

# On Learning more Appropriate Selectional Restrictions

Francesc Ribas\*

Departament de Llenguatges i Sistemes Informàtics  
Universitat Politècnica de Catalunya  
Pau Gargallo, 5  
08028 Barcelona  
Spain  
ribas@lsi.upc.es

## Abstract

We present some variations affecting the association measure and thresholding on a technique for learning Selectional Restrictions from on-line corpora. It uses a wide-coverage noun taxonomy and a statistical measure to generalize the appropriate semantic classes. Evaluation measures for the Selectional Restrictions learning task are discussed. Finally, an experimental evaluation of these variations is reported.

**Subject Areas:** corpus-based language modeling, computational lexicography

## 1 Introduction

In recent years there has been a common agreement in the NLP research community on the importance of having an extensive coverage of selectional restrictions (SRs) tuned to the domain to work with. SRs can be seen as semantic type constraints that a word sense imposes on the words with which it combines in the process of semantic interpretation. SRs may have different applications in NLP, specifically, they may help a parser with Word Sense Selection (WSS, as in (Hirst, 1987)), with preferring certain structures out of several grammatical ones (Whittemore et al., 1990) and finally with deciding the semantic role played by a syntactic complement (Basili et al., 1992). Lexicography is also interested in the acquisition of SRs (both defining in context approach and lexical semantics work (Levin, 1992)).

The aim of our work is to explore the feasibility of using an statistical method for extracting SRs from on-line corpora. Resnik (1992) developed a method for automatically extracting class-based SRs from on-line corpora. Ribas (1994a)

performed some experiments using this basic technique and drew up some limitations from the corresponding results.

In this paper we will describe some substantial modifications to the basic technique and will report the corresponding experimental evaluation. The outline of the paper is as follows: in section 2 we summarize the basic methodology used in (Ribas, 1994a), analyzing its limitations; in section 3 we explore some alternative statistical measures for ranking the hypothesized SRs; in section 4 we propose some evaluation measures on the SRs-learning problem, and use them to test the experimental results obtained by the different techniques; finally, in section 5 we draw up the final conclusions and establish future lines of research.

## 2 The basic technique for learning SRs

### 2.1 Description

The technique functionality can be summarized as:

**Input** The training set, i.e. a list of complement co-occurrence triples, (*verb-lemma*, *syntactic-relationship*, *noun-lemma*) extracted from the corpus.

**Previous knowledge used** A semantic hierarchy (WordNet<sup>1</sup>) where words are clustered in semantic classes, and semantic classes are organized hierarchically. Polysemous words are represented as instances of different classes.

**Output** A set of syntactic SRs, (*verb-lemma*, *syntactic-relationship*, *semantic-class*, *weight*). The final SRs must be mutually disjoint. SRs are weighted according to the statistical evidence found in the corpus.

**Learning process** 3 stages:

1. Creation of the space of candidate classes.

\*This research has been made in the framework of the Aquilex-II Esprit Project (7315), and has been supported by a grant of Departament d'Ensenyament, Generalitat de Catalunya, 91-DOGC-1491.

<sup>1</sup>WordNet is a broad-coverage lexical database, see (Miller et al., 1991)

Acquired SR	Type	Assoc	Examples of nouns in Treebank
< <i>suit, suing</i> >	Senses	0.41	suit
< <i>suit_of_clothes</i> >	Senses	0.41	suit
< <i>suit</i> >	Senses	0.40	suit
< <i>group</i> >	↑Abs	0.35	administration, agency, bank, ...
< <i>legal_action</i> >	Ok	0.28	suit
< <i>person, individual</i> >	Ok	0.23	advocate, buyer, carrier, client, ...
< <i>radical</i> >	Senses	0.16	group
< <i>city</i> >	Senses	0.15	proper_name
< <i>admin_district</i> >	Senses	0.14	proper_name
< <i>social_control</i> >	Senses	0.11	administration, government
< <i>status</i> >	Senses	0.087	government, leadership
< <i>activity</i> >	Senses	-0.01	administration, leadership, provision
< <i>cognition</i> >	Senses	-0.04	concern, leadership, provision, science

Table 1: SRs acquired for the subject of *seek*

- Evaluation of the appropriateness of the candidates by means of a statistical measure.
- Selection of the most appropriate subset in the candidate space to convey the SRs.

