

# Tagging and Chunking with Bigrams

Ferran Pla, Antonio Molina and Natividad Prieto

Universitat Politècnica de València

Departament de Sistemes Informàtics i Computació

Camí de Vera s/n

46020 València

{fpla,amolina,nprieto}@dsic.upv.es

## Abstract

In this paper we present an integrated system for tagging and chunking texts from a certain language. The approach is based on stochastic finite-state models that are learnt automatically. This includes bigram models or finite-state automata learnt using grammatical inference techniques. As the models involved in our system are learnt automatically, this is a very flexible and portable system.

In order to show the viability of our approach we present results for tagging and chunking using bigram models on the Wall Street Journal corpus. We have achieved an accuracy rate for tagging of 96.8%, and a precision rate for NP chunks of 94.6% with a recall rate of 93.6%.

## 1 Introduction

Part of Speech Tagging and Shallow Parsing are two well-known problems in Natural Language Processing. A Tagger can be considered as a translator that reads sentences from a certain language and outputs the corresponding sequences of part of speech (POS) tags, taking into account the context in which each word of the sentence appears. A Shallow Parser involves dividing sentences into non-overlapping segments on the basis of very superficial analysis. It includes discovering the main constituents of the sentences (NPs, VPs, PPs, ...) and their heads. Shallow Parsing usually identifies non-recursive constituents, also called chunks (Abney, 1991) (such as non-recursive Noun Phrases or base NP, base VP, and so on). It can include determining syntactical relationships such as subject-verb, verb-object, etc., Shallow parsing which always follows the tagging process, is used as a fast and reliable pre-processing phase for full or partial parsing. It can be used for Information Retrieval Systems, Information Extraction, Text Summarization and Bilingual Alignment. In addition, it is also used to solve computational linguistics tasks such as disambiguation problems.

### 1.1 POS Tagging Approaches

The different approaches for solving this problem can be classified into two main classes depending

on the tendencies followed for establishing the Language Model (LM): the linguistic approach, based on hand-coded linguistic rules and the learning approach derived from a corpora (labelled or non-labelled). Other approximations that use hybrid methods have also been proposed (Voutilainen and Padró, 1997).

In the linguistic approach, an expert linguist is needed to formalise the restrictions of the language. This implies a very high cost and it is very dependent on each particular language. We can find an important contribution (Voutilainen, 1995) that uses Constraint Grammar formalism. Supervised learning methods were proposed in (Brill, 1995) to learn a set of transformation rules that repair the error committed by a probabilistic tagger. The main advantage of the linguistic approach is that the model is constructed from a linguistic point of view and contains many and complex kinds of knowledge.

In the learning approach, the most extended formalism is based on n-grams or HMM. In this case, the language model can be estimated from a labelled corpus (supervised methods) (Church, 1988)(Weischedel et al., 1993) or from a non-labelled corpus (unsupervised methods) (Cutting et al., 1992). In the first case, the model is trained from the relative observed frequencies. In the second one, the model is learned using the Baum-Welch algorithm from an initial model which is estimated using labelled corpora (Merialdo, 1994). The advantages of the unsupervised approach are the facility to build language models, the flexibility of choice of categories and the ease of application to other languages. We can find some other machine-learning approaches that use more sophisticated LMs, such as Decision Trees (Màrquez and Rodríguez, 1998)(Magerman, 1996), memory-based approaches to learn special decision trees (Daelemans et al., 1996), maximum entropy approaches that combine statistical information from different sources (Ratnaparkhi, 1996), finite state automata inferred using Grammatical Inference (Pla and Prieto, 1998), etc.

The comparison among different approaches is difficult due to the multiple factors that can be consid-

ered: the language, the number and type of the tags, the size of the vocabulary, the ambiguity, the difficulty of the test set, etc. The best results reported on the Wall Street Journal (WSJ) Treebank (Marcus et al., 1993), using statistical language models, have an accuracy ratio between 95% and 97% (depending on the different factors mentioned above). For the linguistic approach the results are better. For example, in (Voutilainen, 1995) an accuracy of 99.7% is reported, but certain ambiguities in the output remain unsolved. Some works have recently been published (Brill and Wu, 1998) in which a set of taggers are combined in order to improve their performance. In some cases, these methods achieve an accuracy of 97.9% (Halteren et al., 1998).

