

Statistical Recognition of Noun Phrases in Unrestricted Text. *

José I. Serrano¹ and Lourdes Araujo²

¹Instituto de Automática Industrial. CSIC. Spain. nachosm@iai.csic.es

²Departamento de Sistemas Informáticos y Programación. Universidad Complutense de Madrid. Spain. lurdes@sip.ucm.es

Abstract. This paper presents a new model for flexible noun phrase detection, which is able to recognize noun phrases similar enough to the ones given by the inferred noun phrase grammar. To allow this flexibility, we use a very accurate set of probabilities for the transitions between the part-of-speech tag sequence which defines a noun phrase. These accurate probabilities are obtained by means of an evolutionary algorithm, which works with both, positive and negative examples of the language, thus improving the system coverage, while maintaining its precision. We have tested the system on different corpora and compare the results with other systems, what has revealed a clear improvement of the performance.

1 Introduction

There are many Natural Language Processing (NLP) applications for which noun phrase (NP) detection is useful. For instance, NPs can be used in the identification of multiword terms, which are mainly noun phrases, or as a preprocessing step for a subsequent complete syntactic analysis. However, probably the most direct and relevant application nowadays is in information recovering. NPs recover most of the information content of a document, helping to detect the topics it is about. In fact, document indexing by NPs has been shown to improve on other kinds of indexing when retrieval is carried out over a very large corpus [1]. Thus, the distribution of noun phrases can guide search engines in collecting relevant documents according to user queries, or they can be used in text summarizing, in machine translation, etc.

Different approaches have been proposed for the NP identification problem. Some of them rely on linguistic knowledge and use a hand-coded language model. Bourigault [2] uses a handcrafted NP grammar along with some heuristics for identifying NPs of maximal length, and Voutilainen [3] uses a constraint grammar formalism. Other proposals follow a learning approach based on the use of corpora. Church [4] uses a probabilistic model automatically trained on the Brown corpus to detect NPs as well as to assign parts of speech. Ramshaw & Marcus [5] uses the supervised learning methods proposed by Brill[6] to learn

* Supported by projects TIC2003-09481-C04 and FIT150500-2003-373

a set transformation rules for the problem. Pla and Prieto [7] use grammatical inference to obtain a FSA which recognizes NPs, what has inspired this work.

This paper presents a new model for a *flexible* identification of basic (non-recursive) NPs in an arbitrary text. The system generates a probabilistic finite-state automaton (FSA) able to recognize the sequences of lexical tags which form an NP. The FSA is generated with the Error Correcting Grammatical Inference (ECGI) algorithm of grammatical inference and initial probabilities are assigned using the collected bigram probabilities. The FSA probabilities provide a method for a flexible recognition of input chains, which are considered to be NPs even if they are not accepted by the FSA but are similar enough to an accepted one. Thus, the system is able to recognize NPs not present in the training examples, what has proven very advantageous for the performance of the system. The FSA probabilities, which are crucial for the flexibility in recognition, are optimized by applying an evolutionary algorithm (EA), what produces a highly robust recognizer. The EA uses both, positive and negative training examples, what contributes to improve the coverage of the system while maintaining a high precision.

Though it is difficult to compare different approaches because they differ in multiple elements, some attempts have been made. Pla [8] gathers results of a number of parses trained on the data set used by Ramshaw [5]: the one proposed by Cardie & Pierce [9], whose technique consists in matching part-of-speech (POS) tag sequences from an initial NP grammar extracted from an annotated corpus and then ranking and pruning the rules according to their achieved precision; the memory-based approach presented in [10], which introduces cascaded chunking, a process in two stages in which the classification of the first stage is used to improve the performance of the second one; the Memory-Based Sequence Learning (MBSL) algorithm [11], which learns sequences of POS tags and brackets, and the hybrid approach of Tjong-Kim-Sang [12], which uses a weighted voting to combine the output of different models. We have also applied our model to the same test set in order to compare the results.

The rest of the paper proceeds as follows: section 2 describes the general scheme of the system, presenting their main elements and its relationships; section 3 is devoted to describe the evolutionary algorithm used to training the FSA; section 4 presents and discusses the experimental results, and section 5 draws the main conclusions of this work.