The appropriateness of a class for expressing SRs (stage 2) is quantified from the strength of co-occurrence of verbs and classes of nouns in the corpus (Resnik, 1992). Given the verb  $v$ , the syntactic-relationship  $s$  and the candidate class  $c$ , the Association Score,  $Assoc$ , between  $v$  and  $c$  in  $s$  is defined:

$$\begin{aligned}
 Assoc(v, s, c) &= p(c|v, s)I(v; c|s) \\
 &= p(c|v, s) \log \frac{p(c|v, s)}{p(c|s)}
 \end{aligned}$$

The two terms of  $Assoc$  try to capture different properties:

- Mutual information ratio,  $I(v; c|s)$ , measures the strength of the statistical association between the given verb  $v$  and the candidate class  $c$  in the given syntactic position  $s$ . It compares the prior distribution,  $p(c|s)$ , with the posterior distribution,  $p(c|v, s)$ .
- $p(c|v, s)$  scales up the strength of the association by the frequency of the relationship.

Probabilities are estimated by Maximum Likelihood Estimation (MLE), i.e. counting the relative frequency of the considered events in the corpus<sup>2</sup>. However, it is not obvious how to calculate class frequencies when the training corpus is not semantically tagged as is the case. Nevertheless, we take a simplistic approach and calculate them in the following manner:

$$freq(v, s, c) = \sum_{n \in c} freq(v, s, n) \times w \quad (1)$$

<sup>2</sup>The utility of introducing smoothing techniques on class-based distributions is dubious, see (Resnik, 1993).

Where  $w$  is a constant factor used to normalize the probabilities<sup>3</sup>

$$w = \frac{\sum_{v \in V} \sum_{s \in S} \sum_{n \in N} freq(v, s, n)}{\sum_{v \in V} \sum_{s \in S} \sum_{n \in N} freq(v, s, n) |senses(n)|} \quad (2)$$

When creating the space of candidate classes (learning process, stage 1), we use a *thresholding* technique to ignore as much as possible the noise introduced in the training set. Specifically, we consider only those classes that have a higher number of occurrences than the threshold. The selection of the most appropriate classes (stage 3) is based on a global search through the candidates, in such a way that the final classes are mutually disjoint (not related by hyperonymy).

## 2.2 Evaluation

Ribas (1994a) reported experimental results obtained from the application of the above technique to learn SRs. He performed an evaluation of the SRs obtained from a training set of 870,000 words of the Wall Street Journal. In this section we summarize the results and conclusions reached in that paper.

For instance, table 1 shows the SRs acquired for the *subject* position of the verb *seek*. *Type* indicates a manual diagnosis about the class appropriateness (Ok: correct; ↑Abs: over-generalization; Senses: due to erroneous senses). *Assoc* corresponds to the association score (higher values appear first). Most of the induced classes are due to incorrect senses. Thus, although *suit* was used in the WSJ articles only in the sense of < *legal\_action* >, the algorithm not only considered the other senses as well (< *suit, suing* >, <

<sup>3</sup>Resnik (1992) and Ribas (1994a) used equation 1 without introducing normalization. Therefore, the estimated function didn't accomplish probability axioms. Nevertheless, their results should be equivalent (for our purposes) to those introducing normalization because it shouldn't affect the relative ordering of  $Assoc$  among rival candidate classes for the same  $(v, s)$ .

