

HOW MUCH WILL A RE-BASED PREPROCESSOR HELP A STATISTICAL PARSER?

Daniel Zeman

Centrum počítačn lingvistiky

Univerzita Karlova

Malostransk nmst 25, CZ-11800 Praha, Czechia

zeman@ufal.ms.mff.cuni.cz

Abstract

The present paper describes experiments on combining a statistical parser with a preprocessor built of regular expressions. While the syntactic structure of a sentence may be very complex and will never be fully covered by a finite-state machine, there are phenomena such as noun phrases and simple coordinations that do quite well. Even these “chunks” may theoretically be unlimitedly complex but we show that in reality they rarely do. So a shallow finite state parser can be built to preprocess the input text and solve the simple chunks without introducing too many errors. We discuss one implementation of such preprocessor that is very easy to write and covers roughly 20% of input words with an accuracy of over 90%. Then we describe two ways of combining the preprocessor with a parser and show that the performance of the parser improves both in speed and accuracy.

1. Introduction

The goal of the work reported on in this paper was to improve performance of a statistical parser. The idea was that there were very simple phenomena in the language that could be captured by a simple tool — where the second “simple” concerns programming effort as well as computational complexity. Combining such a tool with the parser would eventually speed up the parsing process. Moreover, the hope was that it would increase the parsing precision as well because the precision of the preprocessor was pretty high. The question was how much the preprocessor would increase the overall precision: due to its low recall it was expectable that most of structures it would recognize would also be found by the parser.

The parser works in a dependency framework. Its task is, for each input word to find another word from the same sentence that is its governor. Thus the dependency trees produced by the parser contain no nonterminals; rather, every node except of the root corresponds to a word of the sentence. There is no grammar (context-free or other) generating directly the dependency trees. While parsers employing probabilistic context-free grammars can be used only indirectly in this framework (see Collins et al., 1999, for an example), there is a more straightforward way similar to n-gram language modeling. It differs from the n-grams in that it collects statistics about pairs of dependent nodes rather than n-tuples of neighboring words. The parsed language material is Czech.

The preprocessor is a finite state tool implemented fully as a set of regular expressions in Perl. While the parser can be trained on any language (as long as it has a dependency treebank), the preprocessor benefits from some properties of Czech (but it can be applied to other similar – e.g. Slavic – languages). It exploits the fact that the governing and the dependent node often neighbor in the sentence and agree in gender, number, case etc. This property simplifies recognizing of base noun phrases and other patterns.

2. The statistical parser

The parser uses a dependency treebank to learn dependencies between words. It is not lexicalized which means that it only relies on the morphological properties of words, encoded in morphological (POS) tags. So we have the tag on mind when we talk about a ‘word’ if not stated otherwise.

The parser maintains a table of all dependencies between words of particular morphological category (e.g. between a masculine singular noun in nominative case (NNMS1---A-----) and a masculine singular adjective in nominative case (AAMS1---A1-----)) together with their relative frequencies. When applied to a raw text it first determines the possible dependencies. Then it estimates the probabilities of the dependen-

cies as their relative frequencies. And then it performs a Viterbi search to find a tree with the highest possible product of dependency probabilities.

This simple mechanism itself is too weak to build correct dependency structures. Several constraints have to be defined to make the parser work. One of such constraints is projectivity.

2.1. Projectivity

Projectivity is a property that combines tree structure and the word order in sentence. A dependency A-B (where A is the governing node) is **projective** if and only if all the words that are placed between A and B are included in the subtree of A. If you display the tree so that the x-coordinates of nodes correspond to the word order, each non-projective dependency will cross at least one perpendicular from another node (it will not necessarily cross another dependency.)

Projectivity is an important attribute of the word order. It does not depend on the length of the phrases but it still allows us to distinguish “normal” sentence structures from the “wild” ones. We deliberately do not call them “incorrect” because not all non-projective dependencies are errors in Czech.¹ Nevertheless such dependencies are quite rare: we counted only 1.8 % of dependencies in a 19126-sentences corpus being non-projective,² so it seems to be useful to require that all the dependencies generated by the parser are projective.

