

# Towards the Acquisition and Representation of a Broad-Coverage Lexicon

Rebecca Bruce and Janyce Wiebe

Computing Research Laboratory  
and

Department of Computer Science  
New Mexico State University  
Las Cruces, NM 88003

rbruce@crl.nmsu.edu, wiebe@cs.nmsu.edu

## Abstract

Statistical techniques for NLP typically do not take advantage of existing domain knowledge and require large amounts of tagged training data. This paper presents a partial remedy to these shortcomings by introducing a richer class of statistical models, *graphical models*, along with techniques for: (1) establishing the form of the model in this class that best describes a given set of training data, (2) estimating the parameters of graphical models from untagged data, (3) combining constraints formulated in propositional logic with those derived from training data to produce a graphical model, and (4) simultaneously resolving interdependent ambiguities. The paper also describes how these tools can be used to produce a broad-coverage lexicon represented as a probabilistic model, and presents a method for using such a lexicon to simultaneously disambiguate all words in a sentence.

## Introduction

The specification and acquisition of lexical knowledge is a central problem in natural language processing. Statistical techniques have been used to partially automate these tasks, but they typically do not take advantage of existing domain knowledge and require large amounts of tagged training data. These shortcomings have resulted in statistical models that are too impoverished to capture the kind of information that is typically thought to be necessary to resolve word meaning ambiguity. We present a partial remedy to both of these shortcomings by introducing a richer class of statistical models, *graphical models*, along with techniques for: (1) establishing the form of the graphical model that best describes a given set of training data, (2) estimating the parameters of graphical models from untagged data, (3) combining constraints formulated in propositional logic with those derived from training data to produce a graphical model, and (4) simultaneously resolving interdependent ambiguities. We describe how these tools can be used to produce a broad-coverage lexicon represented as a probabilistic model. In this model, information gathered empirically from training data is integrated with knowledge present in

domain theory by expressing both in terms of conditional independence. We also describe how such a lexicon can be used to simultaneously disambiguate all words in a sentence.

While the class of models and the methods for formulating such models (specifically, techniques (1), (2), and (4)) are well known in the field of statistics, this particular application is largely unimplemented, and the lexicon that we describe has not yet been produced. The contribution of this paper is the presentation of a theoretically feasible approach to automatically acquiring and using a broad-coverage lexicon that integrates both empirically and analytically acquired information in a single formalism.

Throughout this presentation, it is assumed that there exists a predefined and finite set of sense distinctions for each word to be included in the lexicon and that the desired level of syntactic analysis can be performed as a preprocessing step. Given these restrictions, the proposed representation is not specific to a particular set of sense distinctions or syntactic theory as long as the two are compatible both with each other and any domain theory to be incorporated in the lexicon.

We are not presently addressing either sense productivity or syntactic analysis. However, we are proposing a representation capable of expressing diverse types of information within a single formalism, and it might be possible to incorporate the information necessary to address such tasks within this representation. In fact, there already exists work supporting such a goal; a mapping of context free grammars to statistical models similar in form to the models described here has been proposed by Eizirik et al. (1993), and mappings between theories in first-order predicate logic (such as might be required to express sense productivity) and the same class of models used by Eizirik et al. have been proposed by Charniak and Goldman (1993), Poole (1993), and Breese (1992).

The remainder of this paper is organized as follows. Graphical models are first described along with tools for formulating them. A specific approach to producing a broad-coverage lexicon is then presented, which

makes use of the tools described in the previous section. The paper concludes with a description of how semantic lexical disambiguation can be performed using such a lexicon.

## The Tools

### Graphical Models

Recall that the lexicon will be represented by a probabilistic model. We make use of a more general class of models than previously used in most NLP applications. Models in this class are capable of characterizing a rich set of relationships between a large number of variables, where these variables may be spatially separated and may also be of different types (e.g., class-based vs. specific collocations). The models we refer to are called "graphical models" (Whittaker 1990).

The probability distribution of a graphical model corresponds to a Markov field. This is because a graphical model is a probability model for multivariate random observations whose dependency structure has a valid graphical representation, a *dependency graph*. The dependency graph of a model is formed by mapping each variable in the model to a node in the graph, and then drawing undirected edges between the nodes corresponding to variables that are interdependent. A valid dependency graph has the property that all variables that are not directly connected in the graph are conditionally independent given the values of the variables mapping to nodes connecting them. For example, if node *A* separates node *B* from node *C* in a dependency graph, then the variable mapping to node *B* is conditionally independent of the variable mapping to node *C*, given the value of the variable mapping to node *A*. Not only does a dependency graph provide a clear interpretation of the interactions among variables, it also has been shown (Pearl 1988, Lauritzen & Spiegelhalter 1988) to form the basis of a distributed computational architecture for probabilistic constraint propagation.

