguation have had the most success to date, when evaluated on unseen test data. The "analytical" Bayesian network component of our method is actually built on top of "empirical" probabilistic classifiers induced statistically from training data. In particular, an empirical classifier is induced for each word in the current sentence to be disambiguated (i.e., for each *target word*). Each empirical classifier is developed independently of the empirical classifiers for other target words. A Bayesian network is constructed from the segment of the WordNet *is-a* hierarchy that is connected to the target words. The results of the empirical classifiers are fed as evidence into the Bayesian network, thus initiating belief propagation. All of the information is represented in a formal, uniform framework: a probabilistic model embodying conditional independence relationships among the variables that form the joint distribution. Conditional independence relationships simplify the formulation of the joint distribution making it possible to work with a large number of variables. Further, models that characterize conditional independence relationships have desirable computational properties (e.g., see the discussion on decomposable models in (Pearl, 1988)). These properties form the basis of the evidence propagation scheme used for Bayesian networks discussed in Section 7. We also make use of these properties in formulating the empirical classifiers as described in (Bruce and Wiebe, 1994). Bayesian networks are a very rich and complex representational framework. They support easy integration of diverse information sources and form

* This research was supported in part by the Office of Naval Research under grant number N00014-95-1-0776.

the basis for much of the current work on reasoning under uncertainty (Pearl, 1988).

This paper explores the representational issues that must be addressed when mapping the lexical information in WordNet to a Bayesian network. The implications of the various choices are analyzed in depth. In section 2, we introduce the basic concepts and illustrate them with an example in section 3, which also includes a brief description of the empirical component. The Bayesian network representations of lexical items and lexical relations are discussed in sections 4 and 5, respectively. In section 6, we describe the integration of the empirical component into the Bayesian network. The process of sense disambiguation is described in section 7. Section 8 discusses related work followed by conclusions in section 9.

2 Bayesian Networks: Background

Bayesian networks model dependencies among nodes through the use of conditional probabilities. Specifically, if a node (*Cause2*) is considered as a cause for another node (*Symptom1*), then the second node is defined relative to the first (i.e., $P(\textit{Symptom1}|\textit{Cause2})$). Some nodes don't have associated causes, so they are just defined via unconditional probabilities (e.g., $P(\textit{Cause2})$). Taken together, the set of all the conditional and unconditional probabilities determine a joint distribution for all the nodes being modeled (e.g., $P(\textit{Symptom1}, \dots, \textit{SymptomN}, \textit{Cause1}, \dots, \textit{CauseM})$). Such global distributions are usually difficult to assess directly; hence, the Bayesian network provides a convenient formalism for specifying the same distribution via local distributions, under conditional independence assumptions. Furthermore, without the conditional independence relations, the full joint distribution for cases with hundreds of senses would be infeasible to process—the independence assumptions are key. Pearl (1988) presents an in-depth coverage of the theory of Bayesian networks and provides an efficient algorithm for evaluating them.

In a Bayesian approach to statistical inference, we distinguish between *prior* and *posterior* probabilities. Prior probabilities express the beliefs that we hold about the likelihood of events *prior* to being given any evidence, *posterior* probabilities express our beliefs in the likelihood of events given all the evidence that is currently known. Thus, the posterior probability of an event changes as new evidence is learned. The conditional and unconditional probabilities mentioned above are the prior probabilities. The posterior probabilities are calculated using the Bayesian network propagation algorithm each time new evidence is added. We discuss propagation in greater detail in Section 7. Intuitively, the posterior probability of a node, say the node GATHERING#1

(switching to a word-sense disambiguation example), is a combination of the beliefs received from its children and the beliefs received from its parents. Once a node has calculated its own belief, it calculates outgoing messages to send to its parents and to its children, which enable them, in turn, to calculate their posterior probabilities. In this way information is propagated throughout the network.

3 An Example

In this section, we illustrate how a simple Bayesian network can be constructed to model the interdependencies among words. This identifies the basic steps in the overall process and helps to motivate the representational issues discussed later.

Suppose that the words “community” and “town” appear in a single sentence, and that their correct senses in that context are COMMUNITY#1 and TOWN#2, respectively. Our task is to assign the correct word senses to both of them, considering information automatically derived from the corpus and gathered individually for each word, as well as information derived from the WordNet *is-a* hierarchy and represented in a Bayesian network. The basic strategy is to add the corpus-derived information to the Bayesian network representations of “community” and “town,” in such a way that it initiates propagation.