2.2. Reduced set of morphological tags

Some of the morphological information encoded in the tags has no influence on syntax. For instance, all adjectives, verbs and adverbs can have positive and negative forms, whereas the negative form is constructed completely regularly by simply inserting the prefix “ne-” to the word beginning. The negative forms behave syntactically exactly the same way the positive ones do, so using such information in syntax training may damage the results. Suppose we have seen a positive form of a word depending on a given other word n -times but we have not ever seen the negative form in this context. Then the edge with the negative form will get a probability of zero (or close to zero when smoothed) but in fact it should have the same probability as the one with positive form. There are many other tags that have been merged together. The original tagset contained over 3000 possible tags, 1279 really appeared in the data. After the reduction this number decreased to about 452 so the reduction rate is approximately 65 %.

2.3. Modeling number of child nodes

Some of the words tend to have many dependents while some others are almost always leaves, and the prepositions take mostly exactly one dependent. The parser described so far does not reflect these observations and, e.g., if there is one preposition in the sentence but several nouns, it often wants to hang all the nouns to the single preposition. To avoid this, an additional model is introduced which knows how often a given tag has a given number of children (0, 1, 2, or 3+). The parser multiplies the edge probability estimations by the probability of the governing node having $n+1$ children where n is the current number of children.

2.4. Direction and distance

Finally, there are two important features that reflect the mutual positions of the two dependency members in the sentence. The first one is **direction**. The model assigns separate probabilities to a dependency where the depending node is (in the sentence) to the left of the governing node, and to the same dependency where it is to the right of the governing node. For instance, a preposition stands always to the left of its dependent.

The second thing is distance, or better, **adjacency**. Now the parser asks whether the two nodes are adjacent in the sentence or not. Separate probability estimates are kept for both options.

3. The finite state preprocessor

We used regular expressions to model the local syntax of coordinations, simple noun phrases and other simple small chunks of sentences. The method is closely related to *local grammars* that have been applied to a similar task for a similar language (Serbo-Croatian; see Nenadić and Vitas 1998, Nenadić 2000).

¹ Note however that non-projective constructions cannot be described by an ordinary context-free grammar.

² The corpus we used to measure non-projectivity was a part of the Prague Dependency Treebank (PDT, see Böhmová et al.). The average sentence length in this part is 16.3 words (dependencies). About 79.4 % of all trees have all dependencies projective. The number of trees that contain one or no crossing (non-projective) dependency is 93.8 % and the number of trees with at most two such dependencies is 98.3 %.

The preprocessor is written in Perl (see Wall et al., 1996). There is a procedure that scans the corpus and repeatedly applies a set of regular expressions (RE's) to each sentence. The RE's can be easily modified and new RE's can be added. Every RE describes a local configuration that signalizes a dependency between two chunks of the input. A chunk is a word or a sequence of chunks that has been recognized by a previous application of an RE (and replaced by a representative word).

For instance, a coordination of adjectives is defined by the following three sub expressions:

```
"<t>A. ([^X]) ([^X]) (\\d) "  
"<l>(a|i|nebo|ani)<"  
"<t>A. \\1\\2\\3"
```

The Perl procedure combines these sub expressions into one RE. The resulting RE finds a sequence of words such that the description of the first word contains a match of the first sub-RE, the second word description matches the second sub-RE and so on. So in this case we are looking for a sequence of three words where the first one has a morphological tag beginning with A (adjective), having any character at second position (sub-part-of-speech), having known gender and number (unknown is encoded as X, which is prohibited by the RE), and having known case encoded as a digit. The second word must have the lemma *a* ("and"), *i* ("as well as"), *nebo* ("or") or *ani* ("nor"). The third word must again be an adjective (the tag beginning with A) and must agree with the first one in gender, number and case. The agreement is enforced by the sub-RE's \1, \2, and \3. They refer to the first, the second and the third parenthesized sub-RE's, respectively. In our case all referred parentheses occurred in the first word's tag, marking its gender, number and case.

The example RE would thus match for instance the coordination *přímým i nepřímým* ("direct as well as indirect"; instrumental case) that has the following (simplified) representation in PDT:

```
<f>přímým<l>přímý<t>AAMS7---A1-----  
<f>i<l>i<t>J^-----  
<f>nepřímým<l>přímý<t>AAMS7---N1-----
```

On the other hand, the RE would not match the sequence *žlutého a černou* ("yellow / Masc. and black / Fem.") because the adjectives do not agree in gender:

```
<f>žlutého<l>žlutý<t>AAMS4---A1-----  
<f>a<l>a<t>J^-----  
<f>černou<l>černý<t>AAMS4---A1-----
```

Two special parameters are passed to the procedure. They are interpreted as relative indices of the phrase head and of the phrase representative (1 and 2 in our case, the first word has the index of zero).

