

# Class-Based Probability Estimation Using a Semantic Hierarchy

Stephen Clark\*  
University of Edinburgh

David Weir†  
University of Sussex

*This article concerns the estimation of a particular kind of probability, namely, the probability of a noun sense appearing as a particular argument of a predicate. In order to overcome the accompanying sparse-data problem, the proposal here is to define the probabilities in terms of senses from a semantic hierarchy and exploit the fact that the senses can be grouped into classes consisting of semantically similar senses. There is a particular focus on the problem of how to determine a suitable class for a given sense, or, alternatively, how to determine a suitable level of generalization in the hierarchy. A procedure is developed that uses a chi-square test to determine a suitable level of generalization. In order to test the performance of the estimation method, a pseudo-disambiguation task is used, together with two alternative estimation methods. Each method uses a different generalization procedure; the first alternative uses the minimum description length principle, and the second uses Resnik's measure of selectional preference. In addition, the performance of our method is investigated using both the standard Pearson chi-square statistic and the log-likelihood chi-square statistic.*

## 1. Introduction

This article concerns the problem of how to estimate the probabilities of noun senses appearing as particular arguments of predicates. Such probabilities can be useful for a variety of natural language processing (NLP) tasks, such as structural disambiguation and statistical parsing, word sense disambiguation, anaphora resolution, and language modeling. To see how such knowledge can be used to resolve structural ambiguities, consider the following prepositional phrase attachment ambiguity:

### Example 1

Fred ate strawberries with a spoon.

The ambiguity arises because the prepositional phrase *with a spoon* can attach to either *strawberries* or *ate*. The ambiguity can be resolved by noting that the correct sense of *spoon* is more likely to be an argument of “*ate-with*” than “*strawberries-with*” (Li and Abe 1998; Clark and Weir 2000).

The problem with estimating a probability model defined over a large vocabulary of predicates and noun senses is that this involves a huge number of parameters, which results in a sparse-data problem. In order to reduce the number of parameters, we propose to define a probability model over senses in a semantic hierarchy and

---

\* Division of Informatics, University of Edinburgh, 2 Buccleuch Place, Edinburgh, EH8 9LW, UK. E-mail: stephenc@cogsci.ed.ac.uk.

† School of Cognitive and Computing Sciences, University of Sussex, Brighton, BN1 9QH, UK. E-mail: david.weir@cogs.susx.ac.uk.

to exploit the fact that senses can be grouped into classes consisting of semantically similar senses. The assumption underlying this approach is that the probability of a particular noun sense can be approximated by a probability based on a suitably chosen class. For example, it seems reasonable to suppose that the probability of (the food sense of) *chicken* appearing as an object of the verb *eat* can be approximated in some way by a probability based on a class such as FOOD.

There are two elements involved in the problem of using a class to estimate the probability of a noun sense. First, given a suitably chosen class, how can that class be used to estimate the probability of the sense? And second, given a particular noun sense, how can a suitable class be determined? This article offers novel solutions to both problems, and there is a particular focus on the second question, which can be thought of as how to find a suitable level of generalization in the hierarchy.<sup>1</sup>

The semantic hierarchy used here is the noun hierarchy of WordNet (Fellbaum 1998), version 1.6. Previous work has considered how to estimate probabilities using classes from WordNet in the context of acquiring selectional preferences (Resnik 1998; Ribas 1995; Li and Abe 1998; McCarthy 2000), and this previous work has also addressed the question of how to determine a suitable level of generalization in the hierarchy. Li and Abe use the minimum description length principle to obtain a level of generalization, and Resnik uses a simple technique based on a statistical measure of selectional preference. (The work by Ribas builds on that by Resnik, and the work by McCarthy builds on that by Li and Abe.) We compare our estimation method with those of Resnik and Li and Abe, using a pseudo-disambiguation task. Our method outperforms these alternatives on the pseudo-disambiguation task, and an analysis of the results shows that the generalization methods of Resnik and Li and Abe appear to be overgeneralizing, at least for this task.

Note that the problem being addressed here is the engineering problem of estimating predicate argument probabilities, with the aim of producing estimates that will be useful for NLP applications. In particular, we are not addressing the problem of acquiring selectional restrictions in the way this is usually construed (Resnik 1993; Ribas 1995; McCarthy 1997; Li and Abe 1998; Wagner 2000). The purpose of using a semantic hierarchy for generalization is to overcome the sparse data problem, rather than find a level of abstraction that best represents the selectional restrictions of some predicate. This point is considered further in Section 5.

The next section describes the noun hierarchy from WordNet and gives a more precise description of the probabilities to be estimated. Section 3 shows how a class from WordNet can be used to estimate the probability of a noun sense. Section 4 shows how a chi-square test is used as part of the generalization procedure, and Section 5 describes the generalization procedure. Section 6 describes the alternative class-based estimation methods used in the pseudo-disambiguation experiments, and Section 7 presents those experiments.

## 2. The Semantic Hierarchy

The noun hierarchy of WordNet consists of senses, or what Miller (1998) calls *lexicalized concepts*, organized according to the “is-a-kind-of” relation. Note that we are using *concept* to refer to a lexicalized concept or sense and not to a set of senses; we use *class* to refer to a set of senses. There are around 66,000 different concepts in the noun hierarchy

---

<sup>1</sup> A third element of the problem, namely, how to obtain arguments of predicates as training data, is not considered here. We assume the existence of such data, obtained from a treebank or shallow parser.

of WordNet version 1.6. A concept in WordNet is represented by a “synset,” which is the set of synonymous words that can be used to denote that concept. For example, the synset for the concept  $\langle \text{cocaine} \rangle^2$  is  $\{ \text{cocaine}, \text{cocain}, \text{coke}, \text{snow}, \text{C} \}$ . Let  $\text{syn}(c)$  be the synset for concept  $c$ , and let  $\text{cn}(n) = \{ c \mid n \in \text{syn}(c) \}$  be the set of concepts that can be denoted by noun  $n$ .

The hierarchy has the structure of a directed acyclic graph (although only around 1% of the nodes have more than one parent), where the edges of the graph constitute what we call the “direct-isa” relation. Let *isa* be the transitive, reflexive closure of *direct-isa*; then  $c' \text{ isa } c$  implies  $c'$  is a kind of  $c$ . If  $c' \text{ isa } c$ , then  $c$  is a *hypernym* of  $c'$  and  $c'$  is a *hyponym* of  $c$ . In fact, the hierarchy is not a single hierarchy but instead consists of nine separate subhierarchies, each headed by the most general kind of concept, such as  $\langle \text{entity} \rangle$ ,  $\langle \text{abstraction} \rangle$ ,  $\langle \text{event} \rangle$ , and  $\langle \text{psychological\_feature} \rangle$ . For the purposes of this work we add a common root dominating the nine subhierarchies, which we denote  $\langle \text{root} \rangle$ .

There are some important points that need to be clarified regarding the hierarchy. First, every concept in the hierarchy has a nonempty synset (except the notional concept  $\langle \text{root} \rangle$ ). Even the most general concepts, such as  $\langle \text{entity} \rangle$ , can be denoted by some noun; the synset for  $\langle \text{entity} \rangle$  is  $\{ \text{entity}, \text{something} \}$ . Second, there is an important distinction between an individual concept and a set of concepts. For example, the individual concept  $\langle \text{entity} \rangle$  should not be confused with the set or class consisting of concepts denoting kinds of entities. To make this distinction clear, we use  $\bar{c} = \{ c' \mid c' \text{ isa } c \}$  to denote the set of concepts dominated by concept  $c$ , including  $c$  itself. For example,  $\overline{\langle \text{animal} \rangle}$  is the set consisting of those concepts corresponding to kinds of animals (including  $\langle \text{animal} \rangle$  itself).

The probability of a concept appearing as an argument of a predicate is written  $p(c \mid v, r)$ , where  $c$  is a concept in WordNet,  $v$  is a predicate, and  $r$  is an argument position.<sup>3</sup> The focus in this article is on the arguments of verbs, but the techniques discussed can be applied to any predicate that takes nominal arguments, such as adjectives. The probability  $p(c \mid v, r)$  is to be interpreted as follows: This is the probability that some noun  $n$  in  $\text{syn}(c)$ , when denoting concept  $c$ , appears in position  $r$  of verb  $v$  (given  $v$  and  $r$ ). The example used throughout the article is  $p(\langle \text{dog} \rangle \mid \text{run}, \text{subj})$ , which is the conditional probability that some noun in the synset of  $\langle \text{dog} \rangle$ , when denoting the concept  $\langle \text{dog} \rangle$ , appears in the subject position of the verb *run*. Note that, in practice, no distinction is made between the different senses of a verb (although the techniques do allow such a distinction) and that each use of a noun is assumed to correspond to exactly one concept.<sup>4</sup>

### 3. Class-Based Probability Estimation

This section explains how a set of concepts, or class, from WordNet can be used to estimate the probability of an individual concept. More specifically, we explain how a set of concepts  $\bar{c}'$ , where  $c'$  is some hypernym of concept  $c$ , can be used to estimate  $p(c \mid v, r)$ . (Recall that  $\bar{c}'$  denotes the set of concepts dominated by  $c'$ , including  $c'$  itself.) One possible approach would be simply to substitute  $\bar{c}'$  for the individual concept  $c$ . This is a poor solution, however, since  $p(\bar{c}' \mid v, r)$  is the conditional probability that

<sup>2</sup> Angled brackets are used to denote concepts in the hierarchy.