2 General Scheme of the System

This work presents the design and implementation of a flexible recognizer of basic NPs given as chains of POS tags. Because the system is constructed from examples of objective syntagmas, an algorithm of grammar inference will be used (ECGI). This algorithm extracts from these examples a FSA, which represents the grammar defined by them. In its turn, this algorithm uses the Viterbi algorithm, a dynamic programming one, so that the FSA has an optimal number of states and paths. In order to extend the recognition capabilities of the FSA

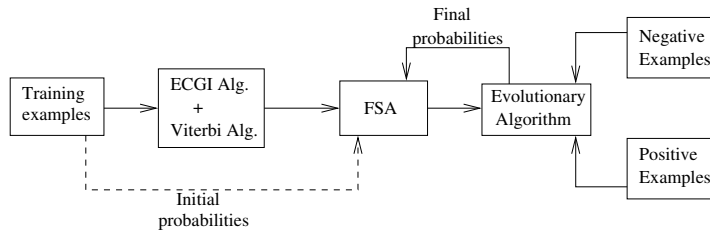


Fig. 1. Scheme of the process to generate the NP flexible recognizer. The FSA is constructed from a set of training examples by applying the ECGI algorithm and using Viterbi to find the most similar chain of tags for a given input. Initial probabilities for the FSA are also obtained from this set of training examples. Then, the probabilities are tuned by applying an evolutionary algorithm which uses new examples, both positive and negative.

beyond the set of training examples, it is necessary to endow the system with a mechanism to recognize also other NPs not appearing in the examples, but very similar to the ones gathered in the grammar. This mechanism to introduce flexibility amounts to allowing the FSA to recognize NPs with a percentage of error. If the number of errors incurred in parsing an NP is below a threshold value, then the new NP is also recognized. An error appears whenever the input tag does not produce any transition from the actual state of the FSA. In order to allow the FSA to find NPs similar to those of the input, their edges are labeled with the probability of each transition between tags within an NP. These probabilities are initially extracted from the NPs examples, and afterwards an evolutionary algorithm is applied to optimize the probabilities. This algorithm uses NP examples different from those used in the FSA construction, as well as negative examples, i.e. syntagmas which must not be recognized as NPs. Figure 1 shows a scheme of the relationships between the different components and phases involved in the flexible recognizer. Let us now describe the different elements of the system.

2.1 Building the FSA for NP recognition

The technique used to obtain the FSA is an algorithm of Grammar Inference, i.e. a process which, from a set of examples, obtains structural patterns of the language and represents them by means of a grammar. Dupont [13] proposes a general scheme for selecting the most appropriate algorithm of grammar inference under different conditions. According to these ideas, we have chosen the Error Correcting Grammatical Inference (ECGI) algorithm [14] since we are interested in a heuristic construction of a regular grammar such that the final result allows for a flexible recognition. This technique involves a progressive construction of the FSA from positive examples. The FSA is modified in order to correct the errors appearing in the parsing of a new example, obtaining non recursive grammars (a FSA without cycles).

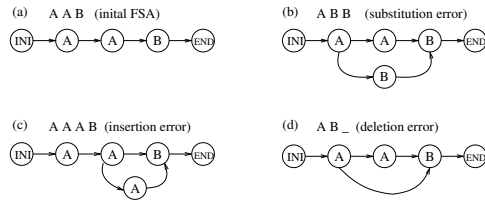


Fig. 2. Examples of application of the ECGI algorithm.

The algorithm constructs the grammar by adding in an incremental way the training examples. In order to avoid introducing noise in the recognizer which could drive it to increase the number of false positives (non-nominal syntagmas recognized as nominal), it is necessary to remove this noise by applying a preprocessing step to the training examples. Thus, the typically non-nominal syntagmas which behave as nominal in the training set (the nominalization of a verb, for example) are handily detected and eliminated from the training list. After the preprocessing, the algorithm proceeds as follows. First, a simple FSA is generated which recognizes the first example of the training set. This FSA is then extended with the other examples. To introduce a new example, the Viterbi algorithm is used to find the sequence of states which recognizes the example with less errors. Then, the FSA is extended by adding states and/or transitions which correct the produced errors. In this way, when the process of parsing all the examples finishes, the resulting FSA is able to recognize all of them and does not present cycles.