## 1.2 Shallow Parsing Approaches

Since the early 90's, several techniques for carrying out shallow parsing have been developed. These techniques can also be classified into two main groups: based on hand-coded linguistic rules and based on learning algorithms. These approaches have a common characteristic: they take the sequence of lexical tags proposed by a POS tagger as input, for both the learning and the chunking processes.

### 1.2.1 Techniques based on hand-coded linguistic rules

These methods use a hand-written set of rules that are defined using POS as terminals of the grammar. Most of these works use finite state methods for detecting chunks or for accomplishing other linguistic tasks (Ejehed, 1988), (Abney, 1996), (At-Mokhtar and Chanod, 1997). Other works use different grammatical formalisms, such as constraint grammars (Voutilainen, 1993), or combine the grammar rules with a set of heuristics (Bourigault, 1992). These works usually use a small test set that is manually evaluated, so the achieved results are not significant. The regular expressions defined in (Ejehed, 1988) identified both non-recursive clauses and non-recursive NPs in English text. The experimentation on the Brown corpus achieved a precision rate of 87% (for clauses) and 97.8 % (for NPs). Abney introduced the concept of chunk (Abney, 1991) and presented an incremental partial parser (Abney, 1996). This parser identifies chunks base on the parts of speech, and it then chooses how to combine them for higher level analysis using lexical information. The average precision and recall rates for chunks were 87.9% and 87.1%, respectively, on a test set of 1000 sentences. An incremental architecture of finite-state transducers for French is presented in (At-Mokhtar and Chanod, 1997). Each transducer performs a linguistic task such as identifying segments or syntactic structures and detecting subjects and objects. The system was evaluated on various

corpora for subject and object detection. The precision rate varied between 99.2% and 92.6%. The recall rate varied between 97.8% and 82.6%.

The NPTool parser described in (Voutilainen, 1993) identified maximal-length noun phrases. NPtool gave a precision rate of 95-98% and a recall rate of 98.5-100%. These results were criticised in (Ramshaw and Marcus, 1995) due to some inconsistencies and apparent mistakes which appeared on the sample given in (Voutilainen, 1993). Bourigault developed the LECTER parser for French using grammatical rules and some heuristics (Bourigault, 1992). It achieved a recall rate of 95% identifying maximal length terminological noun phrases, but he did not give a precision rate, so it is difficult to evaluate the actual performance of the parser.

### 1.2.2 Learning Techniques

These approaches automatically construct a language model from a labelled and bracketed corpus. The first probabilistic approach was proposed in (Church, 1988). This method learnt a bigram model for detecting simple noun phrases on the Brown corpus. Given a sequence of parts of speech as input, the Church program inserts the most probable openings and endings of NPs, using a Viterbi-like dynamic programming algorithm. Church did not give precision and recall rates. He showed that 5 out of 243 NP were omitted, but in a very small test with a POS tagging accuracy of 99.5%.

Transformation-based learning (TBL) was used in (Ramshaw and Marcus, 1995) to detect base NP. In this work chunking was considered as a tagging technique, so that each POS could be tagged with I (inside baseNP), O (outside baseNP) or B (inside a baseNP, but the preceding word was in another baseNP). This approach resulted in a precision rate of 91.8% and a recall rate of 92.3%. This result was automatically evaluated on a test set (200,000 words) extracted from the WSJ Treebank. The main drawback to this approach are the high requirements for time and space which are needed to train the system; it needs to train 100 templates of combinations of words.

There are several works that use a memory-based learning algorithm. These approaches construct a classifier for a task by storing a set of examples in memory. Each example is defined by a set of features that have to be learnt from a bracketed corpus. The Memory-Based Learning (MBL) algorithm (Daelemans et al., 1999) takes into account lexical and POS information. It stores the following features: the word form and POS tag of the two words to the left, the focus word and one word to the right. This system achieved a precision rate of 93.7% and a recall rate of 94.0% on the WSJ Treebank. However, when only POS information was used the performance decreased achieving a precision rate of 90.3% and a