Each RE is applied repeatedly to each sentence until all matches are found. This can also accommodate some recursive constructions. Suppose we have an RE for adjective-noun NP's. And suppose there is a noun preceded by two adjectives that both modify it, e.g. *velký zelený strom* ("big green tree"). First the sequence *zelený strom* is recognized and replaced by the representative noun *strom*. This way the words *velký* and *strom* become neighbors and the sequence *zelený strom* can be recognized in the next run.

To keep the process simple, the order in which different RE's are applied is fixed. Nevertheless some RE's can be tried at several points to see whether the replacements created new material for them. To illustrate this, let's have two RE's called AN (adjective-noun NP), and NaN (coordination of nouns), and let's apply them in the order NaN-AN-NaN. Then if the input contains the sequence *čeští sportovci a slovenští sportovci a umělci* ("Czech sportsmen and Slovak sportsmen and artists"), the first NaN reduces it to *čeští sportovci a slovenští umělci*, this sequence is further reduced by the AN to *sportovci a umělci*, which is finally recognized by the second NaN. This example shows that the covered constructions can be richer than one might think. On the other hand it also reveals the main weakness of the approach: in other configurations such as *čeští hráči i Němci* ("Czech players as well as the Germans"), the RE sequence AN-NaN would be appropriate. So the method can be relatively successful only if there exists an ordering of RE's that matches the data in much more cases than any other ordering.

The programming simplicity and efficiency is a significant issue. The driving procedure was finished in a few days, a regular expression can be formulated and added in a minute or two and the Perl interpreter applies it to a corpus of 19126 sentences often in less than one minute³. On the other hand it is obvious that one can never cover the language this way, partly because a human cannot think of all possible configurations of the data, and partly because adding more complex patterns would make the expressions unintelligible and unmaintainable. That's why the preprocessor can only *help* the parser, not replace it.

In the rest of this section we give examples of patterns that the regular expressions recognize.

³ On a 266 MHz PC running Linux.

3.1. Coordinations

We found 55 different words and punctuation marks that take the function of coordinating conjunctions in our corpus. Not all of them occur in the same pattern and not all are suitable for the finite state tool. Some coordinating elements such as the comma occur very frequently in non-coordinated configurations so it is better to rely on the statistical parser's reasoning what is a coordination and what is not. However many others are most likely seen in simple noun phrase patterns that are easily described by regular expressions.

As an example, we selected the pattern "word coord word" where `word` is a noun or an adjective and `coord` is one of the conjunctions *a (and)*, *i (as well as)*, *nebo (or)*, *ani (nor)*. The pattern is described by two RE's, called NaN and AaA. The first RE models coordinations of nouns, the second coordinations of adjectives. AaA requires that the coordinated members agree in gender, number and case. NaN requires only agreement in case. Only simple coordinations of two elements were tested but an RE for larger coordinations (e.g. *Peter, Paul and Mary*) could easily be created and added. One RE could even capture coordinations of arbitrary (unlimited) size. However, large noun coordinations are quite rare in the real data, so even our toy RE can solve a lot.

Of course the coordination processor will be more efficient if we are able to process simple noun or adjective phrases first. See below the description of the respective RE's. We searched the text for DA's (adverb-adjective pairs) before AaA, and for AN's (adjective-noun pairs) before NaN.

3.2. Adjective – noun NP's

The following RE, called AN, recognizes noun phrases that consist of an adjective followed by a noun; the adjective agrees with the noun in gender, number and case. The noun is the head and the representative of the phrase.

```
"<t>A. ([^X]) ([^X]) (\\d) "  
"<t>N. \\1\\2\\3 "
```

3.3. Adverb – adjective AP's

The following RE, called DA, recognizes adjective phrases consisting of an adjective preceded by an adverb.

```
"<t>Dg "  
"<t>A "
```

3.4. Noun – genitive noun NP's

One common way of forming noun phrases in Czech (and other Slavic languages — see Nenadić (1998) for Serbo-Croatian) is accompanying an existing NP with another NP in genitive case. The genitive NP has a similar function as a PP with the preposition *of* in English. So, for instance *nový prezidentský mandát Václava Havla* means “new presidential term of Václav Havel”. The whole NP consists of two NP's, *nový prezidentský mandát* and *Václava Havla*. The members of each of these NP's agree in gender, number and case. However, the second NP is frozen in genitive regardless in which case the first NP is.