Let us consider this process in more detail. The words “community” and “town” have the following senses in WordNet:

community:

1. people living in a particular local area
2. an association of people with similar interests
3. common ownership
4. the body of people in a learned occupation

town:

1. an urban area with a fixed boundary that is smaller than a city
2. the people living in a municipality smaller than a city
3. an administrative division of a county

These senses are represented as sets of synonyms, or *synsets*. In the *is-a* hierarchy, each synset is linked to its *hypernym*, i.e., the synset representing its conceptual parent. For example, the synset corresponding to

{occupation, vocation, occupational group}
is the hypernym of the synset corresponding to
{profession, community}.

A new Bayesian network is created for each sentence. It includes all of the synsets for the target words in the sentence, together with all of the

synsets reachable from them in the WordNet *is-a* hierarchy. Extracting this information from WordNet is straightforward.

Figure 1 illustrates one way that the Bayesian network for the example sentence containing “town” and “community” can be constructed. In this representation, each word sense is mapped to a node in the network, and there is an edge from X to Y iff word sense X is a hypernym (i.e., a superordinate) of word sense Y (please ignore the octagonal nodes at the bottom for now). Notice that the relation between COMMUNITY#1 and TOWN#2 is mediated by GATHERING#1, a type of GROUP#1. Our goal is for the contextual evidence provided by the empirical classifiers to propagate along this path in such a way that the correct senses of the target words reinforce one another.

After the topology of the network has been established, the conditional probability tables required for each node must be defined. As will be discussed later in section 5, we can make independence assumptions that make estimating the necessary probabilities more easier.

Next, an empirical classifier is developed for each ambiguous word, in this case, “town” and “community”. Each classifier defines a probability distribution describing the likelihood of each sense of the targeted word given the automatically derived features of the context. An example of the type of feature used is the part-of-speech of the word to the right; see (Bruce and Wiebe, 1994) for the other ones we use.

The distributions determined by the empirical classifiers are added as evidence to the Bayesian network, initiating belief propagation. Once the network reaches equilibrium, the posterior probabilities of the nodes for “town” and “community” determine the senses assigned to each ambiguous word.

4 Representing Lexical Items: What does a Node mean?

There are two basic approaches to representing WordNet synsets in a Bayesian network. Since the lexical relations are among synsets and not words, a natural approach is to represent the synsets as nodes. Alternatively, one node could be used to represent all senses of a word.

4.1 The One Node Per Word Approach

When nodes correspond to words, the possible values for each node are *sense*₀ through *sense* _{N} , where N is the number of WordNet synsets representing senses of the target word. *Sense*₀ represents the composite of all other meanings, i.e., of all meanings that are not represented by WordNet synsets. Figure 2 shows the graph for the Bayesian network when word nodes are used for the relations. It also illus-

trates the use of logical links, which are described in the next section. This involves more than just a change in link direction.

4.2 The One Node Per Sense Approach

Figure 1 illustrates the approach in which each synset (each sense) is mapped to a node. An important advantage of using the node per sense approach is that it facilitates handling dependencies among the senses of a word. In the node per word approach, single node cycles are produced when modeling the dependencies of words that have a meaning that is defined in terms of other meanings for that same word.

A disadvantage of this approach is that modeling mutual exclusion among the senses of a single word becomes more difficult. The most straightforward approach modeling mutual exclusion is to create a dependency from each sense node to a separate node with a CPT enforcing mutual exclusion. But since the table must have 2^N entries, this approach becomes impractical for words with a large number of senses. To get around this problem, two levels of mutual-exclusion dependencies could be introduced: one at which mutual exclusion among small groups of senses is enforced, and another enforcing mutual exclusion of the groups.

5 Representing Lexical Relations: What does an edge mean?

Here, we address issues concerning the representation of WordNet *is-a* relationships as causal dependencies. The two primary issues to be addressed are: (1) expressing the Hypernym/Hyponym relationship as a causal dependency, and (2) quantifying the causal dependencies with conditional probability distributions.