The errors which can appear between an input chain and the recognized ones are of three types: insertion, substitution and deletion. Figure 2 shows examples of each of them. When the FSA of Figure 2(a) tries to recognize the input chain (A B B) (Figure 2(b)), a substitution error is detected and solved by adding a new state to the FSA with the label of the error. If the new chain is (A A A B) (Figure 2(c)) an insertion error appears, which is solved by adding a new state with the missing label. Finally, if the chain is (A B) (Figure 2(d)), the deletion error is managed by adding a transition between the states neighbor to the deleted one. The result of the algorithm is a regular grammar which generates a finite language, because by construction it has no cycles.

We have introduced a simplification of this algorithm already proposed in the literature [15]. It aims at reducing the complexity of the algorithm to find in the FSA the closest path to the input chain. This simplification amounts to ignoring the deletion errors. This can reduce the recognition complexity below 10%. However, this simplification also reduces the generality of the language, since it discards chains with lengths different from the training examples. Thus, these examples must be sufficiently varied and heterogenous.

The selection of the most similar chain to the input one is carried out by applying the Viterbi algorithm. This algorithm was initially designed to find the most probable path in a Markov chain which produces a given sequence of

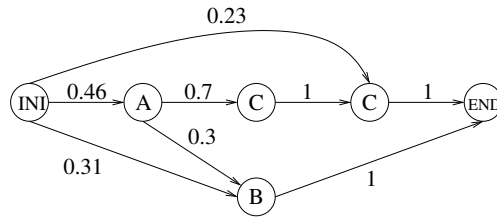


Fig. 3. Probabilistic FSA for flexible NPs recognition.

tags [16]. This algorithm has also been applied to search the minimum path in a multi-stage graph. In our case, the goal is to minimize the number of errors in the FSA for an input chain.

We have modified the Viterbi algorithm to process input chains with a length different from the examples chains. Thus, if the input chain is longer than the longest chain the FSA recognizes, an edge from each state to every other one, including itself, is added and those edges are considered to produce transition errors.

2.2 Flexible Recognition of NPs

The FSA built so far basically recognizes those NPs present in the training examples (some other sequences can also be recognized, though in general they are not NPs). Because we are interested in a recognizer as general as possible and because the training examples usually include only a small sample of the language, it is necessary to endow the system with a mechanism of generalization. This mechanism amounts to allowing the FSA to recognize a chain if it is sufficiently similar to an accepted chain, i.e. if the number of differences is below a certain threshold value. This procedure does not introduces too many chains outside the language, as the experiments have shown. For example, the FSA of Figure 3 does not recognize the chain (D C C). However, if the error threshold value is 1, the chain will be recognized, since the chain (A C C) is in fact recognized. Obviously, the error threshold value is a determining factor for the generalization of the model, and it is object of a detailed study in the experiments.

Now the question is how to parse a chain which presents errors. In this case, the parse arrives at a situation in which no transition is possible, since the next input tag does not matches any of the state for which there is a transition. At this point, we use the statistics derived from the training examples. Each transition is labeled with a real value, which represents its probability in the target language. Then, to process an error tag, the parse follows the most probable transition. The probabilities are relative to each state, in such a way that the values of all the transitions leaving a same state add up to one. Each probability is initially computed as the number of occurrences of a transition relative to the others of

the same state according to the training examples. For example, in the FSA of Figure 3, if transition AC appears 7 times in the examples and transition AB appears 3 times, 10 transitions exist from A to C or B, and thus the probability of the edges that leave state A are 0.7(C) and 0.3(B) ($AC = 7/10$; $AB = 3/10$). With the probabilities of the edges so determined, if the chain (D C C) is parsed, as there is not transition from the initial state to D, the parse will take the most probable transition, which leads to A, and afterwards the process continues with the other tags of the input, (C C), which are correct transitions. Thus, the chain (A C C) is obtained with one error with respect to the input. In this way, the parse produces a chain very similar to the input one and which occurs with high probability in the language. Thus, we can expect that a chain which has been only partially recognized, belongs to the language if the error is small.