*suit\_of\_clothes*>, <*suit*>), but the *Assoc* score ranked them higher than the appropriate sense. We can also notice that the ↑Abs class, <*group*>, seems too general for the example nouns, while one of its daughters, <*people*> seems to fit the data much better.

Analyzing the results obtained from different experimental evaluation methods, Ribas (1994a) drew up some conclusions:

- a. The technique achieves a good coverage.
- b. Most of the classes acquired result from the accumulation of incorrect senses.
- c. No clear co-relation between *Assoc* and the manual diagnosis is found.
- d. A slight tendency to over-generalization exists due to incorrect senses.

Although the performance of the presented technique seems to be quite good, we think that some of the detected flaws could possibly be addressed. Noise due to polysemy of the nouns involved seems to be the main obstacle for the practicality of the technique. It makes the association score prefer incorrect classes and jump on over-generalizations. In this paper we are interested in exploring various ways to make the technique more robust to noise, namely, (a) to experiment with variations of the association score, (b) to experiment with thresholding.

### 3 Variations on the association statistical measure

In this section we consider different variations on the association score in order to make it more robust. The different techniques are experimentally evaluated in section 4.2.

#### 3.1 Variations on the prior probability

When considering the prior probability, the more independent of the context it is the better to measure actual associations. A sensible modification of the measure would be to consider  $p(c)$  as the prior distribution:

$$Assoc'(v, s, c) = p(c|v, s) \log \frac{p(c|v, s)}{p(c)}$$

Using the chain rule on mutual information (Cover and Thomas, 1991, p. 22) we can mathematically relate the different versions of *Assoc*,

$$Assoc'(v, s, c) = p(c|v, s) \log \frac{p(c|s)}{p(c)} + Assoc(v, s, c)$$

The first advantage of *Assoc'* would come from this (information theoretical) relationship. Specifically, the *Assoc'* takes into account the preference (selection) of syntactic positions for particular classes. In intuitive terms, typical subjects (e.g. <person, individual, ...>) would be preferred

(to atypical subjects as <*suit\_of\_clothes*>) as SRs on the subject in contrast to *Assoc*. The second advantage is that as long as the prior probabilities,  $p(c)$ , involve simpler events than those used in *Assoc*,  $p(c|s)$ , the estimation is easier and more accurate (ameliorating data sparseness).

A subsequent modification would be to estimate the prior,  $p(c)$ , from the counts of all the nouns appearing in the corpus independently of their syntactic positions (not restricted to be heads of verbal complements). In this way, the estimation of  $p(c)$  would be easier and more accurate.

#### 3.2 Estimating class probabilities from noun frequencies

In the global weighting technique presented in equation 2 very polysemous nouns provide the same amount of evidence to every sense as non-ambiguous nouns do –while less ambiguous nouns could be more informative about the correct classes as long as they do not carry ambiguity.

The weight introduced in (1) could alternatively be found in a local manner, in such a way that more polysemous nouns would give less evidence to each one of their senses than less ambiguous ones. Local weight could be obtained using  $p(c|n)$ . Nevertheless, a good estimation of this probability seems quite problematic because of the lack of tagged training material. In absence of a better estimator we use a rather poor one as the uniform distribution,

$$w(n, c) = \tilde{p}(c|n) = \frac{|senses(n) \in c|}{|senses(n)|}$$

Resnik (1993) also uses a local normalization technique but he normalizes by the total number of classes in the hierarchy. This scheme seems to present two problematic features (see (Ribas, 1994b) for more details). First, it doesn't take dependency relationships introduced by hyperonymy into account. Second, nouns categorized in lower levels in the taxonomy provide less weight to each class than higher nouns.

#### 3.3 Other statistical measures to score SRs

In this section we propose the application of other measures apart from *Assoc* for learning SRs: log-likelihood ratio (Dunning, 1993), relative entropy (Cover and Thomas, 1991), mutual information ratio (Church and Hanks, 1990),  $\phi^2$  (Gale and Church, 1991). In section (4) their experimental evaluation is presented.