Theoretically, an unlimited number of genitives can be recursively cumulated this way — consider phrases like “the break-down of the new car of the father of a friend of mine”. We need to recognize the rightmost pair of genitives first to get the whole tree correct. So the RE recognizing two consecutive genitives must make sure that there are no more genitives to the right. The (obvious) base RE called NgNg is as follows:

```
"<t>N. . . 2 "  
"<t>N. . . 2 "
```

After the procedure creates one RE from the above pieces, we must add the following negative look-ahead condition. It assures that the next word is not a noun in genitive:

```
" (?! .*? \\n. *?<t>N. . . 2. *? \\n) "
```

When all genitive pairs are found and collapsed, the noun-in-any-case – noun-genitive pairs can collapse. This is done by the RE called NNg:⁴

```
"<t>N "  
"<t>N. . . 2 "
```

⁴ Both RE's could be combined into one (NNg with the look-ahead condition).

3.5. Prepositional phrases⁵

The following RE, called RN, combines prepositions with nouns whose case is compatible with the respective preposition:

```
"<t>R...(\d)"
"<t>N...1"
```

4. Combining the regular expressions with the parser

There are two or three ways how to combine the two parsing tools. The first one assumes that the preprocessor would replace each recognized phrase by a representative word (typically by its head but coordinations – whose head is a conjunction – should be represented by one of their members, converted to plural). The parser then would see the preprocessed phrases neither during the training phase, nor during testing. Such phrases become atomic items for the parser.

There is a possible drawback of this method. If the preprocessor fails to find all members of a phrase, the error can be corrected by the parser only if the forgotten member depends on the head of the phrase. If it ought to be nested more deeply in the phrase, its real governor is now invisible because the phrase is atomic. This leads to the second approach where the phrase would even for the parser be a structure rather than a monolith. The parser would get some dependencies for free but would be allowed to add new dependencies at any place in the structure. The third approach is a special case of the second one. It applies the finite state tool as a postprocessor to the statistical parser output, overriding all parser decisions concerning the recognized chunks.

We refer to the three methods later in this paper as to *transparent preprocessing*, *non-transparent preprocessing*, and *postprocessing* respectively.

5. Results

We ran several experiments to evaluate the finite state tool and its contribution to the performance of the parser. For each regular expression, we figured out its precision (measured as the number of correctly proposed dependencies divided by the total number of proposed dependencies). The test was done on 19126 sentences of the Prague Dependency Treebank (Böhmová et al.); this portion contains 346719 words (tokens). The results are shown in the following table:

Expression	AaA	NaN	AN	DA	NgNg	NNg	RN
Precision	99.2	95.0	98.8	86.8	91.9		95.7

We tested all described RE's in the following order (note that NaN was applied twice):

DA-AaA-AN-NaN-NgNg-NNg-NaN-RN

Also this combined RE tool was tested on the same corpus. It proposed 69262 correct dependencies and 4424 false ones while leaving 273033 tokens unresolved. So the recall is 20.0 %, which means that the preprocessor is able to (correctly) help in one fifth of cases. The precision of the preprocessor is 94.0 %; only 1.3 % of all dependencies are errors that the parser cannot repair.

The above results are for data where morphological tags were assigned manually. Such evaluation is important because it shows the power of the RE's without biasing it by errors of other components. However, an application will hardly have manually annotated data available, so we are adding two other statistics. First, the RE set was run on tags and lemmas assigned by a statistical tagger (Hajič and Hladká, 1998). Every <t> in RE's was automatically replaced by <MDt [^] * >; the <l> markups were treated similarly. The resulting machine proposed 60093 good dependencies and 5001 bad dependencies, 281625 remained unresolved; the precision was 92.3 %, the recall 17.3 %. The decrease in both the precision and the recall can be explained by the errors the tagger does.

The last test used ambiguous morphological input. Every <t> in RE's was automatically replaced by <MMt>; the <l> markups were treated similarly. The resulting machine proposed 69566 good dependencies and 8146 bad dependencies, 269007 remained unresolved; the precision was 89.5 %, the recall 20.1 %. The preprocessor had all morphological hypotheses available, so the recall slightly increased. But occasionally it combined two compatible hypotheses that occurred at neighboring positions but were wrong. That explains the drop in precision.