### Statistical Inference Via Stochastic Simulation

In this section we describe three tools for formulating and using a graphical model. Each of these tools makes use of a stochastic simulation technique to perform statistical inference, where statistical inference is the inference of population characteristics from a sample of training data. In all cases the inference techniques have been chosen to meet the demands of large-scale applications by minimizing the amount of hand-tagged data required. The first tool is a method for selecting the form of a graphical model, given a sample of training data. The second tool is a method for estimating the parameters of a graphical model from untagged data once the form of the model is specified. The final tool is a technique for resolving multiple, mutually-constraining ambiguities in graphical models.

**Model Selection** Part of the lexicon will be derived statistically from training data. To do this, an appropriate probabilistic model must be formulated. In previous work (Bruce and Wiebe 1994ab), we used a method for model formulation that selects informative contextual features and identifies the *decomposable model* (a kind of graphical model) that is the best approximation to the joint distribution of word senses and these features. Here we describe improvements to our method that make large-scale applications feasible.

The objective of statistical modeling is to find the simplest model that fits the training data. Rather than making assumptions about how the variables are related, we explore this empirically. In our previous work (Bruce & Wiebe 1994ab), both feature selection and model formulation are accomplished via a process of hypothesis testing: patterns of interdependencies among variables are expressed by decomposable models, and then these models are evaluated as to how well they fit the training data. A model is said to fit the training data if the distribution it defines differs from the training data by an amount consistent with sampling variation. The difference between the distribution in the training data and that defined by a model can be measured by the likelihood ratio statistic  $G^2$ . In order to know if a model fits the training data, one must know the probability of randomly selecting a data sample, having a given  $G^2$  value, from the population defined by the model. The significance of the fit of a model is then equal to the probability of randomly selecting a data sample with a  $G^2$  value that is as big or bigger than that of the training data. In our previous work, large-sample approximations for the distribution of  $G^2$  were used to determine the significance of a model; this meant that very large data samples were required.

Our new method for determining the significance of a model eliminates the use of large-sample approximations of the distribution of  $G^2$ . This is accomplished by stochastically simulating the exact conditional distribution of  $G^2$  in order to determine the significance of the  $G^2$  value describing the fit of a model. The exact conditional distribution of  $G^2$  for each model is the distribution of  $G^2$  values that is observed for comparable data samples randomly generated from the model being tested. A comparable data sample is one having the same sufficient statistics as the actual training data, where the sufficient statistics are the marginal distributions of the variables that are specified as interdependent in the model. Algorithms exist for generating comparable data samples from models describing conditional independence between two variables (Agresti 1992, Verbeek and Kroonenberg 1985). Using these algorithms, the significance of an arbitrary graphical model can be evaluated by determining the significance of each conditional independence relationship in that model. The advantage of simulating the exact conditional distribution of  $G^2$  to assess the statistical sig-

nificance of a model is that it eliminates (1) the need for parameter estimation in selecting the form of the model, and (2) the need to use large-sample approximations of the distribution of  $G^2$ . It can, therefore, be used to establish the form of a model with only a small amount of sense-tagged data.

**Estimating Model Parameters from Untagged Data** Once the form of the model has been established as described above, it should be possible to obtain reliable estimates of the model parameters from untagged text. The idea is to treat the missing tags as "missing data" and draw on the wealth of statistical techniques for augmenting incomplete data given some knowledge of the form of the complete data density. The scheme we will describe is a stochastic simulation technique referred to as "The Gibbs sampler" (Geman & Geman 1984); it should be noted that there are many possible variations on this general algorithm (Rubin 1991).

The procedure is a Bayesian approach to data augmentation that is very similar in spirit to the EM algorithm (Dempster et al. 1977). The basic idea is straightforward. The observed data,  $y$ , is augmented by  $z$ , the missing data (i.e., the missing tags), and the problem is set up so that: (1) if the missing data,  $z$  were known, then the posterior density  $P(\theta|y, z)$ , could be easily calculated (where  $\theta$  is the vector of model parameters), and (2) if the model parameters,  $\theta$  were known, then the posterior density  $P(z|\theta, y)$ , could be easily calculated. An initial value for  $z$ ,  $z_1$ , is assumed and the process iterates between drawing  $\theta_i$  from  $P(\theta|y, z_i)$  and drawing  $z_{i+1}$  from  $P(z|y, \theta_i)$ . The process continues to iterate until the average value for  $\theta$ , the estimate of the model parameters, converges.