The statistical measures used to detect associations on the distribution defined by two random variables X and Y work by measuring the deviation of the conditional distribution,  $P(X|Y)$ , from the expected distribution if both variables were considered independent, i.e. the marginal distribution,  $P(X)$ . If  $P(X)$  is a good approximation

	$c$	$\neg c$
$v\_s$	$p(c v\_s)$	$p(\neg c v\_s)$
$\neg v\_s$	$p(c \neg v\_s)$	$p(\neg c \neg v\_s)$
	$p(c)$	$p(\neg c)$

Table 2: Conditional and marginal distributions

of  $P(X|Y)$ , association measures should be low (near zero), otherwise deviating significantly from zero.

Table 2 shows the cross-table formed by the conditional and marginal distributions in the case of  $X = \{c, \neg c\}$  and  $Y = \{v\_s, \neg v\_s\}$ . Different association measures use the information provided in the cross-table to different extents. Thus, *Assoc* and mutual information ratio consider only the deviation of the conditional probability  $p(c|v, s)$  from the corresponding marginal,  $p(c)$ .

On the other hand, *log-likelihood* ratio and  $\phi^2$  measure the association between  $v\_s$  and  $c$  considering the deviation of the four conditional cells in table 2 from the corresponding marginals. It is plausible that the deviation of the cells not taken into account by *Assoc* can help on extracting useful SRs.

Finally, it would be interesting to only use the information related to the selectional behavior of  $v\_s$ , i.e. comparing the conditional probabilities of  $c$  and  $\neg c$  given  $v\_s$  with the corresponding marginals. Relative entropy,  $D(P(X|v\_s)||P(X))$ , could do this job.

## 4 Evaluation

### 4.1 Evaluation methods of SRs

Evaluation on NLP has been crucial to fostering research in particular areas. Evaluation of the SR learning task would provide grounds to compare different techniques that try to abstract SRs from corpus using WordNet (e.g. section 4.2). It would also permit measuring the utility of the SRs obtained using WordNet in comparison with other frameworks using other kinds of knowledge. Finally it would be a powerful tool for detecting flaws of a particular technique (e.g. (Ribas, 1994a) analysis).

However, a related and crucial issue is which linguistic tasks are used as a reference. SRs are useful for both lexicography and NLP. On the one hand, from the point of view of lexicography, the goal of evaluation would be to measure the quality of the SRs induced, (i.e., how well the resulting classes correspond to the nouns as they were used in the corpus). On the other hand, from the point of view of NLP, SRs should be evaluated on their utility (i.e., how much they help on performing the reference task).

#### 4.1.1 Lexicography-oriented evaluation

As far as lexicography (quality) is concerned, we think the main criteria SRs acquired from corpora should meet are: (a) correct categorization –inferred classes should correspond to the correct senses of the words that are being generalized–, (b) appropriate generalization level and (c) good coverage –the majority of the noun occurrences in the corpus should be successfully generalized by the induced SRs.

Some of the methods we could use for assessing experimentally the accomplishment of these criteria would be:

- **Introspection** A lexicographer checks if the SRs accomplish the criteria (a) and (b) above (e.g., the manual diagnosis in table 1). Besides the intrinsic difficulties of this approach, it does not seem appropriate when comparing across different techniques for learning SRs, because of its qualitative flavor.
- **Quantification of generalization level appropriateness** A possible measure would be the percentage of sense occurrences included in the induced SRs which are effectively correct (from now on called *Abstraction Ratio*). Hopefully, a technique with a higher abstraction ratio learns classes that fit the set of examples better. A manual assessment of the ratio confirmed this behavior, as testing sets with a lower ratio seemed to be inducing less  $\uparrow$ Abs cases.
- **Quantification of coverage** It could be measured as the proportion of triples whose correct sense belongs to one of the SRs.