<sup>3</sup> The term *predicate* is used loosely here, in that the predicate does not have to be a semantic object but can simply be a word form.

<sup>4</sup> A recent paper that extends the acquisition of selectional preferences to sense-sense relationships is Agirre and Martinez (2001).

some noun denoting a concept in  $\bar{c}'$  appears in position  $r$  of verb  $v$ . For example,  $p(\langle \text{animal} \rangle | \text{run, subj})$  is the probability that some noun denoting a kind of animal appears in the subject position of the verb *run*. Probabilities of sets of concepts are obtained by summing over the concepts in the set:

$$p(\bar{c}' | v, r) = \sum_{c'' \in \bar{c}'} p(c'' | v, r) \quad (1)$$

This means that  $p(\langle \text{animal} \rangle | \text{run, subj})$  is likely to be much greater than  $p(\langle \text{dog} \rangle | \text{run, subj})$  and thus is not a good approximation of  $p(\langle \text{dog} \rangle | \text{run, subj})$ .

What can be done, though, is to *condition* on sets of concepts. If it can be shown that  $p(v | \bar{c}', r)$ , for some hypernym  $c'$  of  $c$ , is a reasonable approximation of  $p(v | c, r)$ , then we have a way of estimating  $p(c | v, r)$ . The probability  $p(v | c, r)$  can be obtained from  $p(c | v, r)$  using Bayes' theorem:

$$p(c | v, r) = p(v | c, r) \frac{p(c | r)}{p(v | r)} \quad (2)$$

Since  $p(c | r)$  and  $p(v | r)$  are conditioned on the argument slot only, we assume these can be estimated satisfactorily using relative frequency estimates. Alternatively, a standard smoothing technique such as Good-Turing could be used.<sup>5</sup> This leaves  $p(v | c, r)$ . Continuing with the  $\langle \text{dog} \rangle$  example, the proposal is to estimate  $p(\text{run} | \langle \text{dog} \rangle, \text{subj})$  using a relative-frequency estimate of  $p(\text{run} | \langle \text{animal} \rangle, \text{subj})$  or an estimate based on a similar, suitably chosen class. Thus, assuming this choice of class,  $p(\langle \text{dog} \rangle | \text{run, subj})$  would be approximated as follows:

$$p(\langle \text{dog} \rangle | \text{run, subj}) \approx p(\text{run} | \langle \text{animal} \rangle, \text{subj}) \frac{p(\langle \text{dog} \rangle | \text{subj})}{p(\text{run} | \text{subj})} \quad (3)$$

The following derivation shows that if  $p(v | \bar{c}'_i, r) = k$  for each child  $c'_i$  of  $c'$ , and  $p(v | c', r) = k$ , then  $p(v | \bar{c}', r)$  is also equal to  $k$ :

$$p(v | \bar{c}', r) = p(\bar{c}' | v, r) \frac{p(v | r)}{p(\bar{c}' | r)} \quad (4)$$

$$= \frac{p(v | r)}{p(\bar{c}' | r)} \left( \sum_i p(\bar{c}'_i | v, r) + p(c' | v, r) \right) \quad (5)$$

$$= \frac{p(v | r)}{p(\bar{c}' | r)} \left( \sum_i p(v | \bar{c}'_i, r) \frac{p(\bar{c}'_i | r)}{p(v | r)} + p(v | c', r) \frac{p(c' | r)}{p(v | r)} \right) \quad (6)$$

$$= \frac{1}{p(\bar{c}' | r)} \left( \sum_i k p(\bar{c}'_i | r) + k p(c' | r) \right) \quad (7)$$

$$= \frac{k}{p(\bar{c}' | r)} \left( \sum_i p(\bar{c}'_i | r) + p(c' | r) \right) \quad (8)$$

$$= k \quad (9)$$

<sup>5</sup> Unsmoothed estimates were used in this work.

Note that the proof applies only to a tree, since the proof assumes that  $\bar{c}'$  is partitioned by  $c'$  and the sets of concepts dominated by each of the daughters of  $c'$ , which is not necessarily true for a directed acyclic graph (DAG). WordNet is a DAG but is a close approximation to a tree, and so we assume this will not be a problem in practice.<sup>6</sup>

The derivation in (4)–(9) shows how probabilities conditioned on sets of concepts can remain constant when moving up the hierarchy, and this suggests a way of finding a suitable set,  $\bar{c}'$ , as a generalization for concept  $c$ : Initially set  $c'$  equal to  $c$  and move up the hierarchy, changing the value of  $c'$ , until there is a significant change in  $p(v | \bar{c}', r)$ . Estimates of  $p(v | \bar{c}'_i, r)$ , for each child  $c'_i$  of  $c'$ , can be compared to see whether  $p(v | \bar{c}', r)$  has significantly changed. (We ignore the probability  $p(v | c', r)$  and consider the probabilities  $p(v | \bar{c}'_i, r)$  only.) Note that this procedure rests on the assumption that  $p(v | \bar{c}, r)$  is close to  $p(v | c, r)$ . (In fact,  $p(v | \bar{c}, r)$  is equal to  $p(v | c, r)$  when  $c$  is a leaf node.) So when finding a suitable level for the estimation of  $p(\langle \text{sandwich} \rangle | \text{eat}, \text{obj})$ , for example, we first assume that  $p(\text{eat} | \langle \text{sandwich} \rangle, \text{obj})$  is a good approximation of  $p(\text{eat} | \langle \text{sandwich} \rangle, \text{obj})$  and then apply the procedure to  $p(\text{eat} | \langle \text{sandwich} \rangle, \text{obj})$ .

A feature of the proposed generalization procedure is that comparing probabilities of the form  $p(v | C, r)$ , where  $C$  is a class, is closely related to comparing ratios of probabilities of the form  $p(C | v, r)/p(C | r)$  (for a given verb and argument position):

$$p(v | C, r) = \frac{p(C | v, r)}{p(C | r)} p(v | r) \quad (10)$$

Note that, for a given verb and argument position,  $p(v | r)$  is constant across classes. Equation (10) is of interest because the ratio  $p(C | v, r)/p(C | r)$  can be interpreted as a measure of association between the verb  $v$  and class  $C$ . This ratio is similar to point-wise mutual information (Church and Hanks 1990) and also forms part of Resnik's association score, which will be introduced in Section 6. Thus the generalization procedure can be thought of as one that finds “homogeneous” areas of the hierarchy, that is, areas consisting of classes that are associated to a similar degree with the verb (Clark and Weir 1999).

Finally, we note that the proposed estimation method does not guarantee that the estimates form a probability distribution over the concepts in the hierarchy, and so a normalization factor is required:

$$p_{sc}(c | v, r) = \frac{\hat{p}(v | [c, v, r], r) \frac{\hat{p}(c|r)}{\hat{p}(v|r)}}{\sum_{c' \in C} \hat{p}(v | [c', v, r], r) \frac{\hat{p}(c'|r)}{\hat{p}(v|r)}} \quad (11)$$

We use  $p_{sc}$  to denote an estimate obtained using our method (since the technique finds sets of semantically similar senses, or “similarity classes”) and  $[c, v, r]$  to denote the class chosen for concept  $c$  in position  $r$  of verb  $v$ ;  $\hat{p}$  denotes a relative frequency estimate, and  $C$  denotes the set of concepts in the hierarchy.

Before providing the details of the generalization procedure, we give the relative-frequency estimates of the relevant probabilities and deal with the problem of am-

<sup>6</sup> Li and Abe (1998) also develop a theoretical framework that applies only to a tree and turn WordNet into a tree by copying each subgraph with multiple parents. One way to extend the experiments in Section 7 would be to investigate whether this transformation has an impact on the results of those experiments.

biguous data. The relative-frequency estimates are as follows:

$$\hat{p}(c | r) = \frac{f(c,r)}{f(r)} = \frac{\sum_{v' \in \mathcal{V}} f(c, v', r)}{\sum_{v' \in \mathcal{V}} \sum_{c' \in \mathcal{C}} f(c', v', r)} \quad (12)$$

$$\hat{p}(v | r) = \frac{f(v,r)}{f(r)} = \frac{\sum_{c' \in \mathcal{C}} f(c', v, r)}{\sum_{v' \in \mathcal{V}} \sum_{c' \in \mathcal{C}} f(c', v', r)} \quad (13)$$

$$\hat{p}(v | \bar{c}, r) = \frac{f(\bar{c}, v, r)}{f(\bar{c}, r)} = \frac{\sum_{c'' \in \bar{\mathcal{C}}} f(c'', v, r)}{\sum_{v' \in \mathcal{V}} \sum_{c'' \in \bar{\mathcal{C}}} f(c'', v', r)} \quad (14)$$

where  $f(c, v, r)$  is the number of  $(n, v, r)$  triples in the data in which  $n$  is being used to denote  $c$ , and  $\mathcal{V}$  is the set of verbs in the data. The problem is that the estimates are defined in terms of frequencies of senses, whereas the data are assumed to be in the form of  $(n, v, r)$  triples: a noun, verb, and argument position. All the data used in this work have been obtained from the British National Corpus (BNC), using the system of Briscoe and Carroll (1997), which consists of a shallow-parsing component that is able to identify verbal arguments.