In the Gibbs sampler, the simulation of  $P(z|y, \theta)$  is further simplified. This is done by partitioning the missing data as  $z = (z_1, z_2, \dots, z_n)$ , in order to make the drawing of each part of  $z$  (i.e.,  $z_j$ ) easier. Specifically, what is done for  $z$  is to sample  $z_1$  from  $P(z_1|z_2, \dots, z_n, \theta, y)$ ,  $z_2$  from  $P(z_2|z_1, z_3, \dots, z_n, \theta, y)$ , and so on up to  $z_n$ . The Gibbs sampler uses some starting value for  $z$  and then generates a value for  $\theta$ , followed by new values for  $z_1$  through  $z_n$ . Each time a new value for a partition of  $z$ , say  $z_j$ , is generated, the updated values of all other partitions of  $z$  are used in selecting that value. As described above, the process iterates between generating  $\theta$  and generating  $z_1$  through  $z_n$  until the average value of  $\theta$  converges.

The process of stochastic simulation is potentially slow. Hundreds of cycles may be required to achieve reasonable accuracy, with each cycle requiring time on the order of the sum of the number of variables and the number of interdependencies between variables (Pearl 1988). Fortunately, the algorithm can be easily executed in parallel by concurrent processors. A further reduction in computation time can be realized by incorporating into the algorithm the simulated annealing heuristic of Kirkpatrick, Gelatt and Vecchi (1983) as

demonstrated by Geman and Geman (1984).

**Resolution of Interdependent Ambiguities** In this section, we present a method for disambiguation of multiple words when the senses of those words are interdependent. Stated another way, we are describing a method for finding the set of word-sense classifications that, when taken as a whole, most completely satisfy the constraints among variables. Methods for doing this have been developed for various classes of models. The Viterbi algorithm (Forney 1973) can be used to find the classification sequence that maximizes the probability of a Markov chain. In this work, we are not restricted to Markov chain models and therefore use techniques that are applicable to the broader class of graphical models. Specifically, we use the Gibbs sampler, presented earlier, to approximate the desired distribution. The difference is that now we know the values of the model parameters and are only interested in estimating the values of the missing data. In this case, the missing data are the sense tags of all words in a sentence and the values of the variables that specify relationships among word senses. These relationships among word senses can be formulated from propositional logic, as described in the next section.

### Combining Empirical and Analytical Constraints

There are many domain theories describing relationships among word senses, and many of these theories can be expressed in propositional logic. In this section, we describe methods for (1) representing, as conditional independence relationships, constraints that are formulated in propositional logic, (2) assigning probabilities to these relationships, and (3) combining such relationships with those derived from training data, to produce a graphical model. The methods outlined are specific to the type of theories used in this application. These are existing concept taxonomies, such as WordNet (Miller 1990), that specify interdependencies among the senses of the ambiguous words in a sentence.

Formulating a probabilistic model involves three major sub-tasks: defining the random variables in the model, identifying the dependencies among those variables (i.e., specifying the form of the model), and quantifying those dependency relationships (i.e., estimating the parameters of the model). In keeping with previous approaches (Bacchus 1990), propositions are, in general, mapped to discrete random variables. An exception to this mapping is made for propositions corresponding to word senses. Each ambiguous word will map to a single random variable, with values corresponding to the possible senses of that word<sup>1</sup>. This treatment is necessary to capture the implicit relationships that exist between a word and its senses, i.e.,  $word \rightarrow wordSense_1 \vee wordSense_2 \vee \dots \vee wordSense_n$

<sup>1</sup>The possible values for each word could also include the null value corresponding to the absence of that word.

as well as the mutual exclusion relationship that exists among the senses of a word, i.e.,  $wordSense_1 \rightarrow \neg wordSense_2 \wedge \neg wordSense_3 \wedge \dots \wedge \neg wordSense_n$ .