#### 4.1.2 NLP evaluation tasks

The NLP tasks where SRs utility could be evaluated are diverse. Some of them have already been introduced in section 1. In the recent literature ((Grishman and Sterling, 1992), (Resnik, 1993), ...) several task oriented schemes to test Selectional Restrictions (mainly on syntactic ambiguity resolution) have been proposed. However, we have tested SRs on a WSS task, using the following scheme. For every triple in the testing set the algorithm selects as most appropriate that noun-sense that has as hyperonym the SR class with highest association score. When more than one sense belongs to the highest SR, a random selection is performed. When no SR has been acquired, the algorithm remains undecided. The results of this WSS procedure are checked against a testing-sample manually analyzed, and precision and recall ratios are calculated. Precision is calculated as the ratio of manual-automatic matches / number of noun occurrences disambiguated by the procedure. Recall is computed as the ratio of manual-automatic matches / total number of noun occurrences.

Technique	Coverage (%)
<i>Assoc &amp; All nouns</i>	95.7
<i>Assoc &amp; p(c s)</i>	95.5
<i>Assoc &amp; Head-nouns</i>	95.3
<i>D</i>	93.7
<i>log - likelihood</i>	92.9
<i>Assoc &amp; Normalizing</i>	92.7
$\phi^2$	88.2
<i>I</i>	74.1

Table 3: Coverage Ratio

## 4.2 Experimental results

In order to evaluate the different variants on the association score and the impact of thresholding we performed several experiments. In this section we analyze the results. As training set we used the 870,000 words of WSJ material provided in the ACL/DCI version of the Penn Treebank. The testing set consisted of 2,658 triples corresponding to four average common verbs in the Treebank: *rise*, *report*, *seek* and *present*. We only considered those triples that had been correctly extracted from the Treebank and whose noun had the correct sense included in WordNet (2,165 triples out of the 2,658, from now on, called the *testing-sample*).

As evaluation measures we used coverage, abstraction ratio, and recall and precision ratios on the WSS task (section 4.1). In addition we performed some evaluation by hand comparing the SRs acquired by the different techniques.

### 4.2.1 Comparing different techniques

Coverage for the different techniques is shown in table 3. The higher the coverage, the better the technique succeeds in correctly generalizing more of the input examples. The labels used for referring to the different techniques are as follows: “*Assoc & p(c|s)*” corresponds to the basic association measure (section 2), “*Assoc & Head-nouns*” and “*Assoc & All nouns*” to the techniques introduced in section 3.1, “*Assoc & Normalizing*” to the local normalization (section 3.2), and finally, *log-likelihood*, *D* (relative entropy) and *I* (mutual information ratio) to the techniques discussed in section 3.3.

The abstraction ratio for the different techniques is shown in table 4. In principle, the higher abstraction ratio, the better the technique succeeds in filtering out incorrect senses (less  $\uparrow$ Abs).

The precision and recall ratios on the noun WSS task for the different techniques are shown in table 5. In principle, the higher the precision and recall ratios the better the technique succeeds in inducing appropriate SRs for the disambiguation task.

As far as the evaluation measures try to account for different phenomena the goodness of a particular technique should be quantified as a trade-off.

Technique	Abs Ratio (%)
<i>I</i>	66.6
<i>log - likelihood</i>	64.6
$\phi^2$	64.4
<i>Assoc &amp; All nouns</i>	64.3
<i>Assoc &amp; Head-nouns</i>	63.9
<i>Assoc &amp; p(c s)</i>	63
<i>D</i>	62.3
<i>Assoc &amp; Normalizing</i>	58.5

Table 4: Abstraction Ratio

Technique	Prec. (%)	Rec. (%)
<i>Assoc &amp; All nouns</i>	80.3	78.5
<i>Assoc &amp; p(c s)</i>	79.9	77.9
<i>Assoc &amp; Head-nouns</i>	78.5	76.7
<i>log - likelihood</i>	77.2	74.4
<i>D</i>	75.9	74.1
<i>Assoc &amp; Normalizing</i>	75.9	73.3
$\phi^2$	67.8	63
<i>I</i>	50.4	45.7
Guessing Heuristic	62.7	62.7