3 Optimizing the Automata Probabilities: Evolutionary Algorithm

Now the goal is to find a set of probabilities for the edges of the FSA described in the previous section, which allow it to recognize as many NPs as possible, avoiding at the same time recognizing false NPs. Accordingly, the search space is composed of the different sets of probabilities for the edges of the FSA and the algorithm looks for those which optimize the recognition of a set of training examples. This complex search is performed with an evolutionary algorithm.

Systems based on evolutionary algorithms maintain a population of potential solutions and are provided with some selection process based on the fitness of individuals, as natural selection does. The population is renewed by replacing individuals with those obtained by applying “genetic” operators to selected individuals. The most usual “genetic” operators are *crossover* and *mutation*. Crossover obtains new individuals by mixing, in some problem dependent way, two individuals, called parents. Mutation produces a new individual by performing some kind of random change on an individual. The production of new generations continues until resources are exhausted or until some individual in the population is fit enough. Evolutionary algorithms have proven to be very useful search and optimization methods, and have previously been applied to different issues of natural language processing [17], such as text categorization [18], tagging [19] and parsing [20].

In our case, individuals represent probabilistic FSAs, all of them with the same structure but different probabilities. They are implemented as an adjacency matrix, what facilitates the application of the genetic operators. Thus, the probability an edge going from state i to state j is the value which appears in row i column j , A_{ij} . Transitions which do not exist are assigned the value -1 . For example, the FSA of Figure 4(a) is implemented as the matrix of Figure 4(b).

The genetic operators applied to renew the population are crossover and mutation. Two variants of crossover have been implemented. The first one is the classic one point crossover, which combines two individuals by combining

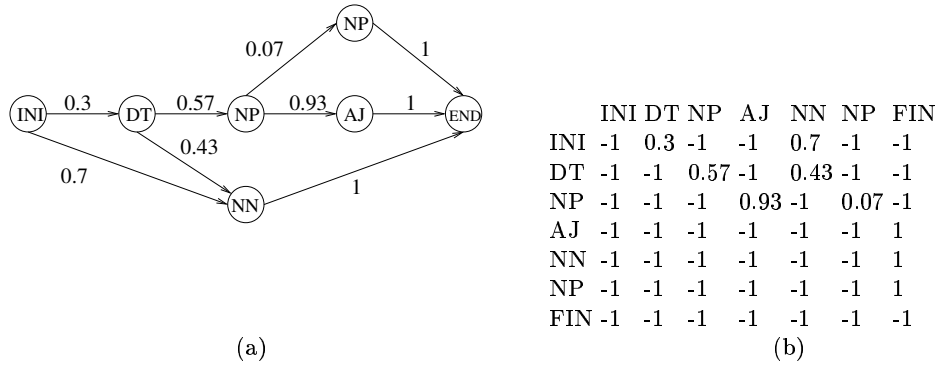


Fig. 4. Individuals representation.

the first part of one parent up to a crossover point with the second part of the other parent and vice versa. The second variant uses two crossover points and exchanges the bit of the two parents between these points. In both variants, the crossover points are randomly selected.

The mutation operator amounts to choosing an edge at random and varying its probability. This variation may be positive or negative, what is decided at random, and the amount is an input parameter to the algorithm, studied in the experiments. Notice that when an edge is modified the remaining edges of the same state must be updated in order to maintain the value of its addition equal to 1. If the variation produces a value smaller than the zero or greater than one, then the edge is assigned zero or one respectively, and the real variation is computed according to this value.

3.1 The Fitness Function

The fitness function must be a measure of the capability of the FSA to recognize NPs and only them. On the one hand, the fitness function must include a measure of the coverage, or *recall* achieved by the individual, i. e. the number of NPs which have been recognized from the set of proposed NPs:

$$recall = \frac{\text{number of recognized NPs}}{\text{number of proposed NPs}}$$

On the other hand, the fitness function must also take into account the precision achieved by the individual:

$$precision = \frac{\text{number of NPs recognized} + \text{number of non NPs discarded}}{\text{number of proposed syntagmas}}$$

We have considered as fitness function two F-measures which combine these two parameters in different ways:

$$F\text{-measure} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

and F2-measure, which gives a higher weight to precision,

$$\text{F2-measure} = \frac{5 \cdot \text{Precision} \cdot \text{Recall}}{4 \cdot \text{Precision} + \text{Recall}}$$

In order to apply these measures to an individual, we must establish the cases in which an NP is considered recognized. Obviously those input chains corresponding to a path from the initial to the final state are recognized by any individual in the population. But in the remaining cases, the path will depend on the probabilities of the FSA, and only those chains with a number of errors below the threshold will be recognized. This threshold value can be defined as an absolute value or as a percentage of the length of the chain. Both alternatives have been evaluated in the experiments.

