

CHAIN ALGORITHM USED FOR PART OF SPEECH RECOGNITION

ANDREEA-DIANA MIHIS⁽¹⁾

ABSTRACT. Dictionary base methods have the advantage that they can be applied to texts written in different languages if there exist an electronically dictionary for that specific language. As a consequence, these kinds of methods can be used to identify the part of speech of words from a text written in a single language. The principal advantage of using a dictionary base approach in part of speech recognition is that it can be applied to different languages, because it does not use a specific grammar. In this paper such a method is used to identify the parts of speech of words from a text using a chain algorithm [7] that disambiguates a text.

1. INTRODUCTION

When learning a language, the use of a dictionary is essential, even though grammar skills are not fully developed. When translating a text from a foreign language, the first step consists in identifying the corresponding word, followed by applying mostly grammar skills of the familiar language in trying to obtain the correct translation. When translating from a familiar language to a foreign language the most used method is word by word (the french mot-a-mot). From beginners point of view the dictionary is a powerful tool, that sometimes can be used successfully to overcome most difficulties. This should be applicable to a computer. Like the beginner, the computer can use the dictionary explanation to choose the correct form of a word and to decide their speech part in a text. This method can be used to many languages, as long as there is an appropriate electronical dictionary.

Pertaining to Natural language processing, this method uses only a dictionary to identify the part of speech of a given word. The basic idea is that in the dictionary, every word will have different definitions, one or more for every part of speech of the specified word. Trying to identify the correct part of speech from

2000 *Mathematics Subject Classification.* 68T35, 68T50, 91F20.

Key words and phrases. chain algorithm, part of speech recognition.

the list, is almost equivalent with trying to match the definitions. And matching the definitions is done using the Chain Algorithm.

2. CHAIN ALGORITHM FOR PART OF SPEECH RECOGNITION

The chain algorithm for part of speech recognition is based on and includes the chain algorithm presented in [7].

The idea of the algorithm is to try to disambiguate all the words, in groups of three, but, because the part of speech (POS) is unknown, the word is searched in dictionary with all POS(s). If a word is not found with a POS, then the senses of that word for that POS does not exists. Finally, the POS and the sense merge from the greatest score.

The algorithm for POS recognition by disambiguating a triplet of words $w_1 w_2 w_3$ for Dice measure is the following:

```

begin
  for each POS  $p^{w_1}$  do
    for each POS  $r^{w_2}$  do
      for each POS  $q^{w_3}$  do
        begin
          for each sense  $s_{p^{w_1}}^i$  do
            for each sense  $s_{r^{w_2}}^j$  do
              for each sense  $s_{q^{w_3}}^k$  do
                 $score(i, j, k) = 3 \times \frac{|D_{p^{w_1}} \cap D_{r^{w_2}} \cap D_{q^{w_3}}|}{|D_{p^{w_1}}| + |D_{r^{w_2}}| + |D_{q^{w_3}}|}$ 
              endfor
            endfor
          endfor
        end
      endfor
    endfor
  endfor
  endfor
  endfor
  endfor
   $(P, R, Q, i^*, j^*, k^*) = argmax_{(p,r,q,i,j,k)} score(i, j, k)$ 
  /* POS of  $w_1$  is P, sense of  $P^{w_1}$  is  $s_{P^{w_1}}^{i^*}$ 
  POS of  $w_2$  is R, sense of  $R^{w_2}$  is  $s_{R^{w_2}}^{j^*}$ 
  POS of  $w_3$  is Q, sense of  $Q^{w_3}$  is  $s_{Q^{w_3}}^{k^*}$  */
end

```

Here, α^w denotes the part of speech α of word w , and $s_{\alpha^w}^i$ is the sense with number i for the word w with POS α . D_{α^w} is the complete dictionary definition for word w with POS α .

For the overlap measure the score is calculated as:

$$score(i, j, k) = \frac{|D_{p^{w_1}} \cap D_{r^{w_2}} \cap D_{q^{w_3}}|}{\min(|D_{p^{w_1}}|, |D_{r^{w_2}}|, |D_{q^{w_3}}|)}$$

For the Jaccard measure the score is calculated as:

$$score(i, j, k) = \frac{|D_p^{w_1} \cap D_r^{w_2} \cap D_q^{w_3}|}{|D_p^{w_1} \cup D_r^{w_2} \cup D_q^{w_3}|}$$

