

Parsing and Subcategorization Data

Jianguo Li

Department of Linguistics

The Ohio State University

Columbus, OH, USA

jianguo@ling.ohio-state.edu

Abstract

In this paper, we compare the performance of a state-of-the-art statistical parser (Bikel, 2004) in parsing written and spoken language and in generating subcategorization cues from written and spoken language. Although Bikel's parser achieves a higher accuracy for parsing written language, it achieves a higher accuracy when extracting subcategorization cues from spoken language. Additionally, we explore the utility of punctuation in helping parsing and extraction of subcategorization cues. Our experiments show that punctuation is of little help in parsing spoken language and extracting subcategorization cues from spoken language. This indicates that there is no need to add punctuation in transcribing spoken corpora simply in order to help parsers.

1 Introduction

Robust statistical syntactic parsers, made possible by new statistical techniques (Collins, 1999; Charniak, 2000; Bikel, 2004) and by the availability of large, hand-annotated training corpora such as WSJ (Marcus et al., 1993) and Switchboard (Godefroy et al., 1992), have had a major impact on the field of natural language processing. There are many ways to make use of parsers' output. One particular form of data that can be extracted from parses is information about subcategorization. Subcategorization data comes in two forms: subcategorization frame (SCF) and subcategorization cue (SCC). SCFs differ from SCCs in that SCFs contain only arguments while SCCs contain both arguments and adjuncts. Both SCFs

and SCCs have been crucial to NLP tasks. For example, SCFs have been used for verb disambiguation and classification (Schulte im Walde, 2000; Merlo and Stevenson, 2001; Lapata and Brew, 2004; Merlo et al., 2005) and SCCs for semantic role labeling (Xue and Palmer, 2004; Punyakanok et al., 2005).

Current technology for automatically acquiring subcategorization data from corpora usually relies on statistical parsers to generate SCCs. While great efforts have been made in parsing written texts and extracting subcategorization data from written texts, spoken corpora have received little attention. This is understandable given that spoken language poses several challenges that are absent in written texts, including disfluency, uncertainty about utterance segmentation and lack of punctuation. Roland and Jurafsky (1998) have suggested that there are substantial subcategorization differences between written corpora and spoken corpora. For example, while written corpora show a much higher percentage of passive structures, spoken corpora usually have a higher percentage of zero-anaphora constructions. We believe that subcategorization data derived from spoken language, if of acceptable quality, would be of more value to NLP tasks involving a syntactic analysis of spoken language, but we do not pursue it here.

The goals of this study are as follows:

1. Test the performance of Bikel's parser in parsing written and spoken language.
2. Compare the accuracy level of SCCs generated from parsed written and spoken language. We hope that such a comparison will shed some light on the feasibility of acquiring SCFs from spoken language using the cur-

rent SCF acquisition technology initially designed for written language.

3. Explore the utility of punctuation¹ in parsing and extraction of SCCs. It is generally recognized that punctuation helps in parsing written texts. For example, Roark (2001) finds that removing punctuation from both training and test data (WSJ) decreases his parser’s accuracy from 86.4%/86.8% (LR/LP) to 83.4%/84.1%. However, spoken language does not come with punctuation. Even when punctuation is added in the process of transcription, its utility in helping parsing is slight. Both Roark (2001) and Engel et al. (2002) report that removing punctuation from both training and test data (Switchboard) results in only 1% decrease in their parser’s accuracy.

2 Experiment Design

Three models will be investigated for parsing and extracting SCCs from the parser’s output:

1. **punc**: leaving punctuation in both training and test data.
2. **no-punc**: removing punctuation from both training and test data.
3. **punc-no-punc**: removing punctuation from only test data.

Following the convention in the parsing community, for written language, we selected sections 02-21 of WSJ as training data and section 23 as test data (Collins, 1999). For spoken language, we designated section 2 and 3 of Switchboard as training data and files of sw4004 to sw4135 of section 4 as test data (Roark, 2001). Since we are also interested in extracting SCCs from the parser’s output, we eliminated from the two test corpora all sentences that do not contain verbs. Our experiments proceed in the following three steps:

1. Tag test data using the POS-tagger described in Ratnaparkhi (1996).
2. Parse the POS-tagged data using Bikel’s parser.