In logic, connectives, such as conjunction and implication, are used to describe relationships among propositions; in a probabilistic model, these relationships are expressed as statements of conditional independence among the variables. We require that the axioms for each proposition be covering, and that the set of axioms, as a whole, be acyclic. Given that these properties hold, all propositions (random variables) in a single axiom can be treated as interdependent, while those propositions not together in any single axiom can be treated as not interdependent. For example, an axiom of the form  $biology \rightarrow science$ , where *biology* and *science* are subject designations in a hierarchy of subject classifications, would be viewed as specifying a dependency between the variables corresponding to *biology* and *science*. If there is also an axiom of the form  $geology \rightarrow science$ , with *geology* referring to a subject designation, then there would also exist a dependency between the variables corresponding to *geology* and *science*. But based on this information alone, no direct dependency between *geology* and *biology* would be established.

Once the dependencies among variables have been identified, they must be quantified. The theory of Markov fields (Pearl 1988, Isham 1981, Darroch et al. 1980) provides a method for assigning complete and consistent probabilities given the dependency structure of the model and numerical measures of the compatibility of the values of interdependent variables, called "compatibility measures." As described in Pearl (1988), the method is referred to as "Gibbs potential." The advantages of this method are that the compatibility measures need only quantify the local interactions within the sets of completely interdependent variables, and that the numerical values assigned to these local interactions need not be globally consistent. Compatibility measures can be assigned in accordance with the semantics of propositional logic. Returning to the previous example of  $biology \rightarrow science$ , a high compatibility measure would be assigned to the co-occurrence of  $biology = true$  and  $science = true$ ,  $biology = false$  and  $science = false$ , as well as  $biology = false$  and  $science = true$ , while a low compatibility measure would be assigned to the pair  $biology = true$  and  $science = false$ .

The Gibbs potential can also be used to merge the model derived from logical constraints with those developed from the training data. In merging multiple models, the compatibility measures are formulated from the parameters of the models to be merged.

## Construction and Application of a Broad-Coverage Lexicon

We now describe the development of a broad-coverage lexicon that integrates information gathered empiri-

cally from training data with analytically derived, domain knowledge. This lexicon can be used to disambiguate a large vocabulary of words with respect to the fine-grained sense distinctions made in a standard dictionary. The statistical techniques described in the first part of this paper can be used to automatically formulate a probabilistic model describing the relationship between each ambiguous word and a select set of unambiguous contextual features, without requiring a large amount of disambiguated corpora. The individual models, formulated from training data, can then be interconnected by a constraint structure specifying the semantic relationships that exist among the senses of words, as described earlier. Fortunately, there are existing concept taxonomies that describe the semantic relationships among word senses. Examples of such theories are: the "subject code" and "box code" hierarchies in the Longman's Dictionary of Contemporary English (Procter et al. 1978), the hyponymy and meronymy taxonomies in WordNet (Miller 1990), and the various taxonomies derived from machine readable dictionaries (Knight 1993, Bruce et al. 1992, Vossen 1990, Chodorow et al. 1985, Michiels & Noel 1982, Amsler 1980). In the remainder of this section, we present a specific approach to producing a broad-coverage lexicon, along with a description of how such a lexicon can be used to disambiguate all targeted content words in a single sentence simultaneously using the Gibbs sampler described in the previous section.

## Proposed Implementation

The first phase of construction is to formulate, from training data, a decomposable model for each word to be disambiguated. The candidate contextual features could include the following: the morphology of the ambiguous word, the part-of-speech (POS) categories of the immediately surrounding words, specific collocations, sentence position, and aspects of the phrase structure of the sentence.

Once a set of contextual features has been chosen, the form of the model for each word is selected from the class of decomposable models, and estimates of the model parameters are then obtained from untagged data. As mentioned earlier, selection of the form of a model requires a small amount of tagged data indicating the senses of the ambiguous word. It may be possible to further reduce the total amount of sense-tagged data needed for this phase of construction by identifying a parametric model or models applicable to a wide range of content words, as described in our previous work (Bruce & Wiebe 1994a).

In the second phase of construction, the statistical models for the individual words are interconnected by a constraint structure specifying the semantic relationships that exist among the senses of words. Several concept taxonomies appear to be applicable, and it should be possible to include multiple theories in a single model, provided that they are based on the same

sense distinctions.

Figure 1 depicts a small fraction of a dependency graph of a probabilistic model that combines hypothetical but representative propositional axioms with models similar in form to those that we developed from training data in previous work. As can be seen in figure 1, the constraints derived from the training data are connected to the global constraint structure, which in this case is derived from a subject classification hierarchy, only at the nodes corresponding to ambiguous words. The sparsity of the graph, as a whole, is a result of the large number of conditional independence relationships manifest in the model. Such a sparse dependency graph is computationally advantageous in that it reduces the number of parameters that need to be estimated and reduces the complexity of the stochastic simulation of the model.