5.1 Hypernym→Hyponym Representations

The Hypernym → Hyponym Representation was illustrated above in section 3: there is an edge from node X to node Y iff X represents a hypernym of node Y in the WordNet *is-a* hierarchy. Consider the node per sense representation (see figure 1). Suppose *Hyper* is a synset that is a hypernym of synsets *Hypo*₁ ... *Hypo* _{k} . Then, the relevant part of the Bayesian network expresses the following:

$$Hyper \rightarrow Hypo_1 \vee \dots \vee Hypo_k$$

As such, we are making a closed world assumption. If, for example, there is a synset ANIMAL#1 with three hyponyms DOG#1, CAT#1, and MOUSE#1, we are assuming that these three are the only kinds of ANIMAL#1's there are.

When using this link representation with either of the node per sense or the node per word representations, the roots of the network are the most superordinate synsets reachable from the target words, and

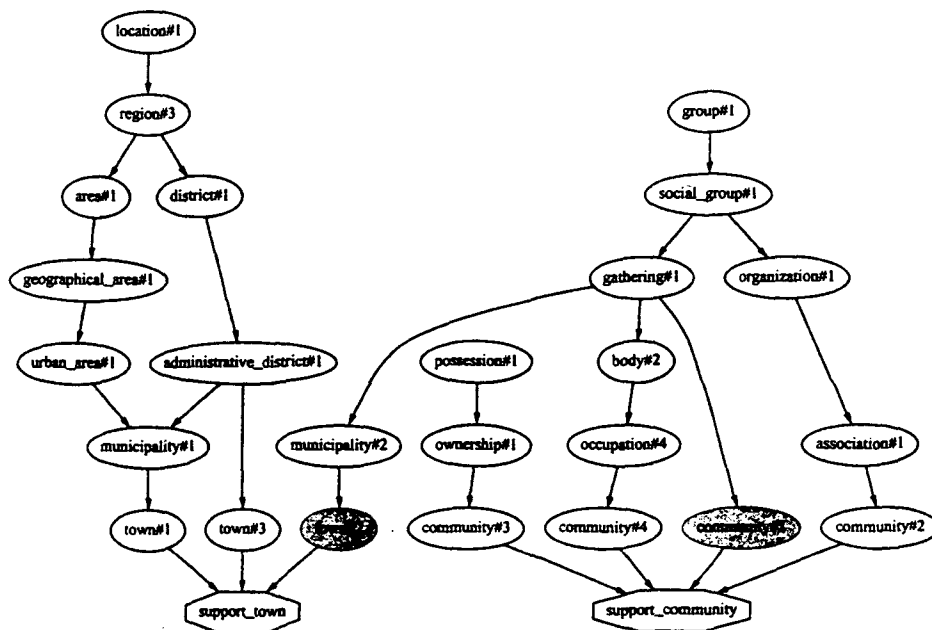


Figure 1: Sense per node Bayesian network with hypernym→hyponym links

the target words are typically (but not necessarily) the leafs of the network.

We now turn to defining the CPT. We discuss this with respect to the node per sense representation (in figure 1) because it is easier to discuss and similar conditional probabilities must be defined under the node per word representation.

To define the CPT for each child node in the Bayesian network, where each child node corresponds to a hyponym node in WordNet, we assign the conditional probability $P(\text{hyponym}|\text{hypernym})$ to be inversely proportional to the number of children that the hypernym has. For instance, MUNICIPALITY#1 has two children in WordNet, so we assign the following conditional probability for TOWN#1 given this hypernym.

P(town#1 municipality#1)	
municipality#1	P(town#1)
F	0.000 + ϵ
T	0.500

In so doing we are: (1) considering each hyponym of a given hypernym to be equally likely, and (2) maintaining the closed world assumption by requiring that these conditional probabilities sum to one. In all CPTs, we add a small positive probability ϵ to all zero probability values in order to allow the realization of all possible configurations of node values (e.g., to handle inconsistent evidence). In future

work, we will consider using frequency of occurrence information in tagged training data to define these CPTs.

For the root nodes, which represent the most superordinate concepts, prior probabilities must be specified. With no evidence to the contrary, uniform prior distributions are assigned to the root nodes; the empirical classifiers are relied upon to provide contextual support (through the leafs of the network).