¹We use punctuation to refer to sentence-internal punctuation unless otherwise specified.

| label | clause type | desired SCCs |
|-------|--------------------------|---|
| S | gerundive | (NP)-GERUND |
| | small clause | NP-NP, (NP)-ADJP |
| | control | (NP)-INF- <i>to</i> |
| SBAR | control | (NP)-INF- <i>wh-to</i> |
| | with a complementizer | (NP)-S- <i>wh</i> , (NP)-S- <i>that</i> |
| | without a complementizer | (NP)-S- <i>that</i> |

Table 1: SCCs for different clauses

3. Extract SCCs from the parser’s output. The extractor we built first locates each verb in the parser’s output and then identifies the syntactic categories of all its sisters and combines them into an SCC. However, there are cases where the extractor has more work to do.

- **Finite and Infinite Clauses**: In the Penn Treebank, **S** and **SBAR** are used to label different types of clauses, obscuring too much detail about the internal structure of each clause. Our extractor is designed to identify the internal structure of different types of clause, as shown in Table 1.
- **Passive Structures**: As noted above, Roland and Jurafsky (Roland and Jurafsky, 1998) have noticed that written language tends to have a much higher percentage of passive structures than spoken language. Our extractor is also designed to identify passive structures from the parser’s output.

3 Experiment Results

3.1 Parsing and SCCs

We used EVALB measures Labeled Recall (LR) and Labeled Precision (LP) to compare the parsing performance of different models. To compare the accuracy of SCCs proposed from the parser’s output, we calculated SCC Recall (SR) and SCC Precision (SP). SR and SP are defined as follows:

$$SR = \frac{\text{number of correct cues from the parser's output}}{\text{number of cues from treebank parse}} \quad (1)$$

$$SP = \frac{\text{number of correct cues from the parser's output}}{\text{number of cues from the parser's output}} \quad (2)$$

$$\text{SCC Balanced F-measure} = \frac{2 * SR * SP}{SR + SP} \quad (3)$$

The results for parsing WSJ and Switchboard and extracting SCCs are summarized in Table 2.

The LR/LP figures show the following trends:

| WSJ | | |
|--------------|---------------|---------------|
| model | LR/LP | SR/SP |
| punc | 87.92%/88.29% | 76.93%/77.70% |
| no-punc | 86.25%/86.91% | 76.96%/76.47% |
| punc-no-punc | 82.31%/83.70% | 74.62%/74.88% |
| Switchboard | | |
| model | LR/LP | SR/SP |
| punc | 83.14%/83.80% | 79.04%/78.62% |
| no-punc | 82.42%/83.74% | 78.81%/78.37% |
| punc-no-punc | 78.62%/80.68% | 75.51%/75.02% |

Table 2: Results of parsing and extraction of SCCs

1. Roark (2001) showed LR/LP of 86.4%/86.8% for punctuated written language, 83.4%/84.1% for unpunctuated written language. We achieve a higher accuracy in both punctuated and unpunctuated written language, and the decrease if punctuation is removed is less
2. For spoken language, Roark (2001) showed LR/LP of 85.2%/85.6% for punctuated spoken language, 84.0%/84.6% for unpunctuated spoken language. We achieve a lower accuracy in both punctuated and unpunctuated spoken language, and the decrease if punctuation is removed is less. The trends in (1) and (2) may be due to parser differences, or to the removal of sentences lacking verbs.
3. Unsurprisingly, if the test data is unpunctuated, but the models have been trained on punctuated language, performance decreases sharply.

In terms of the accuracy of extraction of SCCs, the results follow a similar pattern. However, the utility of punctuation turns out to be even smaller. Removing punctuation from both training and test data results in a less than 0.3% drop in the accuracy of SCC extraction.

Figure 1 exhibits the relation between the accuracy of parsing and that of extracting SCCs. If we consider WSJ and Switchboard individually, there seems to exist a positive correlation between the accuracy of parsing and that of extracting SCCs. In other words, higher LR/LP indicates higher SR/SP. However, Figure 1 also shows that although the parser achieves a higher F-measure value for parsing WSJ, it achieves a higher F-measure value when generating SCCs from Switchboard.