Table 5: Precision and Recall on the WSS task

Most of the results are very similar (differences are not statistically significant). Therefore we should be cautious when extrapolating the results. Some of the conclusions from the tables above are:

1.  $\phi^2$  and *I* get sensibly worse results than other measures (although abstraction is quite good).
2. The local normalizing technique using the uniform distribution does not help. It seems that by using the local weighting we misinform the algorithm. The problem is the reduced weight that polysemous nouns get, while they seem to be the most informative<sup>4</sup>. However, a better informed kind of local weight (section 5) should improve the technique significantly.
3. All versions of *Assoc* (except the local normalization) get good results. Specially the two techniques that exploit a simpler prior distribution, which seem to improve the basic technique.
4. *log-likelihood* and *D* seem to get slightly worse results than *Assoc* techniques, although the results are very similar.

### 4.2.2 Thresholding

We were also interested in measuring the impact of thresholding on the SRs acquired. In figure 1 we can see the different evaluation measures of the basic technique when varying the threshold. Precision and recall coincide when no candidate

<sup>4</sup>In some way, it conforms to Zipf-law (Zipf, 1945): noun frequency and polysemy are correlated.

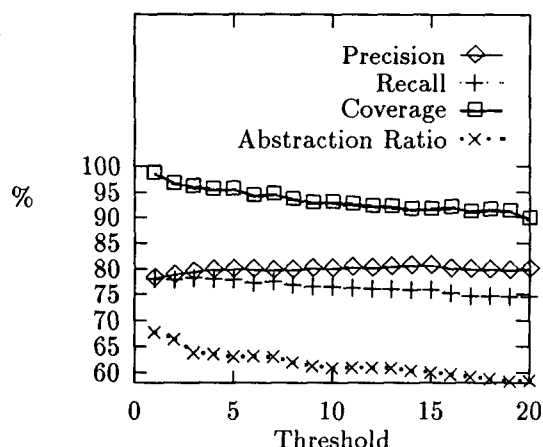


Figure 1: Assoc: Evaluation ratios vs. Threshold

classes are refused ( $threshold = 1$ ). However, as it might be expected, as the threshold increases (i.e. some cases are not classified) the two ratios slightly diverge (precision increases and recall diminishes).

Figure 1 also shows the impact of thresholding on coverage and abstraction ratios. Both decrease when threshold increases, probably because when the rejecting threshold is low, small classes that fit the data well can be induced, learning over-general or incomplete SRs otherwise.

Finally, it seems that precision and abstraction ratios are in inverse co-relation (as precision grows, abstraction decreases). In terms of WSS, general classes may be performing better than classes that fit the data better. Nevertheless, this relationship should be further explored in future work.

## 5 Conclusions and future work

In this paper we have presented some variations affecting the association measure and thresholding on the basic technique for learning SRs from on-line corpora. We proposed some evaluation measures for the SRs learning task. Finally, experimental results on these variations were reported. We can conclude that some of these variations seem to improve the results obtained using the basic technique. However, although the technique still seems far from practical application to NLP tasks, it may be most useful for providing experimental insight to lexicographers. Future lines of research will mainly concentrate on improving the local normalization technique by solving the noun sense ambiguity. We have foreseen the application of the following techniques:

- Simple techniques to decide the best sense  $c$  given the target noun  $n$  using estimates of the n-grams:  $P(c)$ ,  $P(c|n)$ ,  $P(c|v, s)$  and  $P(c|v, s, n)$ , obtained from supervised and

un-supervised corpora.