3.2 Initial Population

Individuals for the initial population are randomly generated from a seed solution, which helps to guide the search. The seed, which is included in the population as another individual, is the set of probabilities obtained from the training examples used to construct the FSA. The individuals of the remaining population are generated by applying the mutation operator several times to the seed. The variation produced by each mutation in this phase is a parameter studied in the experiments. Notice that a small variation would produce individuals too similar to the seed, which in spite of having a high quality according to the fitness function, do not help to explore other areas of the search space where better solutions could be found, while a great variation can lead to lose the advantages of the seed.

4 Experimental Results

The algorithm has been implemented using C++ language and run on a Pentium III processor. The CPU time spent on generating an automaton from 45 examples, with a maximum length of 7, was 0.3 seconds, and from 156 examples, with 9 as maximum length, 1.6 seconds. The time spent on analyzing a text composed of 1108 different syntagmas, being NPs 570 of them, was 1.1 seconds. For the experiments we have used training and test sets from the Brown corpus portion of the Penn Treebank [21], a database of English sentences manually annotated with syntactic information. After a number of experiments, we have selected a set of default criteria and parameters for the EA, which provide a high performance for different settings of the problem, and which have been used in the experiments described below. The EA uses a two point crossover, F1-measure as fitness function, and elitism¹. The parameter values are a population size of 100 individuals, a number of 150 generations, a crossover rate of 60%, a mutation

¹ By elitism we refer to the technique of retaining in the population the best individuals found so far.

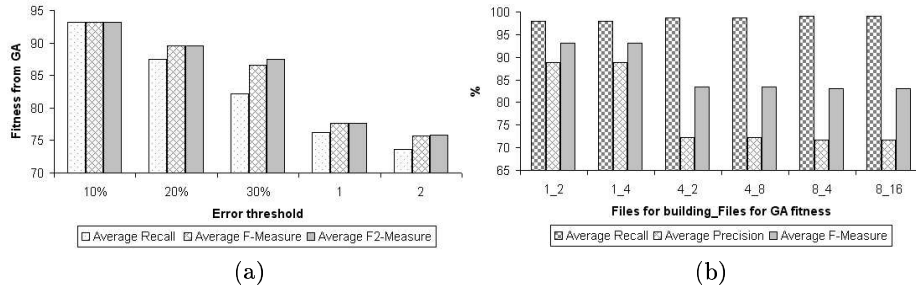


Fig. 5. (a) Error thresholds and fitness function types comparison. (b) F-Measure values for different combinations of number of files for building the recognizer-number of files for EA fitness.

rate of 40%, a variation of FSA probabilities in mutation of 0.5, and a variation of FSA probabilities in the generation of the initial population of 0.7. Each EA has been run five times, and the average and maximum values of these executions are presented for each experiment. Setting the EA parameters to these values, the time spent on ten executions of the EA was over 3 hours and a half, using a test set of 1569 different syntagmas, 543 NPs, for the fitness function calculation.

Before any other empirical trials, we have carried out a test in order to show the benefits of the EA over the recognizer. We have found out that the higher the error threshold the better the improvement. For relative error thresholds of 10%, 20% and 30% the F-Measure over the test increases 4%, 10% and 14%, respectively. For absolute error threshold values of 1 and 2 the increase is about 5% and 12%, respectively.

We also have performed experiments to determine the most appropriate error threshold allowed in the recognition. Figure 5 (a) shows the results obtained with different performance measures. Two different ways of defining the threshold value have been studied: as an absolute number of errors, and as a percentage of the NP length. Best results are obtained with the threshold defined as a percentage and for relative low values, such as 10%. This value has been used in the following experiments.