The fact that the parser achieves a higher accuracy for extracting SCCs from Switchboard than WSJ merits further discussion. Intuitively, it

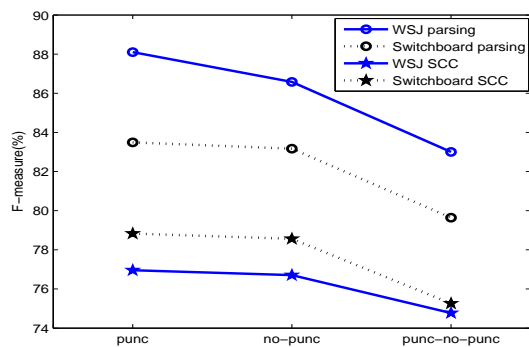


Figure 1: F-measure for parsing and extraction of SCCs

seems to be true that the shorter an SCC is, the more likely that the parser is to get it right. This intuition is confirmed by the data shown in Figure 2. Figure 2 plots the accuracy level of extracting SCCs by SCC's length. It is clear from Figure 2 that as SCCs get longer, the F-measure value drops progressively for both WSJ and Switchboard. Again, Roland and Jurafsky (1998) have suggested that one major subcategorization difference between written and spoken corpora is that spoken corpora have a much higher percentage of the zero-anaphora construction. We then examined the distribution of SCCs of different length in WSJ and Switchboard. Figure 3 shows that SCCs of length 0² account for a much higher percentage in Switchboard than WSJ, but it is always the other way around for SCCs of non-zero length. This observation led us to believe that the better performance that Bikel's parser achieves in extracting SCCs from Switchboard may be attributed to the following two factors:

1. Switchboard has a much higher percentage of SCCs of length 0.
2. The parser is very accurate in extracting shorter SCCs.

3.2 Extraction of Dependents

In order to estimate the effects of SCCs of length 0, we examined the parser's performance in retrieving dependents of verbs. Every constituent (whether an argument or adjunct) in an SCC generated by the parser is considered a dependent of

²Verbs have a length-0 SCC if they are intransitive and have no modifiers.

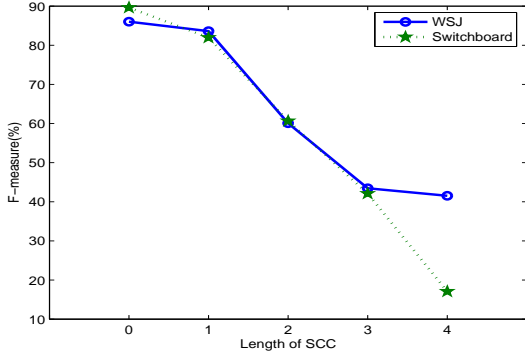


Figure 2: F-measure for SCCs of different length

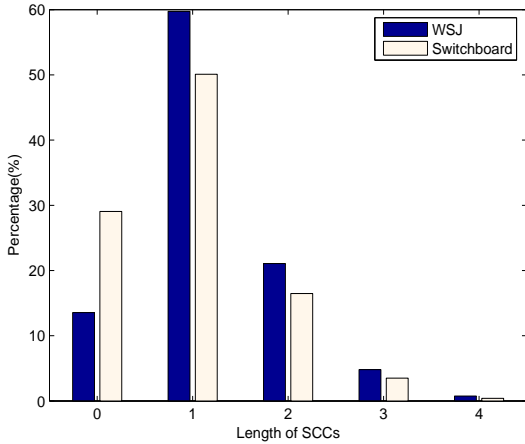


Figure 3: Distribution of SCCs by length

that verb. SCCs of length 0 will be discounted because verbs that do not take any arguments or adjuncts have no dependents³. In addition, this way of evaluating the extraction of SCCs also matches the practice in some NLP tasks such as semantic role labeling (Xue and Palmer, 2004). For the task of semantic role labeling, the total number of dependents correctly retrieved from the parser’s output affects the accuracy level of the task.

To do this, we calculated the number of dependents shared by between each SCC proposed from the parser’s output and its corresponding SCC proposed from Penn Treebank. We based our calculation on a modified version of Minimum Edit Distance Algorithm. Our algorithm works by creating a shared-dependents matrix with one column for each constituent in the target sequence (SCCs proposed from Penn Treebank) and one

³We are aware that subjects are typically also considered dependents, but we did not include subjects in our experiments

$$\text{shared-dependents}[i,j] = \text{MAX}(\text{shared-dependents}[i-1,j], \text{shared-dependents}[i-1,j-1]+1 \text{ if target}[i] = \text{source}[j], \text{shared-dependents}[i-1,j-1] \text{ if target}[i] \neq \text{source}[j], \text{shared-dependents}[i,j-1])$$