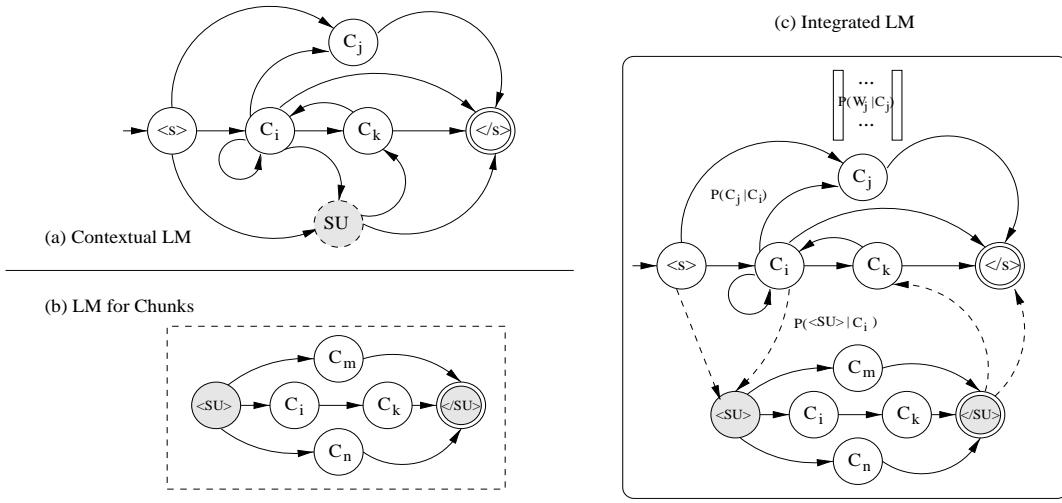


Figure 2: Integrated Language Model for Tagging and Chunking.

categories, the smoothing called "add one" is applied. Afterwards, a renormalization process is carried out.

Once the LMs have been learnt, a regular substitution of the lower model(s) into the upper one is made. In this way, we get a single Integrated LM which shows the possible concatenations of lexical tags and syntactical units, with their own transition probabilities which also include the lexical probabilities as well (see Figure 2c). Note that the models in Figure 2 are not smoothed).

## 2.2 The Decoding Process: Tagging and Parsing

The tagging and shallow parsing process consists of finding out the sequence of states of maximum probability on the Integrated LM for an input sentence. Therefore, this sequence must be compatible with the contextual, syntactical and lexical constraints. This process can be carried out by Dynamic Programming using the Viterbi algorithm, which is conveniently modified to allow for transitions between certain states of the automata without consuming any symbols (epsilon transitions). A portion of the Dynamic Programming trellis for a generic sentence using the Integrated LM shown in Figure 2c can be seen in Figure 3. The states of the automata that can be reached and that are compatible with the lexical constraints are marked with a black circle (i.e., from the state  $C_k$  it is possible to reach the state  $C_i$  if the transition is in the automata and the lexical probability  $P(W_i|C_i)$  is not null). Also, the transitions to initial and final states of the models for chunks (i.e., from  $C_i$  to  $<SU>$ ) are allowed; these states are marked in Figure 3 with a white circle and in this case no symbol is consumed. In all these cases, the transitions to initial and final pro-

duce transitions to their successors (the dotted lines in Figure 3) where now symbols must be consumed.

Once the Dynamic Programming trellis is built, we can obtain the maximum probability path for the input sentence, and thus the best sequence of lexical tags and the best segmentation in chunks.

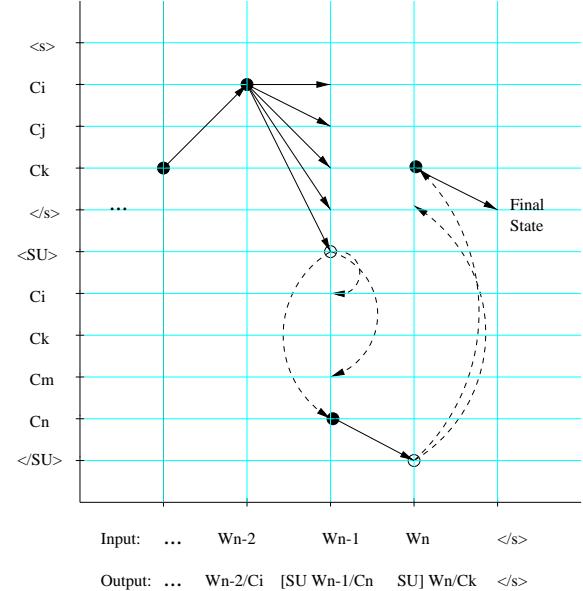


Figure 3: Partial Trellis for Programming Decoding based on the Integrated LM.