5.2 Hyponym→Hypernym Representations

Under the Hyponym → Hypernym Representation, there is an edge from node X to node Y iff X represents a hyponym of node Y in the WordNet *is-a* hierarchy. Consider the node per sense representation (see figure 2). The Bayesian network represents the following:

$$\begin{aligned}
 &(\text{Hypo}_i = s_i \rightarrow \text{Hyper}_j = s_j) \\
 &\wedge \dots \wedge (\text{Hypo}_n = s_n \rightarrow \text{Hyper}_m = s_m)
 \end{aligned}$$

Under the semantics of the WordNet *is-a* hierarchy, all instances of a hyponym are instances of its hypernym. So, a typical CPT for this representation is as follows:

P(municipality#1 town#1)	
town	P(municipality#1)
F	0.0 + ϵ
T	1.0 - ϵ

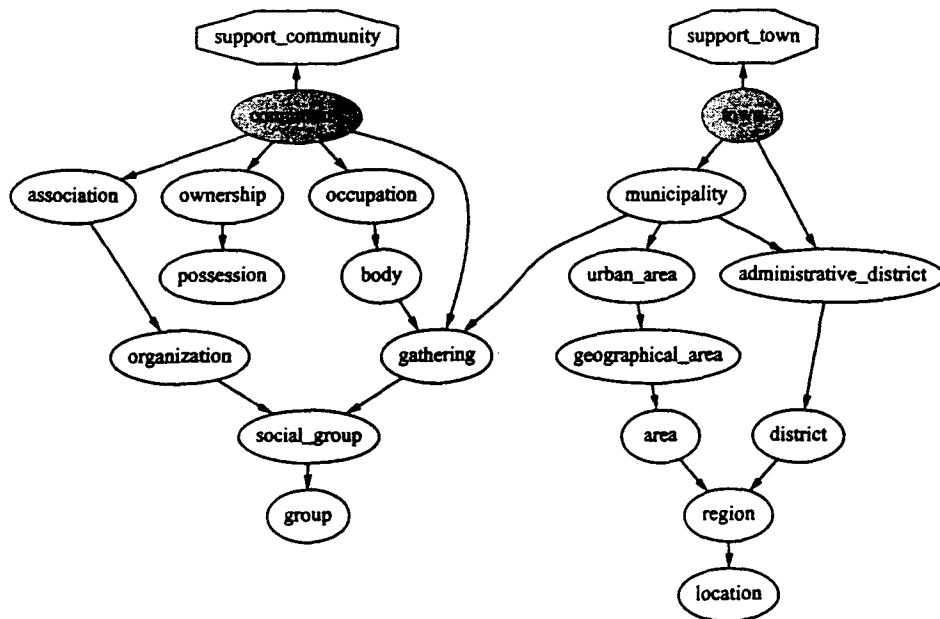


Figure 2: Word per node Bayesian network with hyponym→hyponym links

Note that this case is not illustrated in the graphs shown: these only cover two of the four main possibilities.

Interestingly, in this representation, the root nodes represent the target words. Thus, the root nodes are the sites where evidence from the empirical classifiers is added to then network. In the absence of this evidence, these nodes take on their prior probabilities. As above, we assign uniform distributions as the priors. Recall that, in the case of multiple parents, CPTs must specify the conditional distribution of the child node given the values of all of its parent nodes. The issues involved in working with multiple parent nodes are discussed below.

5.3 CPT Entries when Multiple Parents: Causal Independence

If a node has multiple parents, say n parents, then specifying all of the entries in the CPT for that node can be prohibitive. If no additional independence assumptions are made regarding the interactions among the parent nodes, then the number of probabilities that must be specified is exponential in n , and probabilistic inference is made correspondingly more complex (Heckerman and Breese, 1994). To overcome this problem, the *noisy-OR* model (Pearl, 1988) is often adopted. Under this model, certain independence assumptions are made regarding the interactions among the parent nodes, with the effect that the number of probabilities that

must be specified is linear in n . Basically, one need only specify the conditional probabilities of the child and each parent individually.

As presented in (Pearl, 1988), the noisy-OR model assumes that all of the variables are binary. Heckerman and Breese (1994) present a generalization of the noisy-OR model, *causal independence*. In this model, the parents are assumed to be independent causes for the child. This allows us to formulate a CPT from the specification of only the following conditional probabilities: $P(c|p_{ij})$, where c ranges over the values of the child, and p_{ij} ranges over the values of parent P_i . These values are combined via the constraints of the model to produce the CPT for the child node.