We take a simple approach to the problem of estimating the frequencies of senses, by distributing the count for each noun in the data evenly among all senses of the noun:

$$\hat{f}(c, v, r) = \sum_{n \in \text{syn}(c)} \frac{f(n, v, r)}{|\text{cn}(n)|} \quad (15)$$

where  $\hat{f}(c, v, r)$  is an estimate of the number of times that concept  $c$  appears in position  $r$  of verb  $v$ , and  $|\text{cn}(n)|$  is the cardinality of  $\text{cn}(n)$ . This is the approach taken by Li and Abe (1998), Ribas (1995), and McCarthy (2000).<sup>7</sup> Resnik (1998) explains how this apparently crude technique works surprisingly well. Alternative approaches are described in Clark and Weir (1999) (see also Clark [2001]), Abney and Light (1999), and Ciaramita and Johnson (2000).

#### 4. Using a Chi-Square Test to Compare Probabilities

In this section we show how to test whether  $p(v | \bar{c}, r)$  changes significantly when considering a node higher in the hierarchy. Consider the problem of deciding whether  $p(\text{run} | \langle \text{canine} \rangle, \text{subj})$  is a good approximation of  $p(\text{run} | \langle \text{dog} \rangle, \text{subj})$ . ( $\langle \text{canine} \rangle$  is the parent of  $\langle \text{dog} \rangle$  in WordNet.) To do this, the probabilities  $p(\text{run} | \bar{c}_i, \text{subj})$  are compared using a chi-square test, where the  $\bar{c}_i$  are the children of  $\langle \text{canine} \rangle$ . In this case, the null hypothesis of the test is that the probabilities  $p(\text{run} | \bar{c}_i, \text{subj})$  are the same for each child  $\bar{c}_i$ . By judging the strength of the evidence against the null hypothesis, how similar the true probabilities are likely to be can be determined. If the test indicates that the probabilities are sufficiently unlikely to be the same, then the null hypothesis is rejected, and the conclusion is that  $p(\text{run} | \langle \text{canine} \rangle, \text{subj})$  is not a good approximation of  $p(\text{run} | \langle \text{dog} \rangle, \text{subj})$ .

An example contingency table, based on counts obtained from a subset of the BNC using the system of Briscoe and Carroll, is given in Table 1. (Recall that the frequencies are estimated by distributing the count for a noun equally among the noun's senses; this explains the fractional counts.) One column contains estimates of counts arising

<sup>7</sup> Resnik takes a similar approach but divides the count evenly among the noun's senses and all the hypernyms of those senses.

**Table 1**Contingency table for the children of  $\langle \text{canine} \rangle$  in the subject position of *run*.

$\bar{c}_i$	$\hat{f}(\bar{c}_i, \text{run}, \text{subj})$	$\hat{f}(\bar{c}_i, \text{subj}) - \hat{f}(\bar{c}_i, \text{run}, \text{subj})$	$\hat{f}(\bar{c}_i, \text{subj}) = \sum_{v \in \mathcal{V}} \hat{f}(\bar{c}_i, v, \text{subj})$
$\langle \text{bitch} \rangle$	0.3 (0.5)	26.7 (26.6)	27.0
$\langle \text{dog} \rangle$	12.8 (10.5)	620.4 (622.7)	633.2
$\langle \text{wolf} \rangle$	0.3 (0.6)	38.7 (38.4)	39.0
$\langle \text{jackal} \rangle$	0.0 (0.3)	20.0 (19.7)	20.0
$\langle \text{wild\_dog} \rangle$	0.0 (0.0)	3.0 (3.0)	3.0
$\langle \text{hyena} \rangle$	0.0 (0.2)	10.0 (9.8)	10.0
$\langle \text{fox} \rangle$	0.0 (1.2)	72.3 (71.1)	72.3
	13.4	791.1	804.5

from concepts in  $\bar{c}_i$  appearing in the subject position of the verb *run*:  $\hat{f}(\bar{c}_i, \text{run}, \text{subj})$ . A second column presents estimates of counts arising from concepts in  $\bar{c}_i$  appearing in the subject position of a verb other than *run*. The figures in brackets are the expected values if the null hypothesis is true.

There is a choice of which statistic to use in conjunction with the chi-square test. The usual statistic encountered in textbooks is the Pearson chi-square statistic, denoted  $X^2$ :

$$X^2 = \sum_{i,j} \frac{(o_{ij} - e_{ij})^2}{e_{ij}} \quad (16)$$

where  $o_{ij}$  is the observed value for the cell in row  $i$  and column  $j$ , and  $e_{ij}$  is the corresponding expected value. An alternative statistic is the log-likelihood chi-square statistic, denoted  $G^2$ :<sup>8</sup>

$$G^2 = 2 \sum_{i,j} o_{ij} \log_e \frac{o_{ij}}{e_{ij}} \quad (17)$$

The two statistics have similar values when the counts in the contingency table are large (Agresti 1996). The statistics behave differently, however, when the table contains low counts, and, since corpus data are likely to lead to some low counts, the question of which statistic to use is an important one. Dunning (1993) argues for the use of  $G^2$  rather than  $X^2$ , based on an analysis of the sampling distributions of  $G^2$  and  $X^2$ , and results obtained when using the statistics to acquire highly associated bigrams. We consider Dunning's analysis at the end of this section, and the question of whether to use  $G^2$  or  $X^2$  will be discussed further there. For now, we continue with the discussion of how the chi-square test is used in the generalization procedure.

For Table 1, the value of  $G^2$  is 3.8, and the value of  $X^2$  is 2.5. Assuming a level of significance of  $\alpha = 0.05$ , the critical value is 12.6 (for six degrees of freedom). Thus, for this  $\alpha$  value, the null hypothesis would not be rejected for either statistic, and the conclusion would be that there is no reason to suppose that  $p(\text{run} \mid \langle \text{canine} \rangle, \text{subj})$  is not a reasonable approximation of  $p(\text{run} \mid \langle \text{dog} \rangle, \text{subj})$ .

<sup>8</sup> An alternative formula for  $G^2$  is given in Dunning (1993), but the two are equivalent.

**Table 2**Contingency table for the children of ⟨liquid⟩ in the object position of *drink*.

$\bar{c}_i$	$\hat{f}(\bar{c}_i, \text{drink}, \text{obj})$	$\hat{f}(\bar{c}_i, \text{obj}) - \hat{f}(\bar{c}_i, \text{drink}, \text{obj})$	$\hat{f}(\bar{c}_i, \text{obj}) = \sum_{v \in \mathcal{V}} \hat{f}(\bar{c}_i, v, \text{obj})$
⟨beverage⟩	261.0 (238.7)	2,367.7 (2,390.0)	2,628.7
⟨supernatant⟩	0.0 (0.1)	1.0 (0.9)	1.0
⟨alcohol⟩	11.5 (9.4)	92.0 (94.1)	103.5
⟨ammonia⟩	0.0 (0.8)	8.5 (7.7)	8.5
⟨antifreeze⟩	0.0 (0.1)	1.0 (0.9)	1.0
⟨distillate⟩	0.0 (0.5)	6.0 (5.5)	6.0
⟨water⟩	12.0 (31.6)	335.7 (316.1)	347.7
⟨ink⟩	0.0 (2.9)	32.0 (29.1)	32.0
⟨liquor⟩	0.7 (1.1)	11.6 (11.2)	12.3
	285.2	2,855.5	3,140.7

As a further example, Table 2 gives counts for the children of ⟨liquid⟩ in the object position of *drink*. Again, the counts have been obtained from a subset of the BNC using the system of Briscoe and Carroll. Not all the sets dominated by the children of ⟨liquid⟩ are shown, as some, such as ⟨sheep\_dip⟩, never appear in the object position of a verb in the data. This example is designed to show a case in which the null hypothesis is rejected. The value of  $G^2$  for this table is 29.0, and the value of  $X^2$  is 21.2. So for  $G^2$ , even if an  $\alpha$  value as low as 0.0005 were being used (for which the critical value is 27.9 for eight degrees of freedom), the null hypothesis would still be rejected. For  $X^2$ , the null hypothesis is rejected for  $\alpha$  values greater than 0.005. This seems reasonable, since the probabilities associated with the children of ⟨liquid⟩ and the object position of *drink* would be expected to show a lot of variation across the children.