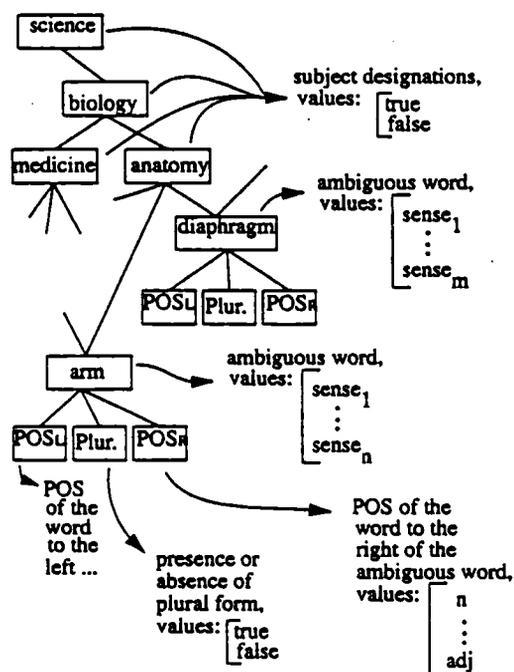


Figure 1

### Resolution of Ambiguity Through Simultaneous Satisfaction of all Constraints

Using the model constructed as described above, the disambiguation of all ambiguous words in a sentence will be accomplished via a stochastic simulation of that model using the Gibbs sampler. The tag sequence identified for each sentence via the stochastic simulation will be the maximum a posteriori probability (MAP) solution to the satisfaction of all probabilistic dependencies among variables.

In performing a stochastic simulation, one is in effect randomly sampling from the probability distribution

defined by the model and the known variable values. During word-sense disambiguation, the values of the contextual features, which are the leaves of the dependency graph, are observable and are therefore known; these values are fixed during the stochastic simulation. As discussed earlier, the value of each unobservable variable in the model (i.e., variables whose values are not known) is sequentially updated using the values of the model parameters and the current values of the interdependent variables. This process of sequential updating (or sampling) continues until convergence or near convergence in the average of the sampled values is achieved.

During the sampling process, the value of a variable changes in response to changes in the current values of all dependent variables (i.e., the directly connected variables in the dependency graph). As depicted in figure 1, the sense tag of an ambiguous word is dependent on the values of the contextual features as well as the values of the variables in the hierarchy structure that connects all ambiguous words. While the values of the contextual features for each input word are observable and therefore known, the values of the variables in the hierarchy are not observable. Their values change in response to changes in the probabilities of the various word senses and the values of other hierarchical variables. If the senses of two or more words in a sentence are directly connected to (dependent on) a common hierarchical variable, then the probability of that variable being "true" increases, which in turn increases the probability of each of the connected word senses. The more indirect the connection (dependency) between the common hierarchical variable and the word sense, the weaker the support. The process is much like spreading activation, but the relationships between variables are expressed as probabilities and the process returns the MAP solution to the constraint satisfaction problem.

There is one restriction on the application of the Gibbs sampler to models formulated from logical constraints. In order for the proof of convergence to hold, the joint distribution of all variables must be strictly positive<sup>2</sup>. This means that logical constraints may not be assigned absolute certainty (i.e., given a probability of 1.0), so that alternative values do not have a zero probability. And because the rate of convergence is effected by small probabilities (Chin & Cooper 1987), the probabilities assigned to logical constraints may have to be adjusted to speed convergence.

The last issue to be addressed in this section is the size of the model and the corresponding complexity of the stochastic simulation. As described so far, the model would contain nodes corresponding to one or more concept taxonomies, in addition to nodes corresponding to a large number of ambiguous words and

<sup>2</sup>Technically, the requirement is that the Gibbs sampler be irreducible, but the simplest way of assuring irreducibility is to assure that the joint distribution is strictly positive.

their contextual features. But, as stated earlier, the time required for the stochastic simulation is proportional to the number of nodes plus the number of edges in the network. In order to assure that the time requirements for the stochastic simulation are not extreme, the model can be limited to only the nodes needed to disambiguate the current input sentence. This can be done by dynamically creating a network corresponding to each input sentence using "marker-passing," as described in Charniak and Goldman (1993), to select the portions of the complete model needed to process each specific input sentence.