Table 3: The algorithm for computing shared dependents

| | | | | | |
|-------|----|--------|-------|-----|----------|
| INF | #5 | 1 | 1 | 2 | 3 |
| ADVP | #4 | 1 | 1 | 2 | 2 |
| PP-in | #3 | 1 | 1 | 2 | 2 |
| NP | #2 | 1 | 1 | 1 | 1 |
| NP | #1 | 1 | 1 | 1 | 1 |
| | #0 | #1 | #2 | #3 | #4 |
| | NP | S-that | PP-in | INF | |

Table 4: An example of computing the number of shared dependents

row for each constituent in the source sequence (SCCs proposed from the parser’s output). Each cell shared-dependents[i,j] contains the number of constituents shared between the first i constituents of the target sequence and the first j constituents of the source sequence. Each cell can then be computed as a simple function of the three possible paths through the matrix that arrive there. The algorithm is illustrated in Table 3.

Table 4 shows an example of how the algorithm works with NP-S-that-PP-in-INF as the target sequence and NP-NP-PP-in-ADVP-INF as the source sequence. The algorithm returns 3 as the number of dependents shared by two SCCs.

We compared the performance of Bikel’s parser in retrieving dependents from written and spoken language over all three models using Dependency Recall (DR) and Dependency Precision (DP). These metrics are defined as follows:

$$DR = \frac{\text{number of correct dependents from parser's output}}{\text{number of dependents from treebank parse}} \quad (4)$$

$$DP = \frac{\text{number of correct dependents from parser's output}}{\text{number of dependents from parser's output}} \quad (5)$$

$$\text{Dependency F-measure} = \frac{2 * DR * DP}{DR + DP} \quad (6)$$

The results of Bikel’s parser in retrieving dependents are summarized in Figure 4. Overall, the parser achieves a better performance for WSJ over all three models, just the opposite of what have been observed for SCC extraction. Interestingly, removing punctuation from both the training and test data actually slightly improves the F-measure.

This holds true for both WSJ and Switchboard. This Dependency F-measure differs in detail from similar measures in (Xue and Palmer, 2004). For present purposes all that matters is the relative value for WSJ and Switchboard.

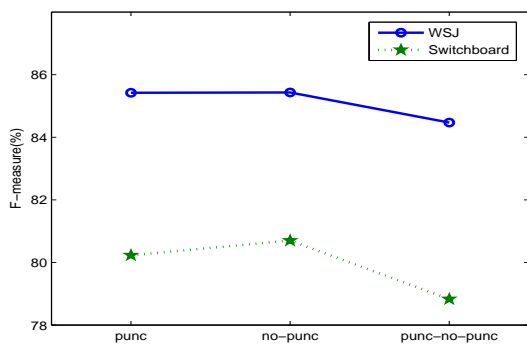


Figure 4: F-measure for extracting dependents

4 Conclusions and Future Work

4.1 Use of Parser's Output

In this paper, we have shown that it is not necessarily true that statistical parsers always perform worse when dealing with spoken language. The conventional accuracy metrics for parsing (LR/LP) should not be taken as the only metrics in determining the feasibility of applying statistical parsers to spoken language. It is necessary to consider what information we want to extract out of parsers' output and make use of.

1. Extraction of SCFs from Corpora: This task usually proceeds in two stages: (i) Use statistical parsers to generate SCCs. (ii) Apply some statistical tests such as the Binomial Hypothesis Test (Brent, 1993) and log-likelihood ratio score (Dunning, 1993) to SCCs to filter out false SCCs on the basis of their reliability and likelihood. Our experiments show that the SCCs generated for spoken language are as accurate as those generated for written language, which suggests that it is feasible to apply the current technology for automatically extracting SCFs from corpora to spoken language.
2. Semantic Role Labeling: This task usually operates on parsers' output and the number of dependents of each verb that are correctly retrieved by the parser clearly affects the accuracy of the task. Our experiments show

that the parser achieves a much lower accuracy in retrieving dependents from the spoken language than written language. This seems to suggest that a lower accuracy is likely to be achieved for a semantic role labeling task performed on spoken language. We are not aware that this has yet been tried.

4.2 Punctuation and Speech Transcription Practice

Both our experiments and Roark's experiments show that parsing accuracy measured by LR/LP experiences a sharper decrease for WSJ than Switchboard after we removed punctuation from training and test data. In spoken language, commas are largely used to delimit disfluency elements. As noted in Engel et al. (2002), statistical parsers usually condition the probability of a constituent on the types of its neighboring constituents. The way that commas are used in speech transcription seems to have the effect of increasing the range of neighboring constituents, thus fragmenting the data and making it less reliable. On the other hand, in written texts, commas serve as more reliable cues for parsers to identify phrasal and clausal boundaries.