A key question is how to select the appropriate value for  $\alpha$ . One solution is to treat  $\alpha$  as a parameter and set it empirically by taking a held-out test set and choosing the value of  $\alpha$  that maximizes performance on the relevant task. For example, Clark and Weir (2000) describes a prepositional phrase attachment algorithm that employs probability estimates obtained using the WordNet method described here. To set the value of  $\alpha$ , the performance of the algorithm on a development set could be compared across different values of  $\alpha$ , and the value that leads to the best performance could be chosen. Note that this approach sets no constraints on the value of  $\alpha$ : The value could be as high as 0.995 or as low as 0.0005, depending on the particular application.

There may be cases in which the conditions for the appropriate application of a chi-square test are not met. One condition that is likely to be violated is the requirement that expected values in the contingency table not be too small. (A rule of thumb often found in textbooks is that the expected values should be greater than five.) One response to this problem is to apply some kind of thresholding and either ignore counts below the threshold, or apply the test only to tables that do not contain low counts. Ribas (1995), Li and Abe (1998), McCarthy (2000), and Wagner (2000) all use some kind of thresholding when dealing with counts in the hierarchy (although not in the context of a chi-square test). Another approach would be to use Fisher's exact test (Agresti 1996; Pedersen 1996), which can be applied to tables regardless of the size of



the counts they contain. The main problem with this test is that it is computationally expensive, especially for large contingency tables.

What we have found in practice is that applying the chi-square test to tables dominated by low counts tends to produce an insignificant result, and the null hypothesis is not rejected. The consequences of this for the generalization procedure are that low-count tables tend to result in the procedure moving up to the next node in the hierarchy. But given that the purpose of the generalization is to overcome the sparse-data problem, moving up a node is desirable, and therefore we do not modify the test for tables with low counts.

The final issue to consider is which chi-square statistic to use. Dunning (1993) argues for the use of  $G^2$  rather than  $X^2$ , based on the claim that the sampling distribution of  $G^2$  approaches the true chi-square distribution quicker than the sampling distribution of  $X^2$ . However, Agresti (1996, page 34) makes the opposite claim: "The sampling distributions of  $X^2$  and  $G^2$  get closer to chi-squared as the sample size  $n$  increases. . . . The convergence is quicker for  $X^2$  than  $G^2$ ."

In addition, Pedersen (2001) questions whether one statistic should be preferred over the other for the bigram acquisition task and cites Cressie and Read (1984), who argue that there are some cases where the Pearson statistic is more reliable than the log-likelihood statistic. Finally, the results of the pseudo-disambiguation experiments presented in Section 7 are at least as good, if not better, when using  $X^2$  rather than  $G^2$ , and so we conclude that the question of which statistic to use should be answered on a per application basis.

## 5. The Generalization Procedure

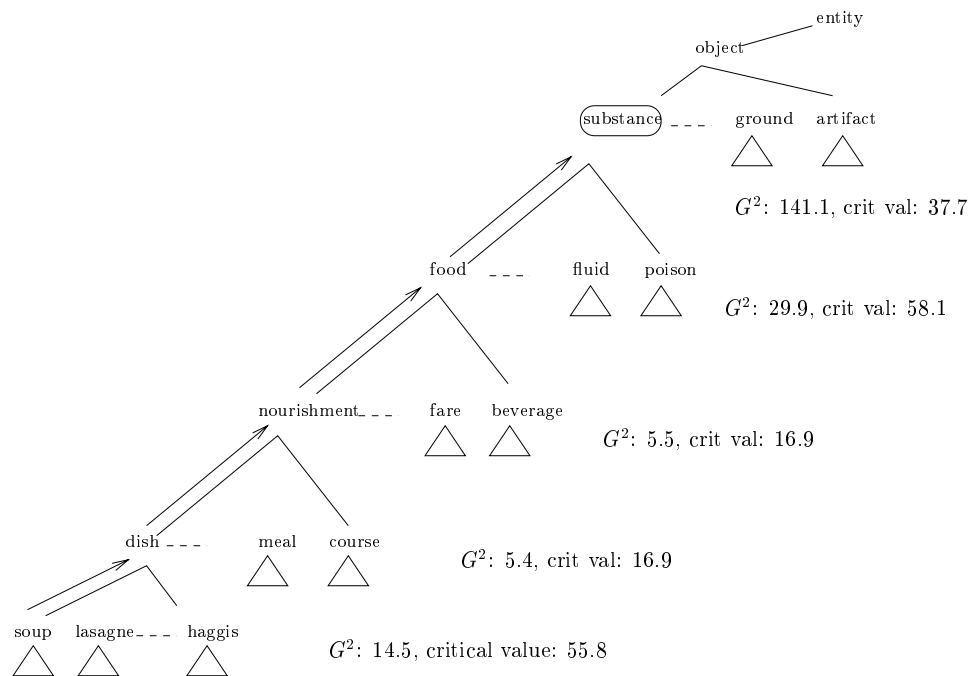
The procedure for finding a suitable class,  $\bar{c}$ , to generalize concept  $c$  in position  $r$  of verb  $v$  works as follows. (We refer to  $\bar{c}$  as the "similarity class" of  $c$  with respect to  $v$  and  $r$  and the hypernym  $c'$  as  $\text{top}(c, v, r)$ , since the chosen hypernym sits at the "top" of the similarity class.) Initially, concept  $c$  is assigned to a variable  $\text{top}$ . Then, by working up the hierarchy, successive hypernyms of  $c$  are assigned to  $\text{top}$ , and this process continues until the probabilities associated with the sets of concepts dominated by  $\text{top}$  and the siblings of  $\text{top}$  are significantly different. Once a node is reached that results in a significant result for the chi-square test, the procedure stops, and  $\text{top}$  is returned as  $\text{top}(c, v, r)$ . In cases where a concept has more than one parent, the parent is chosen that results in the lowest value of the chi-square statistic, as this indicates the probabilities are the most similar. The set  $\overline{\text{top}(c, v, r)}$  is the similarity class of  $c$  for verb  $v$  and position  $r$ . Figure 1 gives an algorithm for determining  $\text{top}(c, v, r)$ .

Figure 2 gives an example of the procedure at work. Here,  $\text{top}(\langle \text{soup} \rangle, \text{stir}, \text{obj})$  is being determined. The example is based on data from a subset of the BNC, with 303 cases of an argument in the object position of *stir*. The  $G^2$  statistic is used, together with an  $\alpha$  value of 0.05. Initially,  $\text{top}$  is set to  $\langle \text{soup} \rangle$ , and the probabilities corresponding to the children of  $\langle \text{dish} \rangle$  are compared:  $p(\text{stir} | \langle \text{soup} \rangle, \text{obj})$ ,  $p(\text{stir} | \langle \text{lasagne} \rangle, \text{obj})$ ,  $p(\text{stir} | \langle \text{haggis} \rangle, \text{obj})$ , and so on for the rest of the children. The chi-square test results in a  $G^2$  value of 14.5, compared to a critical value of 55.8. Since  $G^2$  is less than the critical value, the procedure moves up to the next node. This process continues until a significant result is obtained, which first occurs at  $\langle \text{substance} \rangle$  when comparing the children of  $\langle \text{object} \rangle$ . Thus  $\langle \text{substance} \rangle$  is the chosen level of generalization.

Now we show how the chosen level of generalization varies with  $\alpha$  and how it varies with the size of the data set. A note of clarification is required before presenting the results. In related work on acquiring selectional preferences (Ribas 1995; McCarthy

**Algorithm**  $\text{top}(c, v, r)$ :  
 $\text{top} \leftarrow c$   
 $\text{sig\_result} \leftarrow \text{false}$   
**comment**  $\text{parent}_{\min}$  gives lowest  $G^2$  value,  $G^2_{\min}$   
**while** not  $\text{sig\_result}$  &  $\text{top} \neq \langle \text{root} \rangle$  **do**  
     $G^2_{\min} \leftarrow \infty$   
    **for all** parents of  $\text{top}$  **do**  
        calculate  $G^2$  for sets dominated by children of parent  
        **if**  $G^2 < G^2_{\min}$   
            **then**  $G^2_{\min} \leftarrow G^2$   
                 $\text{parent}_{\min} \leftarrow \text{parent}$   
    **end**  
    **if** chi-square test for  $\text{parent}_{\min}$  is significant  
        **then**  $\text{sig\_result} \leftarrow \text{true}$   
        **else** move up to next node:  $\text{top} \leftarrow \text{parent}_{\min}$   
**end**  
return  $\text{top}$

**Figure 1**  
An algorithm for determining  $\text{top}(c, v, r)$ .



**Figure 2**  
An example generalization: Determining  $\text{top}(\langle \text{soup} \rangle, \text{stir}, \text{obj})$ .

1997; Li and Abe 1998; Wagner 2000), the level of generalization is often determined for a small number of hand-picked verbs and the result compared with the researcher's intuition about the most appropriate level for representing a selectional preference. According to this approach, if  $\langle \text{sandwich} \rangle$  were chosen to represent  $\langle \text{hotdog} \rangle$  in the object position of *eat*, this might be considered an undergeneralization, since  $\langle \text{food} \rangle$  might be considered more appropriate. For this work we argue that such an evaluation is not appropriate; since the purpose of this work is probability estimation, the most appropriate level is the one that leads to the most accurate estimate, and this may or may not agree with intuition. Furthermore, we show in Section 7 that to generalize unnecessarily can be harmful for some tasks: If we already have lots of data regarding  $\langle \text{sandwich} \rangle$ , why generalize any higher? Thus the purpose of this section is not to show that the acquired levels are "correct," but simply to show how the levels vary with  $\alpha$  and the sample size.