Another question which has been investigated is the most appropriate size of the training sets for both, the construction of the FSA (C), and for the evaluation of fitness measure of the EA (F). Figure 5 (b) shows the results of the different considered measures for different combination of these values (C-F). Best results for all measures are obtained when using a small number of files for the FSA construction, and a large number of files for the EA. Using a larger training set in the construction of the FSA produces too large FSAs, which recognize too much false NPs and deteriorate the system performance. We can in fact observe that the recall measures are high when using a larger set in the construction of the FSA, at the price of producing very low precision measures.

Since the system also uses training examples to establish the initial probabilities from which the initial population of EA is generated, we have performed

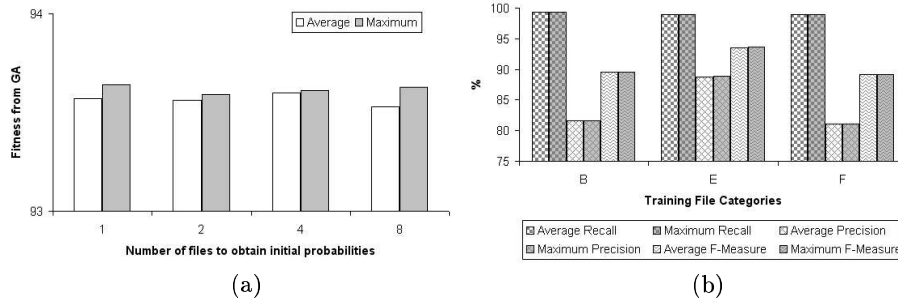


Fig. 6. (a) F-measures values from the GA using different number of files to obtain the initial probabilities or seed. (b) Accuracy values for different categories of the file used to build the recognizer.

a test in order to study the influence of the number of examples used this way in the EA performance. As observed in Figure 6 (a), there exists an improvement the more examples are used, but it is not very significant. However, the decrease in performance when using eight files indicates that it is not convenient to use too many training examples, because the extracted probabilities become overfitted.

We have also performed experiments to determine the impact of using different categories of topics as training sets. Figure 6 (b) shows the results obtained using training files of category B (press editorial), E (skills and hobbies) and F (popular lore). We can observe that the results are quite different depending on the used category, since the kind and frequency of NPs is different in each of them.

In order to compare the overall performance with other related systems, we have applied the system to a standard subset of the Wall Street Journal portion of the Penn Treebank, also used in the other works. Table 1 compares the precision, recall and F-measure values. We can observe that our results improve on those of all the other systems, both in recall and F-measure, and only the Ramshaw95 and Tjong-Kim-Sang00 systems provide a slightly better precision. It proves the usefulness of the mechanism of flexible recognition, since it achieves, as expected,

Table 1. Comparison of precision, recall and F-measure values with similar systems over the Ramshaw standard corpus.

	<i>Precision (%)</i>	<i>Recall (%)</i>	<i>F-Measure (%)</i>
[Ramshaw95]	91,8	92,3	92,0
[Cardie98]	90,7	91,1	90,9
[Veenstra98]	89,0	94,3	91,6
[Argamon98]	91,6	91,6	91,6
[Tjong-Kim-Sang00]	93,6	92,9	93,3
[This work]	91,6	97,8	94,5

a significant improvement of the performance, maintaining at the same time a high level of precision thanks to the accurate probabilities that the EA provides.

5 Conclusions

This paper presents the design of a probabilistic FSA for the identification of basic NPs, which can be used for document indexing in information retrieval and other applications.

The FSA probabilities provide a method for a flexible recognition of input chains, which are considered to be NPs if they are similar enough to an accepted one. An evolutionary algorithm is used to optimize the FSA probabilities, and it uses both, positive and negative training examples, what contributes to improve the coverage of the system maintaining at the same time a high precision.

We have performed experiments to tune the system in several ways, such as the EA parameters, the threshold value of the errors allowed in the recognition, and the size of the training sets used in different phases of the process. A comparison with other systems tested on the same set of texts has shown that the mechanism introduced for the flexible recognition of NP clearly improves the coverage of the system, while the adjusted probabilities provided by the EA yield a high level of precision, thus yielding a better overall performance.