## 3 Experimental Work

In this section we will describe a set of experiments that we carried out in order to demonstrate the capabilities of the proposed approach for tagging and shallow parsing. The experiments were carried out



process. In Table 2, we show the most relevant results that we obtained for tagging and for NP chunking. The first row shows the result when the tagging and the chunking are done in a integrated way. The following rows show the performance of the sequential process using different taggers:

- LEX: it takes into account only lexical probabilities. In this case, the tagging accuracy was 94.3%.
- BIG: it is based on a bigram model that achieved an accuracy of 96.9%.
- IDEAL: it simulates a tagger with an accuracy rate of 100%. To do this, we used the tagged sentences of the WSJ corpus directly.

These results confirm that precision and recall rates increase when the accuracy of the tagger is better. The performance of the sequential process (using the BIG tagger) is slightly better than the performance of the integrated process (BIG-BIG). We think that this is probably because of the way we combined the probabilities of the different models.

## 4 Conclusions and Future Work

In this paper, we have presented a system for Tagging and Chunking based on an Integrated Language Model that uses a homogeneous formalism (finite-state machine) to combine different knowledge sources: lexical, syntactical and contextual models. It is feasible both in terms of performance and also in terms of computational efficiency.

All the models involved are learnt automatically from data, so the system is very flexible and portable and changes in the reference language, lexical tags or other kinds of chunks can be made in a direct way.

The tagging accuracy (96.9% using BIG and 96.8% using BIG-BIG) is higher than other similar approaches. This is because we have used the tag dictionary (including the test set in it) to restrict the possible tags for unknown words, this assumption obviously increase the rates of tagging (we have not done a quantitative study of this factor).

As we have mentioned above, the comparison with other approaches is difficult due among other reasons to the following ones: the definitions of base NP are not always the same, the sizes of the train and the test sets are different and the knowledge sources used in the learning process are also different. The precision for NP-chunking is similar to other statistical approaches presented in section 1, for both the integrated process (94.6%) and the sequential process using a tagger based on bigrams (94.9%). The recall rate is slightly lower than for some approaches using the integrated system (93.6%) and is similar for the

sequential process (94.1%). When we used the sequential system taking an error free input (IDEAL), the performance of the system obviously increased (95.5% precision and 94.7% recall). These results show the influence of tagging errors on the process. Nevertheless, we are studying why the results between the integrated process and the sequential process are different. We are testing how the introduction of some adjustment factors among the models for weighting the different probability distribution can improve the results.

The models that we have used in this work, are bigrams, but trigrams or any stochastic regular model can be used. In this respect, we have worked on a more complex LMs, formalized as a finite-state automata which is learnt using Grammatical Inference techniques. Also, our approach would benefit from the inclusion of lexical-contextual information into the LM.

## 5 Acknowledgments

This work has been partially supported by the Spanish Research Project CICYT (TIC97-0671-C02-01/02).

## References

- S. Abney. 1991. *Parsing by Chunks*. R. Berwick, S. Abney and C. Tenny (eds.) Principle-based Parsing . Kluwer Academic Publishers, Dordrecht.
- S. Abney. 1996. Partial Parsing via Finite-State Cascades. In *Proceedings of the ESSLLI'96 Robust Parsing Workshop*, Prague, Czech Republic.
- S. Argamon, I. Dagan, and Y. Krymolowski. 1998. A Memory-based Approach to Learning Shallow Natural Language Patterns. In *Proceedings of the joint 17th International Conference on Computational Linguistics and 36th Annual Meeting of the Association for Computational Linguistics, COLING-ACL*, pages 67–73, Montréal, Canada.
- S. At-Mokhtar and J.P. Chanod. 1997. Incremental Finite-State Parsing. In *Proceedings of the 5th Conference on Applied Natural Language Processing*, Washington D.C., USA.
- D. Bourigault. 1992. Surface Grammatical Analysis for the Extraction of Terminological Noun Phrases. In *Proceedings of the 15th International Conference on Computational Linguistics*, pages 977–981.
- Eric Brill and Jun Wu. 1998. Classifier Combination for Improved Lexical Disambiguation. In *Proceedings of the joint 17th International Conference on Computational Linguistics and 36th Annual Meeting of the Association for Computational Linguistics, COLING-ACL*, pages 191–195, Montréal, Canada.
- E. Brill. 1995. Transformation-based Error-driven Learning and Natural Language Processing: A