⁵ Some more statistics of PP patterns in Czech – but in other corpus – are in (Petkevič, 1999).

Finally we combined the finite state tool with the parser in both the ways described in Section 4. The resulting systems were tested on a part of the PDT that had not been used to train the parser. This testing portion consisted of 3697 sentences. Every word is now assigned a dependency so the precision is equal to the recall (and we call the measure *accuracy*). All experiments used machine-disambiguated morphology. The following table shows the baseline performance of the parser itself, the performance of RE's where unresolved words would always be attached to the root, and the performances of parser with RE transparent preprocessor, non-transparent preprocessor and the postprocessor. Note also the times in the second line: they demonstrate how the parsing speed increases when some chunks are preprocessed.

	Parser	Sole RE's	Transp. Prepr.	Nontr. Prepr.	Postprocessed
Accuracy (%)	53.7	30.6	55.6	57.2	56.6
Time (min)	53	0	30	46	54

6. Conclusion

We briefly presented a statistical dependency parser and showed that a set of regular expressions helps the parser both in speed and accuracy. While the REs' contribution to the accuracy is to be expected to decrease as the parser itself improves, it hardly ever will damage the parsing results because the precision of the REs is very high – higher than of any known parser for Czech or English. At the same time the parsing is speeded up substantially which should not change with further development of the parser. The regular expressions can be written very easily and quickly.

7. Acknowledgements

This research is being supported by the Ministry of Education of the Czech Republic project No. LN00A063 (Center for Computational Linguistics). The data (PDT) would not be available unless the grant No. 405/96/K214, of the Grant Agency of the Czech Republic, supported the treebank design.

8. References

- Alena Böhmová, Jan Hajič, Eva Hajičová, Barbora Hladká (in press). The Prague Dependency Treebank: Three-Level Annotation Scenario. In: Anne Abeillé (Ed.): *Treebanks: Building and Using Syntactically Annotated Corpora*. Kluwer Academic Publishers, Dordrecht, The Netherlands. See also <http://ufal.mff.cuni.cz/pdt/>.
- Michael Collins, Jan Hajič, Eric Brill, Lance Ramshaw, Christoph Tillmann (1999). A Statistical Parser of Czech. In: *Proceedings of ACL 1999*, pp. 505–512, College Park, Maryland.
- Jan Hajič, Eric Brill, Michael Collins, Barbora Hladká, Douglas Jones, Cynthia Kuo, Lance Ramshaw, Oren Schwartz, Christoph Tillmann, Daniel Zeman (1998). *Core Natural Language Processing Technology Applicable to Multiple Languages. The Workshop 98 Final Report*. At: <http://www.cisp.jhu.edu/ws98/projects/nlp/report/>. JHU, Baltimore, Maryland.
- Jan Hajič, Barbora Hladká (1998). Tagging Inflective Languages: Prediction of Morphological Categories for a Rich, Structured Tagset. In: *Proceedings of the 36th Meeting of the ACL and COLING'98*, pp. 483–490. Université de Montréal, Montréal, Québec.
- Goran Nenadić, Duško Vitas (1998). Using Local Grammars for Agreement Modeling in Highly Inflective Languages. In: *Proceedings of TSD 1998*, pp. 91–96, Brno, Czechia.
- Goran Nenadić (2000). Local Grammars and Parsing Coordination of Nouns in Serbo-Croatian. In: *Proceedings of TSD 2000*, pp. 57–62, Springer LNAI 1902, Brno, Czechia.
- Vladimír Petkevič (1999). Czech Translation of G. Orwell's '1984': Morphology and Syntactic Patterns in the Corpus. In: *Proceedings of TSD 2000*, pp. 77–82, Mariánské Lázně, Czechia.
- Larry Wall, Tom Christiansen, Randal Schwartz (1996). *Programming Perl*. <http://www.perl.org/>.
- Daniel Zeman (1998). *A Statistical Approach to Parsing of Czech*. In: Prague Bulletin of Mathematical Linguistics, vol. 69, pp. 29–37. Univerzita Karlova, Praha, Czechia.
- Daniel Zeman (2001). *Parsing with Regular Expressions: A Minute to Learn, A Lifetime to Master*. In: Prague Bulletin of Mathematical Linguistics, vol. 75. Univerzita Karlova, Praha, Czechia.