We assign the probabilities using a causal independence model which specializes to the noisy-OR model when applied to binary nodes. First consider that the inclusive-or connective can be viewed as outputting a true value iff none of the inputs is false:

$$\text{output} = \neg((\neg v_1) \wedge \dots \wedge (\neg v_n))$$

where each v_i is a logical-valued input variable. The extension to the case where probabilities are associated with each input is relatively straightforward:

$$\begin{aligned} \text{child} &= \neg((\neg v_1) \wedge \dots \wedge (\neg v_n)) \\ P(\text{child}|V_1 = v_1, \dots, V_n = v_n) &= \\ &= 1.0 - \prod(1.0 - P(\text{child}|v_i)), \quad \forall v_i v_i = T. \end{aligned}$$

When extending to the general case, the relationship between the value of the child node and the values of its parent nodes is not necessarily defined by a

truth function. But, the probabilities are assigned analogously:

$$P(\text{Child} = c | V_1 = v_1, \dots, V_n = v_n) = 1.0 - \prod (1.0 - P(\text{child} = c | V_i = v_i)) \forall V_i P(\text{child} = c | V_i = v_i) > \epsilon.$$

In their work on plan recognition, Charniak and Goldman (1993) use the noisy-OR model, specifically for representing the dependencies of observed actions on the potential plans that could explain them.

6 Integrating Empirical and Analytical Information: Virtual Evidence Nodes

Due to space limitations, we consider just one method for integrating the empirical and analytical components. In this technique, support from the empirical classifiers is added to the Bayesian network using *virtual evidence nodes* (Pearl, 1988). The usual way to add evidence to a Bayesian network is to instantiate a node to a particular value (called “clamping”); the influence of this evidence is then propagated through the network. However, that method is not appropriate for our task, because we do not know the sense of any word (so there is no node in the Bayesian network that can be initially instantiated). Virtual evidence nodes provide a way to specify uncertain evidence, in the form of a distribution over node values (i.e., the probability of each node value). They are represented by the octagonal nodes in figures 1 and 2. There is one for each of the target words to be disambiguated. These nodes represent the support for each sense that was derived from the corpus by the empirical component. Each virtual evidence node is implemented as a binary-valued node whose parent is the node for which evidence is being provided. The evidence distribution determines the conditional probability table.

7 Edge Direction and Belief Propagation

There is a very important implication of the choice between the hypernym \rightarrow hyponym and the hyponym \rightarrow hypernym representations. In a Bayesian network, suppose that evidence is added to a node (either by clamping or by virtual evidence nodes). This evidence will propagate to its ancestors in the Bayesian network, and also to the children of its ancestors. For example, in figure 1, evidence introduced at node SUPPORT_COMMUNITY will propagate, among other places, back to COMMUNITY#1, back to GATHERING#1, and then down to MUNICIPALITY#2, and so on. Thus, this representation, hypernym \rightarrow hyponym, supports the kind of propagation described in this paper.

On the other hand, consider the hyponym \rightarrow hypernym representations (figures 1 and 2). In these

sense	before	after
community#1	.20	.70
gathering#1	.55	.87
municipality#2	.25	.33
town#2	.25	.33
community#4	.20	.10
body#2	.20	.10
location#1	.50	.67
municipality#1	.25	.33
town#1	.25	.33

Table 1: Propagation w/ hyponym \rightarrow hypernym links

representations, the targeted words are the roots of the Bayesian network, so the evidence is added to the roots of the network. This evidence will **not** propagate from, say, COMMUNITY#1 to TOWN#2 in figure 1. Information propagates between such nodes **only if evidence were added to their mutual descendants**. As Pearl says, “evidence gathered at a particular node does not influence any of its spouses until their common child gathers diagnostic support” ((Pearl, 1988), p. 182). Thus, if evidence is only added at the virtual evidence nodes in figure 2, evidence will not propagate from COMMUNITY=1 to MUNICIPALITY=2 (so it will not propagate further to TOWN=2). The corresponding nodes are spouses, but their child (GATHERING) has not received diagnostic support, by which Pearl means evidence propagated from below.