To show how the level of generalization varies with changes in  $\alpha$ ,  $\text{top}(c, v, \text{obj})$  was determined for a number of hand-picked  $(c, v, \text{obj})$  triples over a range of values for  $\alpha$ . The triples were chosen to give a range of strongly and weakly selecting verbs and a range of verb frequencies. The data were again extracted from a subset of the BNC using the system of Briscoe and Carroll (1997), and the  $G^2$  statistic was used in the chi-square test. The results are shown in Table 3. The number of times the verb occurred with some object is also given in the table.

The results suggest that the generalization level becomes more specific as  $\alpha$  increases. This is to be expected, since, given a contingency table chosen at random, a higher value of  $\alpha$  is more likely to lead to a significant result than a lower value of  $\alpha$ . We also see that, for some cases, the value of  $\alpha$  has little effect on the level. We would expect there to be less change in the level of generalization for strongly selecting verbs, such as *drink* and *eat*, and a greater range of levels for weakly selecting verbs such as *see*. This is because any significant difference in probabilities is likely to be more marked for a strongly selecting verb, and likely to be significant over a wider range of  $\alpha$  values. The table only provides anecdotal evidence, but provides some support to this argument.

To investigate more generally how the level of generalization varies with changes in  $\alpha$ , and also with changes in sample size, we took 6,000  $(c, v, \text{obj})$  triples and calculated the difference in depth between  $c$  and  $\text{top}(c, v, r)$  for each triple. The 6,000 triples were taken from the first experimental test set described in Section 7, and the training data from this experiment were used to provide the counts. (The test set contains nouns, rather than noun senses, and so the sense of the noun that is most probable given the verb and object slot was used.) An average difference in depth was then calculated. To give an example of how the difference in depth was calculated, suppose  $\langle \text{dog} \rangle$  generalized to  $\langle \text{placental\_mammal} \rangle$  via  $\langle \text{canine} \rangle$  and  $\langle \text{carnivore} \rangle$ ; in this case the difference would be three.

The results for various levels of  $\alpha$  and different sample sizes are shown in Table 4. The figures in each column arise from using the contingency tables based on the complete training data, but with each count in the table multiplied by the percentage at the head of the column. Thus the 50% column is based on contingency tables in which each original count is multiplied by 50%, which is equivalent to using a sample one-half the size of the original training set. Reading across a row shows how the generalization varies with sample size, and reading down a column shows how it varies with  $\alpha$ . The results show clearly that the extent of generalization decreases with an increase in the value of  $\alpha$ , supporting the trend observed in Table 3. The results also show that the extent of generalization increases with a decrease in sample

**Table 3**  
Example levels of generalization for different values of  $\alpha$ .

$(c, v, r), f(v, r)$	$\alpha$	
$(\langle \text{coffee} \rangle, \text{drink}, \text{obj})$ $f(\text{drink}, \text{obj}) = 849$	0.0005	$\langle \text{coffee} \rangle \langle \text{BEVERAGE} \rangle \langle \text{food} \rangle \dots \langle \text{object} \rangle \langle \text{entity} \rangle$
	0.05	$\langle \text{coffee} \rangle \langle \text{BEVERAGE} \rangle \langle \text{food} \rangle \dots \langle \text{object} \rangle \langle \text{entity} \rangle$
	0.5	$\langle \text{coffee} \rangle \langle \text{BEVERAGE} \rangle \langle \text{food} \rangle \dots \langle \text{object} \rangle \langle \text{entity} \rangle$
	0.995	$\langle \text{coffee} \rangle \langle \text{BEVERAGE} \rangle \langle \text{food} \rangle \dots \langle \text{object} \rangle \langle \text{entity} \rangle$
$(\langle \text{hotdog} \rangle, \text{eat}, \text{obj})$ $f(\text{eat}, \text{obj}) = 1,703$	0.0005	$\langle \text{hotdog} \rangle \langle \text{sandwich} \rangle \langle \text{snack\_food} \rangle \langle \text{DISH} \rangle \dots \langle \text{food} \rangle \dots \langle \text{entity} \rangle$
	0.05	$\langle \text{hotdog} \rangle \langle \text{sandwich} \rangle \langle \text{snack\_food} \rangle \langle \text{DISH} \rangle \dots \langle \text{food} \rangle \dots \langle \text{entity} \rangle$
	0.5	$\langle \text{hotdog} \rangle \langle \text{sandwich} \rangle \langle \text{snack\_food} \rangle \langle \text{DISH} \rangle \dots \langle \text{food} \rangle \dots \langle \text{entity} \rangle$
	0.995	$\langle \text{hotdog} \rangle \langle \text{SANDWICH} \rangle \langle \text{snack\_food} \rangle \langle \text{dish} \rangle \dots \langle \text{food} \rangle \dots \langle \text{entity} \rangle$
$(\langle \text{Socrates} \rangle, \text{kiss}, \text{obj})$ $f(\text{kiss}, \text{obj}) = 345$	0.0005	$\langle \text{Socrates} \rangle \dots \langle \text{person} \rangle \langle \text{life\_form} \rangle \langle \text{CAUSAL\_AGENT} \rangle \langle \text{entity} \rangle$
	0.05	$\langle \text{Socrates} \rangle \dots \langle \text{person} \rangle \langle \text{life\_form} \rangle \langle \text{CAUSAL\_AGENT} \rangle \langle \text{entity} \rangle$
	0.5	$\langle \text{Socrates} \rangle \dots \langle \text{person} \rangle \langle \text{life\_form} \rangle \langle \text{CAUSAL\_AGENT} \rangle \langle \text{entity} \rangle$
	0.995	$\langle \text{Socrates} \rangle \dots \langle \text{PERSON} \rangle \langle \text{life\_form} \rangle \langle \text{causal\_agent} \rangle \langle \text{entity} \rangle$
$(\langle \text{dream} \rangle, \text{remember}, \text{obj})$ $f(\text{remember}, \text{obj}) = 1,982$	0.0005	$\langle \text{dream} \rangle \dots \langle \text{preoccupation} \rangle \langle \text{cognitive\_state} \rangle \langle \text{STATE} \rangle$
	0.05	$\langle \text{dream} \rangle \dots \langle \text{preoccupation} \rangle \langle \text{cognitive\_state} \rangle \langle \text{STATE} \rangle$
	0.5	$\langle \text{dream} \rangle \dots \langle \text{preoccupation} \rangle \langle \text{COGNITIVE\_STATE} \rangle \langle \text{state} \rangle$
	0.995	$\langle \text{dream} \rangle \dots \langle \text{PREOCCUPATION} \rangle \langle \text{cognitive\_state} \rangle \langle \text{state} \rangle$
$(\langle \text{man} \rangle, \text{see}, \text{obj})$ $f(\text{see}, \text{obj}) = 16,757$	0.0005	$\langle \text{man} \rangle \dots \langle \text{mammal} \rangle \dots \langle \text{ANIMAL} \rangle \langle \text{life\_form} \rangle \langle \text{entity} \rangle$
	0.05	$\langle \text{man} \rangle \dots \langle \text{MAMMAL} \rangle \dots \langle \text{animal} \rangle \langle \text{life\_form} \rangle \langle \text{entity} \rangle$
	0.5	$\langle \text{man} \rangle \dots \langle \text{MAMMAL} \rangle \dots \langle \text{animal} \rangle \langle \text{life\_form} \rangle \langle \text{entity} \rangle$
	0.995	$\langle \text{MAN} \rangle \dots \langle \text{mammal} \rangle \dots \langle \text{animal} \rangle \langle \text{life\_form} \rangle \langle \text{entity} \rangle$
$(\langle \text{belief} \rangle, \text{abandon}, \text{obj})$ $f(\text{abandon}, \text{obj}) = 673$	0.0005	$\langle \text{belief} \rangle \langle \text{mental\_object} \rangle \langle \text{cognition} \rangle \langle \text{PSYCHOLOGICAL\_FEATURE} \rangle$
	0.05	$\langle \text{belief} \rangle \langle \text{MENTAL\_OBJECT} \rangle \langle \text{cognition} \rangle \langle \text{psychological\_feature} \rangle$
	0.5	$\langle \text{BELIEF} \rangle \langle \text{mental\_object} \rangle \langle \text{cognition} \rangle \langle \text{psychological\_feature} \rangle$
	0.995	$\langle \text{BELIEF} \rangle \langle \text{mental\_object} \rangle \langle \text{cognition} \rangle \langle \text{psychological\_feature} \rangle$
$(\langle \text{nightmare} \rangle, \text{have}, \text{obj})$ $f(\text{have}, \text{obj}) = 93,683$	0.0005	$\langle \text{nightmare} \rangle \langle \text{dreaming} \rangle \langle \text{IMAGINATION} \rangle \dots \langle \text{psychological\_feature} \rangle$
	0.05	$\langle \text{nightmare} \rangle \langle \text{dreaming} \rangle \langle \text{IMAGINATION} \rangle \dots \langle \text{psychological\_feature} \rangle$
	0.5	$\langle \text{nightmare} \rangle \langle \text{DREAMING} \rangle \langle \text{imagination} \rangle \dots \langle \text{psychological\_feature} \rangle$
	0.995	$\langle \text{nightmare} \rangle \langle \text{DREAMING} \rangle \langle \text{imagination} \rangle \dots \langle \text{psychological\_feature} \rangle$

Note: The selected level is shown in upper case.