In addition, our experiment demonstrates that punctuation does not help much with extraction of SCCs from spoken language. Removing punctuation from both the training and test data results in a less than 0.3% decrease in SR/SP. Furthermore, removing punctuation from both the training and test data actually slightly improves the performance of Bikel's parser in retrieving dependents from spoken language. All these results seem to suggest that adding punctuation in speech transcription is of little help to statistical parsers including at least three state-of-the-art statistical parsers (Collins, 1999; Charniak, 2000; Bikel, 2004). As a result, there may be other good reasons why someone who wants to build a Switchboard-like corpus should choose to provide punctuation, but there is no need to do so simply in order to help parsers.

However, segmenting utterances into individual units is necessary because statistical parsers require sentence boundaries to be clearly delimited. Current statistical parsers are unable to handle an input string consisting of two sentences. For example, when presented with an input string as in (1) and (2), if the two sentences are separated by a period (1), Bikel's parser wrongly treats the second sentence as a sentential complement of the

main verb *like* in the first sentence. As a result, the extractor generates an SCC NP-S for *like*, which is incorrect. The parser returns the same parse after we removed the period (2) and let the parser parse it again.

(1) I like the long hair. It was back in high school.

(2) I like the long hair It was back in high school.

Hence, while adding punctuation in transcribing a Switchboard-like corpus is not of much help to statistical parsers, segmenting utterances into individual units is crucial for statistical parsers. In future work, we plan to develop a system capable of automatically segmenting speech utterances into individual units.

5 Acknowledgments

This study was supported by NSF grant 0347799. Our thanks go to Chris Brew, Eric Fosler-Lussier, Mike White and three anonymous reviewers for their valuable comments.

References

- D. Bikel. 2004. Intricacies of Collin’s parsing models. *Computational Linguistics*, 30(2):479–511.
- M. Brent. 1993. From grammar to lexicon: Unsupervised learning of lexical syntax. *Computational Linguistics*, 19(3):243–262.
- E. Charniak. 2000. A maximum-entropy-inspired parser. In *Proceedings of the 2000 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 132–139.
- M. Collins. 1999. *Head-driven statistical models for natural language parsing*. Ph.D. thesis, University of Pennsylvania.
- T. Dunning. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61–74.
- D. Engel, E. Charniak, and M. Johnson. 2002. Parsing and disfluency placement. In *Proceedings of 2002 Conference on Empirical Methods of Natural Language Processing*, pages 49–54.
- J. Godefroy, E. Holliman, and J. McDaniel. 1992. SWITCHBOARD: Telephone speech corpus for research and development. In *Proceedings of ICASSP-92*, pages 517–520.
- M. Lapata and C. Brew. 2004. Verb class disambiguation using informative priors. *Computational Linguistics*, 30(1):45–73.
- M. Marcus, G. Kim, and M. Marcinkiewicz. 1993. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- P. Merlo and S. Stevenson. 2001. Automatic verb classification based on statistical distribution of argument structure. *Computational Linguistics*, 27(3):373–408.
- P. Merlo, E. Joanis, and J. Henderson. 2005. Unsupervised verb class disambiguation based on diathesis alternations. manuscripts.
- V. Punyakanok, D. Roth, and W. Yih. 2005. The necessity of syntactic parsing for semantic role labeling. In *Proceedings of the 2nd Midwest Computational Linguistics Colloquium*, pages 15–22.
- A. Ratnaparkhi. 1996. A maximum entropy model for part-of-speech tagging. In *Proceedings of the Conference on Empirical Methods of Natural Language Processing*, pages 133–142.
- B. Roark. 2001. *Robust Probabilistic Predictive Processing: Motivation, Models, and Applications*. Ph.D. thesis, Brown University.
- D. Roland and D. Jurafsky. 1998. How verb subcategorization frequency is affected by the corpus choice. In *Proceedings of the 17th International Conference on Computational Linguistics*, pages 1122–1128.
- S. Schulte im Walde. 2000. Clustering verbs semantically according to alternation behavior. In *Proceedings of the 18th International Conference on Computational Linguistics*, pages 747–753.
- N. Xue and M. Palmer. 2004. Calibrating features for semantic role labeling. In *Proceedings of 2004 Conference on Empirical Methods in Natural Language Processing*, pages 88–94.