However, there are many other possibilities for adding evidence to the network, under which desired propagation would occur. Thus, our discussion of the hyponym \rightarrow hypernym representations is not just a cautionary tale. For example, one might use Yarowsky’s (1992) unsupervised method for assigning words to thesaural categories to add evidence to a node representing a superordinate concept in the WordNet *is-a* hierarchy. (Virtual evidence nodes could be used for this purpose too.) In the hyponym \rightarrow hypernym representations, this superordinate concept (say GATHERING#i or SOCIALGROUP#1) is a **descendent** of the nodes representing the targeted words. It would thus provide the needed diagnostic support to enable propagation from one target word to another. Note that the hyponym \rightarrow hypernym representation is conceptually appealing, since its semantics is based directly on the semantics of the WordNet *is-a* hierarchy.

As an illustration, consider applying sample evidence of (.70, .10, .10, .10) for the senses of “community” (with no evidence for town). Table 1 shows the posterior probabilities before and after applying this evidence.

As can be seen, the high evidence for COMMU-

sense	before	after
community#1	.054	.562
gathering#1	.126	.631
municipality#2	.063	.312
town#2	.060	.290
community#4	.020	.030
body#2	.064	.296
location#1	.500	.500
municipality#1	.068	.063
town#1	.053	.050

Table 2: Propagation w/ hypernym→hyponym links

NITY#1 increases the support for the hypernym GATHERING#1 (as well as for the other ancestors in the same path not shown). However, no support is reaching MUNICIPALITY#2.

If the hypernym → hyponym representation is used instead (as in figure 1), an appropriate propagation does take place. The propagation occurs in two phases. First, the high evidence for COMMUNITY#1 is propagated “upstream” to the hypernym node. Then, the increased support for this synset is propagated “downstream” to increase the likelihood of the value for the appropriate sense of “town”. Table 2 shows the posterior probabilities in this case.

7.1 Attenuation of Spreading Activation

An important aspect of spreading activation approaches is that the strength of the evidence being propagated is attenuated the further the evidence spreads from the original source. Traditional spreading activation schemes have used various heuristics to model this attenuation, often incorporating a distance factor in terms of number of links. By using probabilistic propagation, we can account for both length of path and fan-out at the nodes along the path (i.e., how many children they have). The length of the path is taken into account by the propagation algorithm. Intuitively, when a node calculates its posterior distribution, it calculates a distribution taking into account all possibilities (e.g., gathering#1=1, municipality#2=1; gathering#1=1, municipality#2=0; and so on). As the evidence is dispersed among the various possibilities at subsequent nodes, the evidence for any single possibility tends to decrease. This is so for either edge direction.

8 Comparison to Related Work

Spreading activation schemes have been common in various forms, starting with Quillian’s (Quillian, 1968) work on semantic memory. Quillian used spreading activation to identify paths between concepts for the purpose of comparison and contrast. To construct the semantic networks, dictionary def-

initions were manually encoded in the form a graph.

Hirst (1988) also used spreading activation to perform word-sense disambiguation. The approach relies on the identification of paths between interdependent word meanings. To avoid extraneous connections, constraints were introduced; for instance, a limit on path length was introduced, and *is-a* links were normally not traversed in reverse direction. Furthermore, heuristics were used to give preference to shorter paths and to avoid connections through nodes with many out-going arcs.

There have been several approaches that have relied upon word-overlap in dictionary definitions to resolve word-sense ambiguities in context, starting with (Lesk, 1986). Cowie et al. (1992) extend the idea by using simulated annealing to optimize a configuration of word senses simultaneously in terms of degree of word overlap.

Veronis and Ide (1990) developed a neural network model to overcome the limitation of addressing only pairwise dependencies in word-overlap approaches. Using dictionary definitions, they constructed a network containing links from each word node to the nodes for each of its senses and links from each of the sense nodes to the nodes of the words used in the definition.

Sussna (1993) produces a semantic network based on several different WordNet relations. His disambiguation method minimizes the pairwise distance among senses via a weighting scheme that accounts for both fan-out and depth in the hierarchy. Of the approaches we have surveyed, his is most similar to our analytical component.