**Table 4**  
Extent of generalization for different values of  $\alpha$  and sample sizes.

$\alpha$	100%	50%	10%	1%
0.0005	3.3	3.9	5.0	5.6
0.05	2.8	3.5	4.6	5.6
0.5	2.1	2.9	4.1	5.4
0.995	1.2	1.5	2.6	3.9

size. Again, this is to be expected, since any difference in probability estimates is less likely to be significant for tables with low counts.

## 6. Alternative Class-Based Estimation Methods

The approaches used for comparison are that of Resnik (1993, 1998), subsequently developed by Ribas (1995), and that of Li and Abe (1998), which has been adopted by McCarthy (2000). These have been chosen because they directly address the question of how to find a suitable level of generalization in WordNet.

The first alternative uses the “association score,” which is a measure of how well a set of concepts,  $C$ , satisfies the selectional preferences of a verb,  $v$ , for an argument position,  $r$ :<sup>9</sup>

$$A(C, v, r) = p(C | v, r) \log_2 \frac{p(C | v, r)}{p(C | r)} \quad (18)$$

An estimate of the association score,  $\hat{A}(C, v, r)$ , can be obtained using relative frequency estimates of the probabilities. The key question is how to determine a suitable level of generalization for concept  $c$ , or, alternatively, how to find a suitable class to represent concept  $c$  (assuming the choice is from those classes that contain all concepts dominated by some hypernym of  $c$ ). Resnik’s solution to this problem (which he neatly refers to as the “vertical-ambiguity” problem) is to choose the class that maximizes the association score.

It is not clear that the class with the highest association score is always the most appropriate level of generalization. For example, this approach does not always generalize appropriately for arguments that are *negatively* associated with some verb. To see why, consider the problem of deciding how well the concept ⟨location⟩ satisfies the preferences of the verb *eat* for its object. Since locations are not the kinds of things that are typically eaten, a suitable level of generalization would correspond to a class that has a low association score with respect to *eat*. However, ⟨location⟩ is a kind of ⟨entity⟩ in WordNet,<sup>10</sup> and choosing the class with the highest association score is likely to produce ⟨entity⟩ as the chosen class. This is a problem, because the association score of ⟨entity⟩ with respect to *eat* may be too high to reflect the fact that ⟨location⟩ is a very unlikely object of the verb.

Note that the solution to the vertical-ambiguity problem presented in the previous sections is able to generalize appropriately in such cases. Continuing with the *eat* ⟨location⟩ example, our generalization procedure is unlikely to get as high as ⟨entity⟩ (assuming a reasonable number of examples of *eat* in the training data), since the probabilities corresponding to the daughters of ⟨entity⟩ are likely to be very different with respect to the object position of *eat*.

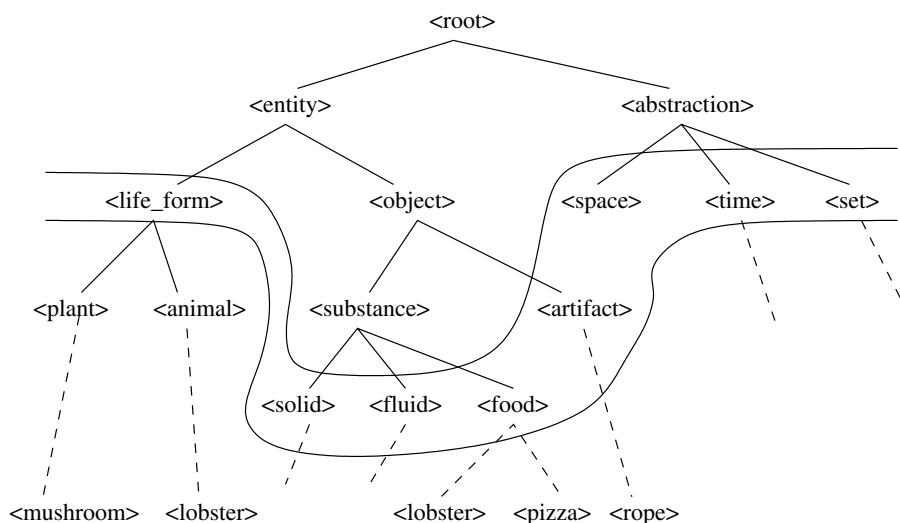
The second alternative uses the minimum description length (MDL) principle. Li and Abe use MDL to select a set of classes from a hierarchy, together with their associated probabilities, to represent the selectional preferences of a particular verb. The preferences and class-based probabilities are then used to estimate probabilities of the form  $p(n | v, r)$ , where  $n$  is a noun,  $v$  is a verb, and  $r$  is an argument slot.

Li and Abe’s application of MDL requires the hierarchy to be in the form of a thesaurus, in which each leaf node represents a noun and internal nodes represent the class of nouns that the node dominates. The hierarchy is also assumed to be in the form of a tree. The class-based models consist of a partition of the set of nouns (leaf nodes) and a probability associated with each class in the partition. The probabilities are the conditional probabilities of each class, given the relevant verb and argument position. Li and Abe refer to such a partition as a “cut” and the cut together with the probabilities as a “tree cut model.” The probabilities of the classes in a cut,  $\Gamma$ , satisfy the following constraint:

$$\sum_{C \in \Gamma} p(C | v, r) = 1 \quad (19)$$

<sup>9</sup> The definition used here is that given by Ribas (1995).

<sup>10</sup> For example, the hypernyms of the concept ⟨Dallas⟩ are as follows: ⟨city⟩, ⟨municipality⟩, ⟨urban\_area⟩, ⟨geographical\_area⟩, ⟨region⟩, ⟨location⟩, ⟨object⟩, ⟨entity⟩.



**Figure 3**  
Possible cut returned by MDL.

In order to determine the probability of a noun, the probability of a class is assumed to be distributed uniformly among the members of that class:

$$p(n | v, r) = \frac{1}{|C|} p(C | v, r) \quad \text{for all } n \in C \quad (20)$$

Since WordNet is a hierarchy with noun senses, rather than nouns, at the nodes, Li and Abe deal with the issue of word sense ambiguity using the method described in Section 3, by dividing the count for a noun equally among the concepts whose synsets contain the noun. Also, since WordNet is a DAG, Li and Abe turn WordNet into a tree by copying each subgraph with multiple parents. And so that each noun in the data appears (in a synset) at a leaf node, Li and Abe remove those parts of the hierarchy dominated by a noun in the data (but only for that instance of WordNet corresponding to the relevant verb).

An example cut showing part of the WordNet hierarchy is shown in Figure 3 (based on an example from Li and Abe [1998]; the dashed lines indicate parts of the hierarchy that are not shown in the diagram). This is a possible cut for the object position of the verb *eat*, and the cut consists of the following classes:  $\langle \text{life\_form} \rangle$ ,  $\langle \text{solid} \rangle$ ,  $\langle \text{fluid} \rangle$ ,  $\langle \text{food} \rangle$ ,  $\langle \text{artifact} \rangle$ ,  $\langle \text{space} \rangle$ ,  $\langle \text{time} \rangle$ ,  $\langle \text{set} \rangle$ . (The particular choice of classes for the cut in this example is not too important; the example is designed to show how probabilities of senses are estimated from class probabilities.) Since the class in the cut containing  $\langle \text{pizza} \rangle$  is  $\langle \text{food} \rangle$ , the probability  $p(\langle \text{pizza} \rangle | \text{eat, obj})$  would be estimated as  $p(\langle \text{food} \rangle | \text{eat, obj}) / |\langle \text{food} \rangle|$ . Similarly, since the class in the cut containing  $\langle \text{mushroom} \rangle$  is  $\langle \text{life\_form} \rangle$ , the probability  $p(\langle \text{mushroom} \rangle | \text{eat, obj})$  would be estimated as  $p(\langle \text{life\_form} \rangle | \text{eat, obj}) / |\langle \text{life\_form} \rangle|$ .

The uniform-distribution assumption (20) means that cuts close to the root of the hierarchy result in a greater smoothing of the probability estimates than cuts near the leaves. Thus there is a trade-off between choosing a model that has a cut near the leaves, which is likely to overfit the data, and a more general (simple) model near the root, which is likely to underfit the data. MDL looks ideally suited to the task of model selection, since it is designed to deal with precisely this trade-off. The simplicity of a model is measured using the *model description length*, which is an information-theoretic

term and denotes the number of bits required to encode the model. The fit to the data is measured using the *data description length*, which is the number of bits required to encode the data (relative to the model). The overall description length is the sum of the model description length and the data description length, and the MDL principle is to select the model with the shortest description length.