- Combining the different n-grams by means of smoothing techniques.
- Calculating  $P(c|v, s, n)$  combining  $P(n|c)$  and  $P(c|v, s)$ , and applying the EM Algorithm (Dempster et al., 1977) to improve the model.
- Using the WordNet hierarchy as a source of backing-off knowledge, in such a way that if n-grams composed by  $c$  aren't enough to decide the best sense (are equal to zero), the tri-grams of ancestor classes could be used instead.

## References

- R. Basili, M.T. Pazienza, and P. Velardi. 1992. Computational lexicons: the neat examples and the odd exemplars. In *Procs 3rd ANLP*, Trento, Italy, April.
- K.W. Church and P. Hanks. 1990. Word association norms, mutual information and lexicography. *Computational Linguistics*, 16(1).
- T.M. Cover and J.A. Thomas, editors. 1991. *Elements of Information Theory*. John Wiley.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society*, 39(B):1-38.
- T. Dunning. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61-74.
- W. Gale and K. W. Church. 1991. Identifying word correspondences in parallel texts. In *DARPA Speech and Natural Language Workshop*, Pacific Grove, California, February.
- R. Grishman and J. Sterling. 1992. Acquisition of selectional patterns. In *COLING*, Nantes, France, march.
- G. Hirst. 1987. *Semantic interpretation and the resolution of ambiguity*. Cambridge University Press.
- B. Levin. 1992. *Towards a lexical organization of English verbs*. University of Chicago Press.
- G. Miller, R. Beckwith, C. Fellbaum, D. Gross, and K. Miller. 1991. Five papers on wordnet. *International Journal of Lexicography*.
- P. S. Resnik. 1992. Wordnet and distributional analysis: A class-based approach to lexical discovery. In *AAAI Symposium on Probabilistic Approaches to NL*, San Jose, CA.
- P. S. Resnik. 1993. *Selection and Information: A Class-Based Approach to lexical relationships*. Ph.D. thesis, Computer and Information Science Department, University of Pennsylvania.

- F. Ribas. 1994a. An experiment on learning appropriate selectional restrictions from a parsed corpus. In *COLING*, Kyoto, Japan, August.
- F. Ribas. 1994b. Learning more appropriate selectional restrictions. Technical report, ESPRIT BRA-7315 ACQUILEX-II WP.
- G. Whittmore, K. Ferrara, and H. Brunner. 1990. Empirical study of predictive powers of simple attachment schemes for post-modifier prepositional phrases. In *Procs. ACL*, Pennsylvania.
- G. K. Zipf. 1945. The meaning-frequency relationship of words. *The Journal of General Psychology*, 33:251-256.

(Acquilex-II Working Papers can be obtained by sending a request to [cide@cup.cam.uk](mailto:cide@cup.cam.uk))

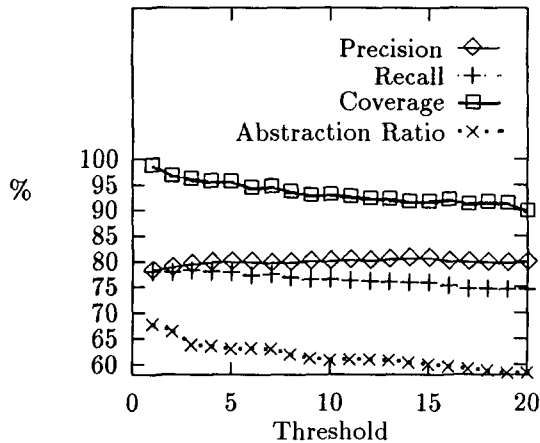


Figure 1: *Assoc*: Evaluation ratios vs. Threshold

classes are refused ( $threshold = 1$ ). However, as it might be expected, as the threshold increases (i.e. some cases are not classified) the two ratios slightly diverge (precision increases and recall diminishes).

Figure 1 also shows the impact of thresholding on coverage and abstraction ratios. Both decrease when threshold increases, probably because when the rejecting threshold is low, small classes that fit the data well can be induced, learning over-general or incomplete SRs otherwise.