## References

- Agresti, A. 1992. A Survey of Exact Inference for Contingency Tables. *Statistical Science* 7(1): 131-177.
- Amsler, R. 1980. The Structure of the Merriam-Webster Pocket Dictionary. *Technical Report TR-164*, University of Texas at Austin.
- Bacchus, F. 1990. *Representing and Reasoning with Uncertain Knowledge*. Cambridge, MA: MIT Press.
- Breese, J. 1992. Construction of Belief and Decision Networks. *Computational Intelligence*, Vol. 8, No. 4, pp 624-644.
- Bruce, R. and Wiebe, J. 1994a. A New Approach to Word-Sense Disambiguation. *Proc. ARPA Human Language Technology Workshop*. Plainsboro, NJ.
- Bruce, R. and Wiebe, J. 1994b. Word-Sense Disambiguation Using Decomposable Models. *Proceedings of the 32nd Annual Meeting of the ACL*. Las Cruces, NM.
- Bruce, R.; Wilks, Y.; Guthrie, L.; Slator, B.; and Dunning, T. 1992. NounSense - A Disambiguated Noun Taxonomy with a Sense of Humour. *Memorandum in Computer and Cognitive Science*, MCCS-92-246, Computing Research Laboratory, New Mexico State University.
- Charniak, E. and Goldman, R. 1993. A Bayesian model of plan recognition. *Artificial Intelligence* (64): 53-79.
- Chin, H. and Cooper, G. 1987. Stochastic Simulation of Bayesian Belief Networks. *Procs. 3rd Workshop on Uncertainty in AI*, Seattle, WA. 106-113.
- Chodorow, M., Byrd, R., and Heidorn, G. 1985. Extracting Semantic Hierarchies from a Large On-Line Dictionary. *Proceedings of the 23rd Annual Meeting of the ACL*, Chicago, IL, 299-304.
- Darroch, J.; Lauritzen, S.; and Speed, T. 1980. Markov Fields and Log-Linear Interaction Models for Contingency Tables. *Annals of Statistics* 8(3):522-539.
- Dempster, A.; Laird, N.; and Rubin, D. 1977. Maximum Likelihood from Incomplete Data Via the EM Algorithm. *Journal of the Royal Statistical Society B* (39): 1-38.
- Eizirik, L.; Barbosa, V.; and Mendes, S. 1993. A Bayesian-Network Approach to Lexical Disambiguation. *Cognitive Science* (17): 257-283.
- Forney, G. D. 1973. The Viterbi Algorithm. *Proc. IEEE* (6): 268-278.
- Geman, S. and Geman, D. 1984. Stochastic Relaxation, Gibbs Distributions and the Bayesian Restoration of Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (6): 721-741.
- Isham, V. 1981. An Introduction to Spatial Processes and Markov Random Fields. *International Statistical Review* (49): 21-43.
- Kirkpatrick, S.; Gelatt, C.; and Vecchi, M. 1983. Optimization by Simulated Annealing. *Science* (220) 671-680.
- Knight, K. 1993. Building a Large Ontology for Machine Translation. *Proc. ARPA Human Language Technology Workshop*. Plainsboro, NJ.
- Lauritzen, S. and Spiegelhalter, D. 1988. Local Computations with Probabilities on Graphical Structures and their Application to Expert Systems. *J. R. Statist. Soc. B* 50(2): 157-224.
- Michiels, A., and Noel, J. 1982. Approaches to Thesaurus Production. *Proceedings of the 9th International Conference on Computational Linguistics (COLING-82)*, Prague, Czechoslovakia, 227-232.
- Miller, G. 1990. Special Issue, WordNet: An on-line lexical database. *International Journal of Lexicography* 3(4).
- Pearl, J. 1988. *Probabilistic Reasoning In Intelligent Systems: Networks of Plausible Inference*. San Mateo, Ca.: Morgan Kaufmann.
- Procter, Paul et al. 1978. *Longman Dictionary of Contemporary English*.
- Poole, D. 1993. Probabilistic Horn abduction and Bayesian Networks. *Artificial Intelligence* (64): 81-129.
- Rubin, D. 1991. EM and Beyond. *Psychometrika* 56(2): 241-254.
- Verbeek, A. and Kroonenberg, P. 1985. A survey of algorithms for exact distributions of the test statistics in  $r \times c$  contingency tables with fixed marginals. *Computational Statistics and Data Analysis* (3): 159-185.
- Vossen, P. 1990. The End of the Chain: Where Does Decomposition of Lexical Knowledge Lead Us Eventually? *Proceedings of the 4th Conference on Functional Grammar*, Copenhagen.
- Whittaker, J. 1990. *Graphical Models In Applied Multivariate Statistics*. New York, NY: John Wiley & Sons.