Voorhees (1993) describes an unsupervised approach that exploits the WordNet hypernym taxonomy. In particular, the hierarchy for a given word is automatically partitioned so that the words occurring in the synsets of a partition (or *hood*) only occur with one of the senses for the word. Disambiguation is based on the selecting the hood which has the highest estimated relative frequency for the context relative to training text.

Resnik (1995) also describes an unsupervised approach that is based on estimating synset frequencies. As with Voorhees, the estimated frequency of a synset is based on the frequency of the word plus the frequencies of all its descendant synsets in a large corpus. Therefore, the top-level synsets have the highest frequencies and thus the highest estimated frequency of occurrence. For each pair of nouns from the text to be disambiguated, the *most-informative-subsumer* is determined by finding the common ancestor with the highest information content, where information content is inversely related to frequency. Then each noun is disambiguated by selecting the synset that receives the most support (i.e., information content) from the all of the most-informative-

subsumers.

Eizirik et al. (1993) also describe a Bayesian network model for word-sense disambiguation, which includes syntactic disambiguation as well as lexical information. However, their networks are not automatically constructed.

9 Conclusion

This paper explores various representational issues that must be addressed when formulating a Bayesian network representation of lexical information such as is expressed in WordNet. We describe two choices for the representation of lexical items and two choices for the representation lexical relations. The effects on evidence propagation in the network is also discussed.

References

- Bruce, R., and Wiebe, J. (1994), "Word-sense disambiguation using decomposable models", in *Proc. of the 32nd Annual Meeting of the Association for Computational Linguistics (ACL-94)*, pp 139-146.
- Charniak, E., and R. Goldman (1993), "A Bayesian Model of Plan Recognition", *Artificial Intelligence* 64:53-79.
- Cowie, J., J. Guthrie, and L. Guthrie (1992), "Lexical Disambiguation Using Simulated Annealing", *Proc. COLING-92*, pp. 359-365.
- Eizirik, L., V. Barbosa, and S. Mendes (1993), "A Bayesian-Network Approach to Lexical Disambiguation", *Cognitive Science*, 17:257-283.
- Heckerman, D. and J. Breese (1994), "Causal Independence for Probability Assessment and Inference Using Bayesian Networks", Technical Report MSR-TR-94-08, Microsoft Research, (Revised October, 1995).
- Hirst, G. (1988), "Resolving Lexical Ambiguity Computationally with Spreading Activation and Polaroid Words", in *Lexical Ambiguity Resolution*, S. Small, G. Cottrell, and M. Tanenhaus (eds), San Mateo, CA: Morgan Kaufmann Publishers, pp. 73-107.
- Lesk, M. (1986), "Automatic Sense Disambiguation Using Machine Readable Dictionaries: How to Tell a Pine Cone from an Ice Cream Cone", in *Proc. SIGDOC*, Toronto.
- Miller, G., (1990), "WordNet: An On-line Lexical Database", *International Journal of Lexicography* 3(4).
- Pearl, J. (1988), *Probabilistic Reasoning in Intelligent Systems*, San Mateo, CA: Morgan Kaufmann.
- Quillian, M. (1968), "Semantic Memory", in *Semantic Information and Processing*, M. Minsky, ed., Cambridge, MA: MIT Press.
- Resnik, P. (1995), "Disambiguating Noun Groupings with Respect to WordNet Senses", in *Proc. Third Workshop on Very Large Corpora*, Cambridge, MA, June 1995.
- Sussna, M. (1993), "Word Sense Disambiguation for Free-text Indexing Using a Massive Semantic Network", in *Proc. Second International Conference on Information and Knowledge Management (CIKM-93)*, Arlington, Virginia.
- Veronis, J., and N. Ide (1990), "Word Sense Disambiguation with Very Large Neural Networks Extracted from Machine Readable Dictionaries", in *Proc. COLING-90*, Helsinki, August 1990.
- Voorhees, E. (1993) "Using WordNet to Disambiguate Word Senses for Text Retrieval", in *Proc. 16th Annual ACM SIGIR Conference on Research and Development in Information Retrieval*, Pittsburgh, pp. 171-180.
- Yarowsky, D. (1992), "Word-Sense Disambiguation Using Statistical Models of Roget's Categories Trained on Large Corpora", in *Proc. COLING-92*, Nantes, Aug 23-28, pp. 454-460.