We used McCarthy's (2000) implementation of MDL. So that every noun is represented at a leaf node, McCarthy does not remove parts of the hierarchy, as Li and Abe do, but instead creates new leaf nodes for each synset at an internal node. McCarthy also does not transform WordNet into a tree, which is strictly required for Li and Abe's application of MDL. This did create a problem with overgeneralization: Many of the cuts returned by MDL were overgeneralizing at the  $\langle \text{entity} \rangle$  node. The reason is that  $\langle \text{person} \rangle$ , which is close to  $\langle \text{entity} \rangle$  and dominated by  $\langle \text{entity} \rangle$ , has two parents:  $\langle \text{life\_form} \rangle$  and  $\langle \text{causal\_agent} \rangle$ . This DAG-like property was responsible for the overgeneralization, and so we removed the link between  $\langle \text{person} \rangle$  and  $\langle \text{causal\_agent} \rangle$ . This appeared to solve the problem, and the results presented later for the average degree of generalization do not show an overgeneralization compared with those given in Li and Abe (1998).

## 7. Pseudo-Disambiguation Experiments

The task we used to compare the class-based estimation techniques is a decision task previously used by Pereira, Tishby, and Lee (1993) and Rooth et al. (1999). The task is to decide which of two verbs,  $v$  and  $v'$ , is more likely to take a given noun,  $n$ , as an object. The test and training data were obtained as follows. A number of verb-direct object pairs were extracted from a subset of the BNC, using the system of Briscoe and Carroll. All those pairs containing a noun not in WordNet were removed, and each verb and argument was lemmatized. This resulted in a data set of around 1.3 million  $(v, n)$  pairs.

To form a test set, 3,000 of these pairs were randomly selected such that each selected pair contained a fairly frequent verb. (Following Pereira, Tishby, and Lee, only those verbs that occurred between 500 and 5,000 times in the data were considered.) Each instance of a selected pair was then deleted from the data to ensure that the test data were unseen. The remaining pairs formed the training data. To complete the test set, a further fairly frequent verb,  $v'$ , was randomly chosen for each  $(v, n)$  pair. The random choice was made according to the verb's frequency in the original data set, subject to the condition that the pair  $(v', n)$  did not occur in the training data. Given the set of  $(v, n, v')$  triples, the task is to decide whether  $(v, n)$  or  $(v', n)$  is the correct pair.<sup>11</sup>

We acknowledge that the task is somewhat artificial, but pseudo-disambiguation tasks of this kind are becoming popular in statistical NLP because of the ease with which training and test data can be created. We also feel that the pseudo-disambiguation task is useful for evaluating the different estimation methods, since it directly addresses the question of how likely a particular predicate is to take a given noun as an argument. An evaluation using a PP attachment task was attempted in Clark and Weir (2000), but the evaluation was limited by the relatively small size of the Penn Treebank.

<sup>11</sup> We note that this procedure does not guarantee that the correct pair is more likely than the incorrect pair, because of noise in the data from the parser and also because a highly plausible incorrect pair could be generated by chance.

**Table 5**  
Results for the pseudo-disambiguation task.

Generalization technique	% correct	av.gen.	sd.gen.
Similarity class			
$\alpha = 0.0005$	73.8	3.3	2.0
$\alpha = 0.05$	73.4	2.8	1.9
$\alpha = 0.3$	73.0	2.4	1.8
$\alpha = 0.75$	73.9	1.9	1.6
$\alpha = 0.995$	73.8	1.2	1.2
Low class	73.6	0.9	1.0
<hr/>			
MDL	68.3	4.1	1.9
Assoc	63.9	4.2	2.1

Note: av.gen. is the average number of generalized levels; sd.gen. is the standard deviation.

Using our approach, the disambiguation decision for each  $(v, n, v')$  triple was made according to the following procedure:

```

if  $\max_{c \in \text{CN}(n)} p_{sc}(c | v, \text{obj}) > \max_{c \in \text{CN}(n)} p_{sc}(c | v', \text{obj})$ 
  then choose  $(v, n)$ 
else if  $\max_{c \in \text{CN}(n)} p_{sc}(c | v', \text{obj}) > \max_{c \in \text{CN}(n)} p_{sc}(c | v, \text{obj})$ 
  then choose  $(v', n)$ 
else choose at random

```

If  $n$  has more than one sense, the sense is chosen that maximizes the relevant probability estimate; this explains the maximization over  $\text{CN}(n)$ . The probability estimates were obtained using our class-based method, and the  $G^2$  statistic was used for the chi-square test. This procedure was also used for the MDL alternative, but using the MDL method to estimate the probabilities.

Using the association score for each test triple, the decision was made according to the following procedure:

```

if  $\max_{c \in \text{CN}(n)} \max_{c' \in \text{h}(c)} \hat{A}(\bar{c}', v, \text{obj}) > \max_{c \in \text{CN}(n)} \max_{c' \in \text{h}(c)} \hat{A}(\bar{c}', v', \text{obj})$ 
  then choose  $(v, n)$ 
else if  $\max_{c \in \text{CN}(n)} \max_{c' \in \text{h}(c)} \hat{A}(\bar{c}', v', \text{obj}) > \max_{c \in \text{CN}(n)} \max_{c' \in \text{h}(c)} \hat{A}(\bar{c}', v, \text{obj})$ 
  then choose  $(v', n)$ 
else choose at random

```

We use  $\text{h}(c)$  to denote the set consisting of the hypernyms of  $c$ . The inner maximization is over  $\text{h}(c)$ , assuming  $c$  is the chosen sense of  $n$ , which corresponds to Resnik's method of choosing a set to represent  $c$ . The outer maximization is over the senses of  $n$ ,  $\text{CN}(n)$ , which determines the sense of  $n$  by choosing the sense that maximizes the association score.

The first set of results is given in Table 5. Our technique is referred to as the "similarity class" technique, and the approach using the association score is referred



**Table 6**

Results for the pseudo-disambiguation task with one-fifth training data.

Generalization technique	% correct	av.gen.	sd.gen.
Similarity class			
$\alpha = 0.0005$	66.7	4.5	1.9
$\alpha = 0.05$	68.4	4.1	1.9
$\alpha = 0.3$	70.2	3.7	1.9
$\alpha = 0.75$	72.3	3.0	1.9
$\alpha = 0.995$	72.4	1.9	1.6
Low class	71.9	1.1	1.1
<hr/>			
MDL	62.9	4.7	1.9
Assoc	62.6	4.1	2.0

Note: av.gen. is the average number of generalized levels; sd.gen. is the standard deviation.

to as “Assoc.” The results are given for a range of  $\alpha$  values and demonstrate clearly that the performance of similarity class varies little with changes in  $\alpha$  and that similarity class outperforms both MDL and Assoc.<sup>12</sup>

We also give a score for our approach using a simple generalization procedure, which we call “low class.” The procedure is to select the first class that has a count greater than zero (relative to the verb and argument position), which is likely to return a low level of generalization, on the whole. The results show that our generalization technique only narrowly outperforms the simple alternative. Note that, although low class is based on a very simple generalization method, the estimation method is still using our class-based technique, by applying Bayes’ theorem and conditioning on a class, as described in Section 3; the difference is in how the class is chosen.

To investigate the results, we calculated the average number of generalized levels for each approach. The number of generalized levels for a concept  $c$  (relative to a verb  $v$  and argument position  $r$ ) is the difference in depth between  $c$  and  $\text{top}(c, v, r)$ , as explained in Section 5. For each test case, the number of generalized levels for both verbs,  $v$  and  $v'$ , was calculated, but only for the chosen sense of  $n$ . The results are given in the third column of Table 5 and demonstrate clearly that both MDL and Assoc are generalizing to a greater extent than similarity class. (The fourth column gives a standard deviation figure.) These results suggest that MDL and Assoc are overgeneralizing, at least for the purposes of this task.

To investigate why the value for  $\alpha$  had no impact on the results, we repeated the experiment, but with one fifth of the data. A new data set was created by taking every fifth pair of the original 1.3 million pairs. A test set of 3,000 triples was created from this new data set, as before, but this time only verbs that occurred between 100 and 1,000 times were considered. The results using these test and training data are given in Table 6.

These results show a variation in performance across values for  $\alpha$ , with an optimal performance when  $\alpha$  is around 0.75. (Of course, in practice, the value for  $\alpha$  would need to be optimized on a held-out set.) But even with this variation, similarity class is still outperforming MDL and Assoc across the whole range of  $\alpha$  values. Note that the

<sup>12</sup> The results given for similarity class are different from those given in Clark and Weir (2001) because the probability estimates used in Clark and Weir (2001) were not normalized.