Finally, it seems that precision and abstraction ratios are in inverse co-relation (as precision grows, abstraction decreases). In terms of WSS, general classes may be performing better than classes that fit the data better. Nevertheless, this relationship should be further explored in future work.

## 5 Conclusions and future work

In this paper we have presented some variations affecting the association measure and thresholding on the basic technique for learning SRs from on-line corpora. We proposed some evaluation measures for the SRs learning task. Finally, experimental results on these variations were reported. We can conclude that some of these variations seem to improve the results obtained using the basic technique. However, although the technique still seems far from practical application to NLP tasks, it may be most useful for providing experimental insight to lexicographers. Future lines of research will mainly concentrate on improving the local normalization technique by solving the noun sense ambiguity. We have foreseen the application of the following techniques:

- Simple techniques to decide the best sense  $c$  given the target noun  $n$  using estimates of the n-grams:  $P(c)$ ,  $P(c|n)$ ,  $P(c|v, s)$  and  $P(c|v, s, n)$ , obtained from supervised and

un-supervised corpora.

- Combining the different n-grams by means of smoothing techniques.
- Calculating  $P(c|v, s, n)$  combining  $P(n|c)$  and  $P(c|v, s)$ , and applying the EM Algorithm (Dempster et al., 1977) to improve the model.
- Using the WordNet hierarchy as a source of backing-off knowledge, in such a way that if n-grams composed by  $c$  aren't enough to decide the best sense (are equal to zero), the tri-grams of ancestor classes could be used instead.

## References

- R. Basili, M.T. Pazienza, and P. Velardi. 1992. Computational lexicons: the neat examples and the odd exemplars. In *Procs 3rd ANLP*, Trento, Italy, April.
- K.W. Church and P. Hanks. 1990. Word association norms, mutual information and lexicography. *Computational Linguistics*, 16(1).
- T.M. Cover and J.A. Thomas, editors. 1991. *Elements of Information Theory*. John Wiley.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society*, 39(B):1-38.
- T. Dunning. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61-74.
- W. Gale and K. W. Church. 1991. Identifying word correspondences in parallel texts. In *DARPA Speech and Natural Language Workshop*, Pacific Grove, California, February.
- R. Grishman and J. Sterling. 1992. Acquisition of selectional patterns. In *COLING*, Nantes, France, march.
- G. Hirst. 1987. *Semantic interpretation and the resolution of ambiguity*. Cambridge University Press.
- B. Levin. 1992. *Towards a lexical organization of English verbs*. University of Chicago Press.
- G. Miller, R. Beckwith, C. Fellbaum, D. Gross, and K. Miller. 1991. Five papers on wordnet. *International Journal of Lexicography*.
- P. S. Resnik. 1992. Wordnet and distributional analysis: A class-based approach to lexical discovery. In *AAAI Symposium on Probabilistic Approaches to NL*, San Jose, CA.
- P. S. Resnik. 1993. *Selection and Information: A Class-Based Approach to lexical relationships*. Ph.D. thesis, Computer and Information Science Department, University of Pennsylvania.



- F. Ribas. 1994a. An experiment on learning appropriate selectional restrictions from a parsed corpus. In *COLING*, Kyoto, Japan, August.
- F. Ribas. 1994b. Learning more appropriate selectional restrictions. Technical report, ESPRIT BRA-7315 ACQUILEX-II WP.
- G. Whittemore, K. Ferrara, and H. Brunner. 1990. Empirical study of predictive powers of simple attachment schemes for post-modifier prepositional phrases. In *Procs. ACL*, Pennsylvania.
- G. K. Zipf. 1945. The meaning-frequency relationship of words. *The Journal of General Psychology*, 33:251-256.

(Acquilex-II Working Papers can be obtained by sending a request to [cide@cup.cam.uk](mailto:cide@cup.cam.uk))