A study of the most frequent errors has revealed that those affecting the recall are mainly due to the presence of rare NPs, with a structure very different from the others. However, as they are infrequent, the reduction of the recall is low. The errors which deteriorate the precision are mainly due to some training examples used in the construction of the FSA in which a sequence of tags, which usually does not correspond to an NP, is working as an NP in that particular case, and also to sequences in which many of their tags are not typical NPs tags. According to this observation, we plan to perform a statistical study of the different kinds of errors, in order to introduce mechanisms to filter the training set in some appropriate and automatic manner. We also consider to investigate other improvements, on different aspects of the system: on the application of the EA, with the definition of other genetic operators and fitness functions, on the recognition process, with the consideration of information from negative examples, and on the construction of the FSA, with the possibility of introducing cycles.

References

1. Zhai, C.: Fast statistical parsing of noun phrases for document indexing. In: Proceedings of the Fifth Conference on Applied Natural Language Processing. (1997)
2. Bourigault, D.: Surface grammatical analysis for the extraction of terminological noun phrases. In: Proc. of the Int. Conf. on Computational Linguistics (COLING-92). (1992) 977–981
3. Voutilainen, A.: Nptool, a detector of english noun phrases. In: Proc. of the Workshop on Very Large Corpora (ACL). (1993) 48–57

4. Church, K.W.: A stochastic parts program and noun phrase parser for unrestricted text. In: Proc. of 1st Conference on Applied Natural Language Processing, ANLP. (1988) 136–143
5. Ramshaw, L., Marcus, M.: Text chunking using transformation-based learning. In: Proc. of the third Workshop on Very Large Corpora (ACL). (1995) 82–94
6. Brill, E.: Transformation-based error-driven learning and natural language processing: A case study in part of speech tagging. *Computational Linguistics* **21** (1995)
7. Pla, F., Molina, A., Prieto, N.: Tagging and chunking with bigrams. In: Proc. of the 17th conference on Computational linguistics. (2000) 614–620
8. Pla, F.: Etiquetado léxico y análisis sintáctico superficial basado en modelos estadísticos (2000)
9. Cardie, C., Pierce, D.: Error-driven pruning of treebank grammars for base noun phrase identification. In: Proc. of COLING-ACL'98. (1998) 218–224
10. Veenstra, J.: Fast np chunking using memory-based learning techniques. In: Proc. of BENELEARN-98: Eighth Belgian-Dutch Conference on Machine Learning. (1998) 71–78
11. Argamon, S., Dagan, I., Krymolowski, Y.: A memory-based approach to learning shallow natural language patterns. In: Proc. of joint International Conference COLING-ACL. (1998) 67–73
12. Tjong-Kim-Sang, E.F.: Noun phrase representation by system combination. In: Proc. of ANLP-NAACL. (2000) 50–55
13. Dupont, P.: Inductive and statistical learning of formal grammars. Technical report, Reseach talk, Departement ingenierie Informatique, Universite Catholique de Louvain (2002)
14. Rulot, H., Vidal, E.: Modelling (sub)string-length-based constraints through a grammatical inference method. In: *Pattern Recognition: Theory and Applications*. Springer-Verlag (1987) 451–459
15. Torró, F., Vidal, E., Rulot, H.: Fast and accurate speaker independent speech recognition using structural models learnt by the ecgi. In: *Signal Processing V: Theories and Applications*. Elsevier Science Publishers B.V. (1990)
16. Forney, G.D.: The viterbi algorithm. *Proceedings of The IEEE* **61** (1973) 268–278
17. Kool, A.: Literature survey. Technical report, Center for Dutch Language and Speech. University of Antwerp (2000)
18. Serrano, J., Castillo, M.D., Sesmero, M.: Genetic learning of text patterns. In: Proc. of CAEPIA03. (2003) 231–234
19. Araujo, L.: Part-of-speech tagging with evolutionary algorithms. In: Proc. of the Int. Conf. on Intelligent Text Processing and Computational Linguistics (CICLing-2002), Lecture Notes in Computer Science 2276, Springer-Verlag (2002) 230–239
20. Araujo, L.: A probabilistic chart parser implemented with an evolutionary algorithm. In: Proc. of the Int. Conf. on Intelligent Text Processing and Computational Linguistics (CICLing-2004), Lecture Notes in Computer Science 2276, Springer-Verlag (2004) 81–92
21. Marcus, M.P., Santorini, B., Marcinkiewicz, M.A.: Building a large annotated corpus of english: The penn treebank. *Computational Linguistics* **19** (1994) 313–330