**Table 7**  
Disambiguation results for  $G^2$  and  $X^2$ .

$\alpha$ value	% correct ( $G^2$ )		% correct ( $X^2$ )	
0.0005	73.8	(3.3)	74.1	(3.0)
0.05	73.4	(2.8)	73.8	(2.5)
0.3	73.0	(2.4)	74.1	(2.2)
0.75	73.9	(1.9)	74.3	(1.8)
0.995	73.8	(1.2)	73.3	(1.2)

$\alpha$  values corresponding to the lowest scores lead to a significant amount of generalization, which provides additional evidence that MDL and Assoc are overgeneralizing for this task. The low-class method scores highly for this data set also, but given that the task is one that apparently favors a low level of generalization, the high score is not too surprising.

As a final experiment, we compared the task performance using the  $X^2$ , rather than  $G^2$ , statistic in the chi-square test. The results are given in Table 7 for the complete data set.<sup>13</sup> The figures in brackets give the average number of generalized levels. The  $X^2$  statistic is performing at least as well as  $G^2$ , and the results show that the average level of generalization is slightly higher for  $G^2$  than  $X^2$ . This suggests a possible explanation for the results presented here and those in Dunning (1993): that the  $X^2$  statistic provides a less conservative test when counts in the contingency table are low. (By a conservative test we mean one in which the null hypothesis is not easily rejected.) A less conservative test is better suited to the pseudo-disambiguation task, since it results in a lower level of generalization, on the whole, which is good for this task. In contrast, the task that Dunning considers, the discovery of bigrams, is better served by a more conservative test.

## 8. Conclusion

We have presented a class-based estimation method that incorporates a procedure for finding a suitable level of generalization in WordNet. This method has been shown to provide superior performance on a pseudo-disambiguation task, compared with two alternative approaches. An analysis of the results has shown that the other approaches appear to be overgeneralizing, at least for this task. One of the features of the generalization procedure is the way that  $\alpha$ , the level of significance in the chi-square test, is treated as a parameter. This allows some control over the extent of generalization, which can be tailored to particular tasks. We have also shown that the task performance is at least as good when using the Pearson chi-square statistic as when using the log-likelihood chi-square statistic.

There are a number of ways in which this work could be extended. One possibility would be to use all the classes dominated by the hypernyms of a concept, rather than just one, to estimate the probability of the concept. An estimate would be obtained for each hypernym, and the estimates combined in a linear interpolation. An approach similar to this is taken by Bikel (2000), in the context of statistical parsing.

There is still room for investigation of the hidden-data problem when data are used that have not been sense disambiguated. In this article, a very simple approach is taken,

<sup>13</sup>  $\chi^2$  performed slightly better than  $G^2$  using the smaller data set also.

which is to split the count for a noun evenly among the noun's senses. Abney and Light (1999) have tried a more motivated approach, using the expectation maximization algorithm, but with little success. The approach described in Clark and Weir (1999) is shown in Clark (2001) to have some impact on the pseudo-disambiguation task, but only with certain values of the  $\alpha$  parameter, and ultimately does not improve on the best performance.

Finally, an issue that has not been much addressed in the literature (except by Li and Abe [1996]) is how the accuracy of class-based estimation techniques compare when automatically acquired classes, as opposed to the manually created classes from WordNet, are used. The pseudo-disambiguation task described here has also been used to evaluate clustering algorithms (Pereira, Tishby, and Lee, 1993; Rooth et al., 1999), but with different data, and so it is difficult to compare the results. A related issue is how the structure of WordNet affects the accuracy of the probability estimates. We have taken the structure of the hierarchy for granted, without any analysis, but it may be that an alternative design could be more conducive to probability estimation.

### Acknowledgments

This article is an extended and updated version of a paper that appeared in the proceedings of NAACL 2001. The work on which it is based was carried out while the first author was a D.Phil. student at the University of Sussex and was supported by an EPSRC studentship. We would like to thank Diana McCarthy for suggesting the pseudo-disambiguation task and providing the MDL software, John Carroll for supplying the data, and Ted Briscoe, Geoff Sampson, Gerald Gazdar, Bill Keller, Ted Pedersen, and the anonymous reviewers for their helpful comments. We would also like to thank Ted Briscoe for presenting an earlier version of this article on our behalf at NAACL 2001.

### References

- Abney, Steven P. and Marc Light. 1999. Hiding a semantic hierarchy in a Markov model. In *Proceedings of the ACL Workshop on Unsupervised Learning in Natural Language Processing*, University of Maryland, College Park, pages 1–8.
- Agirre, Eneko and David Martinez. 2001. Learning class-to-class selectional preferences. In *Proceedings of the Fifth ACL Workshop on Computational Language Learning*, Toulouse, France, pages 15–22.
- Agresti, Alan. 1996. *An Introduction to Categorical Data Analysis*. Wiley.
- Bikel, Daniel M. 2000. A statistical model for parsing and word-sense disambiguation. In *Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pages 155–163, Hong Kong.
- Briscoe, Ted and John Carroll. 1997. Automatic extraction of subcategorization from corpora. In *Proceedings of the Fifth ACL Conference on Applied Natural Language Processing*, pages 356–363, Washington, DC.
- Church, Kenneth W. and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29.
- Ciaramita, Massimiliano and Mark Johnson. 2000. Explaining away ambiguity: Learning verb selectional preference with Bayesian networks. In *Proceedings of the 18th International Conference on Computational Linguistics*, pages 187–193, Saarbrücken, Germany.
- Clark, Stephen. 2001. *Class-Based Statistical Models for Lexical Knowledge Acquisition*. Ph.D. dissertation, University of Sussex.
- Clark, Stephen and David Weir. 1999. An iterative approach to estimating frequencies over a semantic hierarchy. In *Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pages 258–265, University of Maryland, College Park.
- Clark, Stephen and David Weir. 2000. A class-based probabilistic approach to structural disambiguation. In *Proceedings of the 18th International Conference on Computational Linguistics*, pages 194–200, Saarbrücken, Germany.
- Clark, Stephen and David Weir. 2001. Class-based probability estimation using a semantic hierarchy. In *Proceedings of the Second Meeting of the North American Chapter of the Association for Computational Linguistics*, pages 95–102, Pittsburgh.
- Cressie, Noel A. C. and Timothy R. C. Read. 1984. Multinomial goodness of fit tests.

- Journal of the Royal Statistics Society Series B*, 46:440–464.
- Dunning, Ted. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61–74.
- Fellbaum, Christiane, editor. 1998. *WordNet: An Electronic Lexical Database*. MIT Press.
- Li, Hang and Naoki Abe. 1996. Clustering words with the MDL principle. In *Proceedings of the 16th International Conference on Computational Linguistics*, pages 4–9, Copenhagen, Denmark.
- Li, Hang and Naoki Abe. 1998. Generalizing case frames using a thesaurus and the MDL principle. *Computational Linguistics*, 24(2):217–244.
- McCarthy, Diana. 1997. Word sense disambiguation for acquisition of selectional preferences. In *Proceedings of the ACL/EACL Workshop on Automatic Information Extraction and Building of Lexical Semantic Resources for NLP Applications*, pages 52–61, Madrid.
- McCarthy, Diana. 2000. Using semantic preferences to identify verbal participation in role switching. In *Proceedings of the First Conference of the North American Chapter of the Association for Computational Linguistics*, pages 256–263, Seattle.
- Miller, George A. 1998. Nouns in WordNet. In Christiane Fellbaum, editor, *WordNet: An Electronic Lexical Database*. MIT Press, pages 23–46.
- Pedersen, Ted. 1996. Fishing for exactness. In *Proceedings of the South-Central SAS Users Group Conference*, Austin, pages 188–200.
- Pedersen, Ted. 2001. A decision tree of bigrams is an accurate predictor of word sense. In *Proceedings of the Second Meeting of the North American Chapter of the Association for Computational Linguistics*, pages 79–86, Pittsburgh.
- Pereira, Fernando, Naftali Tishby, and Lillian Lee. 1993. Distributional clustering of English words. In *Proceedings of the 31st Annual Meeting of the Association for Computational Linguistics*, pages 183–190, Columbus, OH.
- Resnik, Philip. 1993. *Selection and Information: A Class-Based Approach to Lexical Relationships*. Ph.D. dissertation, University of Pennsylvania.
- Resnik, Philip. 1998. WordNet and class-based probabilities. In Christiane Fellbaum, editor, *WordNet: An Electronic Lexical Database*. MIT Press, pages 239–263.
- Ribas, Francesc. 1995. On learning more appropriate selectional restrictions. In *Proceedings of the Seventh Conference of the European Chapter of the Association for Computational Linguistics*, pages 112–118, Dublin.
- Rooth, Mats, Stefan Riezler, Detlef Prescher, Glenn Carroll, and Franz Beil. 1999. Inducing a semantically annotated lexicon via EM-based clustering. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pages 104–111, University of Maryland, College Park.
- Wagner, Andreas. 2000. Enriching a lexical semantic net with selectional preferences by means of statistical corpus analysis. In *Proceedings of the ECAI-2000 Workshop on Ontology Learning*, Berlin, pages 37–42.