- Case Study in Part-of-speech Tagging. *Computational Linguistics*, 21(4):543–565.
- C. Cardie and D. Pierce. 1998. Error-Driven Pruning of Treebank Grammars for Base Noun Phrase Identification. In *Proceedings of the joint 17th International Conference on Computational Linguistics and 36th Annual Meeting of the Association for Computational Linguistics, COLING-ACL*, pages 218–224, Montréal, Canada, August.
- K. W. Church. 1988. A Stochastic Parts Program and Noun Phrase Parser for Unrestricted Text. In *Proceedings of the 1st Conference on Applied Natural Language Processing, ANLP*, pages 136–143. ACL.
- P. Clarksond and R. Ronnenfeld. 1997. Statistical Language Modeling using the CMU-Cambridge Toolkit. In *Proceedings of Eurospeech*, Rhodes, Greece.
- D. Cutting, J. Kupiec, J. Pederson, and P. Sibun. 1992. A Practical Part-of-speech Tagger. In *Proceedings of the 3rd Conference on Applied Natural Language Processing, ANLP*, pages 133–140. ACL.
- W. Daelemans, J. Zavrel, P. Berck, and S. Gillis. 1996. MBT: A Memory-Based Part-of-speech Tagger Generator. In *Proceedings of the 4th Workshop on Very Large Corpora*, pages 14–27, Copenhagen, Denmark.
- W. Daelemans, S. Buchholz, and J. Veenstra. 1999. Memory-Based Shallow Parsing. In *Proceedings of EMNLP/VLC-99*, pages 239–246, University of Maryland, USA, June.
- E. Ejerhed. 1988. Finding Clauses in Unrestricted Text by Finitary and Stochastic Methods. In *Proceedings of Second Conference on Applied Natural Language Processing*, pages 219–227. ACL.
- H. van Halteren, J. Zavrel, and W. Daelemans. 1998. Improving Data Driven Wordclass Tagging by System Combination. In *Proceedings of the joint 17th International Conference on Computational Linguistics and 36th Annual Meeting of the Association for Computational Linguistics, COLING-ACL*, pages 491–497, Montréal, Canada, August.
- S. M. Katz. 1987. Estimation of Probabilities from Sparse Data for the Language Model Component of a Speech Recognizer. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 35.
- D. M. Magerman. 1996. Learning Grammatical Structure Using Statistical Decision-Trees. In *Proceedings of the 3rd International Colloquium on Grammatical Inference, ICGI*, pages 1–21. Springer-Verlag Lecture Notes Series in Artificial Intelligence 1147.
- M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini. 1993. Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2).
- Lluís Màrquez and Horacio Rodríguez. 1998. Part-of-Speech Tagging Using Decision Trees. In C. Nédellec and C. Rouveiro, editor, *LNAI 1398: Proceedings of the 10th European Conference on Machine Learning, ECML'98*, pages 25–36, Chemnitz, Germany. Springer.
- B. Merialdo. 1994. Tagging English Text with a Probabilistic Model. *Computational Linguistics*, 20(2):155–171.
- F. Pla and N. Prieto. 1998. Using Grammatical Inference Methods for Automatic Part-of-speech Tagging. In *Proceedings of 1st International Conference on Language Resources and Evaluation, LREC*, Granada, Spain.
- L. Ramshaw and M. Marcus. 1995. Text Chunking Using Transformation-Based Learning. In *Proceedings of third Workshop on Very Large Corpora*, pages 82–94, June.
- A. Ratnaparkhi. 1996. A Maximum Entropy Part-of-speech Tagger. In *Proceedings of the 1st Conference on Empirical Methods in Natural Language Processing, EMNLP*.
- Atro Voutilainen and Lluís Padró. 1997. Developing a Hybrid NP Parser. In *Proceedings of the 5th Conference on Applied Natural Language Processing, ANLP*, pages 80–87, Washington DC. ACL.
- Atro Voutilainen. 1993. NPTool, a Detector of English Noun Phrases. In *Proceedings of the Workshop on Very Large Corpora*. ACL, June.
- Atro Voutilainen. 1995. A Syntax-Based Part-of-speech Analyzer. In *Proceedings of the 7th Conference of the European Chapter of the Association for Computational Linguistics, EACL*, Dublin, Ireland.
- R. Weischedel, R. Schwartz, J. Palmucci, M. Meteer, and L. Ramshaw. 1993. Coping with Ambiguity and Unknown Words through Probabilistic Models. *Computational Linguistics*, 19(2):260–269.