The POS recognition algorithm is a chain algorithm. This means that for the first triplet it is used to identify the POS(s) of all tree words, and for the following ones, only for the last word, because the POS(s) of first two are already identified. As a consequence, the complexity of the algorithm decreases:

```

begin
  for each POS  $q^{w_3}$  do
    begin
      for each sense  $s_{q^{w_3}}^k$  do
         $score(i^*, j^*, k) = 3 \times \frac{|D_P^{w_1} \cap D_R^{w_2} \cap D_q^{w_3}|}{|D_P^{w_1}| + |D_R^{w_2}| + |D_q^{w_3}|}$ 
      endfor
    end
  endfor
   $(P, R, Q, i^*, j^*, k^*) = argmax_{(p,r,q,i,j,k)} score(i, j, k)$ 
  /*POS of  $w_3$  is Q, sense of  $Q^{w_3}$  is  $s_{Q^{w_3}}^{k^*}$  */
end

```

As it can be seen in detail in [7], some "stop" words were used to eliminate the trivial cases of glosses intersection.

3. EXPERIMENTAL EVALUATION

To test the POS recognition algorithm were used ten files from Brown corpus thanks to the possibility offered by SemCorpus to evaluate the results. The input files were not tagged, and they contain compound words, underscored.

The algorithm was tested only for English Language, having WordNet as base dictionary, implying the following POS(s): Nouns, Verbs, Adjectives and Adverbs.

The application that was implemented to test the algorithm has as output the list of words tagged with their POS(s) and their corresponding sense, as follows: *noun#n#i*. Here the noun *noun* has the POS noun and the sense i from WordNet. If #.# were missing, than the algorithm failed to find the POS and sense of the corresponding word.

Because in SemCor exist some words without tag (Notag) in computing the precision, they and the stop words are ignored.

The final results can be seen in Figure 1 and in detail in annex.

In the upper part of the figure is displayed the precision for the whole file, ascending. The precision for every part of speech is displayed in the lower part of the figure, using the same order as for the first one. The precision for a specific POS is computed thus if a word has a specific POS, it will be recognized with the precision found in the graph. There are some words that are identified to have a specific POS, even if in that case they have another POS, but these words are not considered in computing the precision shown in the Figure 1.

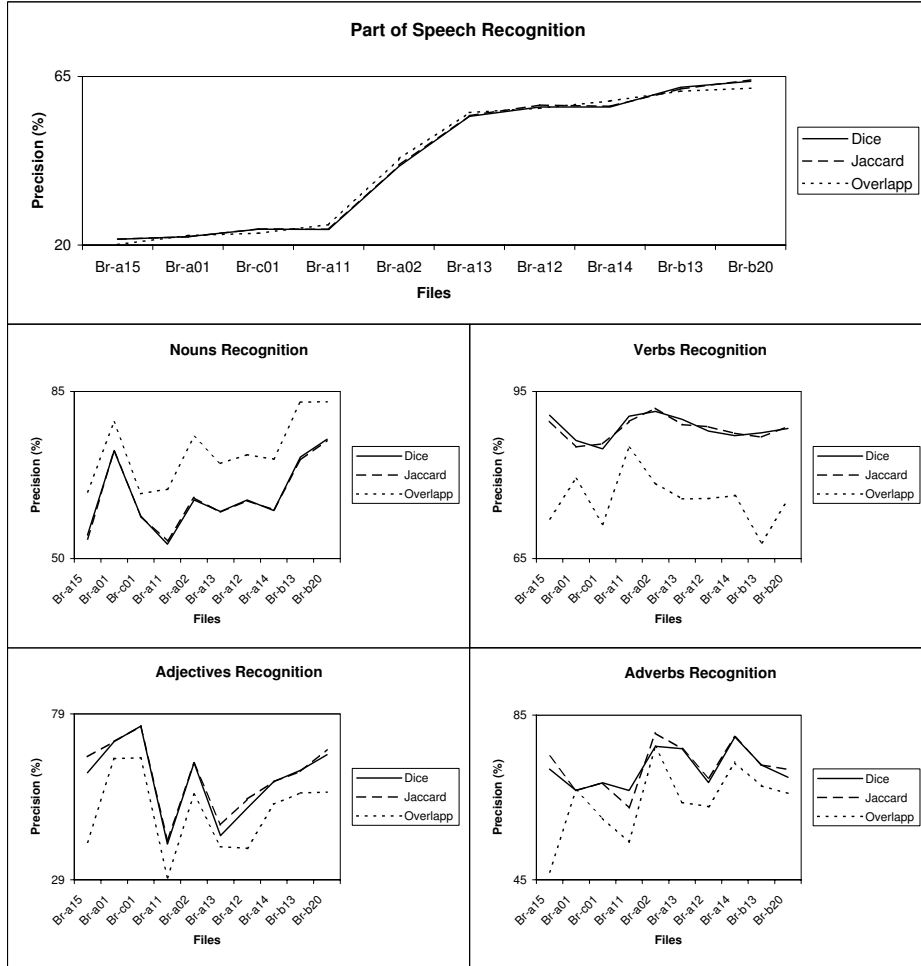


FIGURE 1. Precision of POS Recognition sorted ascending by the global precision

Because not all POS(s) were considered in applying the POS recognition algorithm, the precision for the whole file is lower than the precision for every POS. It seems that the highest precision is for verbs and the lowest one is for adjectives. In case of Nouns, the highest results were obtained using Overlapp measure, but in the other cases, with Overlapp measure the lowest results were obtained.

Table 1 summarizes the results.

Precision	Minimum Value	Maximum Value	Average Value
POS Recognition	20.16	64.13	42.78
Disambiguation	7.00	41.43	21.78
Nouns	53.04	82.83	66.12
Verbs	67.83	92.02	84.45
Adjectives	29.55	75.34	56.39
Adverbs	46.87	80.65	69.41

TABLE 1. Minimum, Maximum and Average Value for Precision

The Precision for disambiguating the whole text is not great, only 21,78% in average. But the precision for a specific POS is by average more then 68%.

4. APPLICATION OF POS RECOGNITION

An application of POS recognition algorithm is in Entity-Relation diagrams automatic construction.

To create an Entity-Relation diagram from a specification written in a specified natural language, first step is to identify the nouns - Entities and verbs - Relations.

For instance, for specification: *A driver drives cars* the following result is obtained:

Word : driver Driver#n#1 : driver

Word : drives Drives#v#28 : drive, take

Word : cars Cars#n#2 : car, railcar, railwaycar, railroadcar

The corresponding Entity-Relation Diagram will have as Entities nouns: *driver* and *car*, and as Relation the verb *drives*.

Next step is to identify Attributes for the Entities. This can be done by searching for the adjectives which determine the noun, and by searching for the nouns bounded to the Entity noun with an membership verb, as in the following: *APeoplehasanAddress*. *Address* is an Attribute for the Entity *Pupil*.

5. A SHORT COMPRESSION TO ANOTHER RELATED METHODS

Most of the methods found in the literature that deals with part of speech identification, or word tagging, are designed only for a single language. They use grammatical notions, corpuses and artificial intelligence techniques such as training on a part of the target text. They have very good results: more than 90% precision.

The advantage of the proposed method is that it is not using grammatical concepts, and so, it can be used for more than one language. And because of this, if a bilingual dictionary is available, the algorithm can be used to identify the language of every word.

6. CONCLUSION AND FURTHER WORK

POS recognition algorithm can be used to identify the POS(s) of words from a text written in a specified language if there is an electronically dictionary, without using grammar notions or another sources.

The precision for POS recognition can be improved. For instance, if at the beginning all the words with only one POS and only one sense are annotated, and than the CHAIN algorithm is applied ascending and descending, using the words with only one POS and sense as anchors. Another way is by not ignoring some POS(s), and this can be easily done by considering all the POS(s) found in the dictionary for the corresponding language. Another possibility is to start with the word with the longest gloss. And of course, combining the Chain algorithm with an artificial intelligence technique, to maximize the sum of all scores.

The POS recognition algorithm is a consequence of CHAIN algorithm for word sense disambiguation found in [7]. I believe that the CHAIN algorithm can have some more interesting consequences.

7. ANNEX

See Figure 2.

8. REFERENCES

- [1] S. Banarjee and T. Pedersen. 2003. *"Extended Gloss Overlaps as a Measure of Semantic Relatedness."* Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, August 9-15, Acapulco, Mexico.
- [2] E. Agirre and P. Edmonds (editors). 2006. *"WSD: Algorithms and Applications."* Springer.
- [3] C. Fellbaum (editor). 1998. *"WordNet An Electronic Lexical Database."* The MIT Press.
- [4] D. Jurafsky and J. Martin. 2000. *"Speech and language processing."* Prentice Hall.
- [5] C. Manning and H. Schutze. 1999. *"Foundation of statistical natural language processing."* MIT.
- [6] D. Tatar and G. Serban. 2001. *"A new algorithm for WSD."* Studia Univ. Babeş-Bolyai, Informatica, 2, 99108.
- [7] D. Tatar, G. Serban, A. Mihis, M. Lupea and M. Frentiu. *"A chain dictionary method for Word Sense Disambiguation and applications"*, Proceedings of KEPT2007, to appear.
- [8] <http://wordnet.princeton.edu/perl/webwn>

⁽¹⁾ COMPUTER SCIENCE DEPARTMENT, BABES-BOLYAI UNIVERSITY, KOGALNICEANU STREET NR. 1, RO-400084, CLUJ-NAPOCA, ROMANIA
E-mail address: mihis@nessie.cs.ubbcluj.ro

