# Part-Of-Speech Taggers Features In French Learner Texts

Nadezhda Barymova
Petrozavodsk State University
Petrozavodsk, Russia
nbarymova@gmail.com

Anna Zhestkova
Petrozavodsk State University
Petrozavodsk, Russia
zhest1806anna@gmail.com

Nadezhda Oulianova
Petrozavodsk State University
Petrozavodsk, Russia
nad.ulyanova2013@yandex.ru

Olga Nikiforova
Petrozavodsk State University
Petrozavodsk, Russia
o_nik75@rambler.ru

*Abstract*—**This paper presents investigation result of the comparative analysis of TreeTagger, CoreNLP and SpaCy as part-of-speech taggers for processing texts in French. Acting in the framework of a Learner Texts Corpus Creation Project a group of researchers from the Institute of Foreign Languages together with the Institute of Mathematics and Information Technologies of Petrozavodsk State University analyzed the above mentioned tagging tools mainly being focused on defining the most effective one.**

## I. INTRODUCTION

Corpus technologies are actively developing and find their application in various fields, from linguistic researches to their use in the foreign languages teaching techniques. [1], [2]

The most important component of the corpus in any language is a tagset, which allows the morphological analyzer (morphologizer) automatically process linguistic objects in a formalized form [3].

The main objective of our research is to select the most suitable tagger for processing a text in French. To tackle this problem we:

- analyzed 30 learner texts in French using the most famous and affordable taggers TreeTagger, CoreNLP and SpaCy.
- highlighted the main errors (separation of words, highlighting parts of speech, etc.) made by each of the specified taggers.
-compared the tagging results and highlight the tagger with the least number of errors
- selected the most appropriate tagger to use in work with the corpus of learner texts in French using theanalysis results

## II. MAIN PART

During the research work we analyzed 30 learner texts of different types such as essays, motivation letters, topics, article evaluation using TreeTagger, CoreNLP and SpaCy taggers. All the texts under analysis were written by the students of the Institute of Foreign Languages of Petrozavodsk State University. The analysis of the texts made it possible to obtain each tagger tools comprehensive evaluation.

TreeTagger [4] is the only tool with tags in French among the three part-of-speech taggers under consideration which proves to be its main advantage comparing to other tools. Considering this part-of-speech tagger as a possible version of a tool for working with texts, it should be taken into account that it is not very efficient from the point of view of visual perception.

The reasons for that can be mentioned as the follows: first, the tagger font color is black, second, the tags are located directly above the words making the text less readable so you have to read the whole text while paying a special attention to the superscripts, and finally, the body of the text is not split into separate sentences, which also complicates the process of checking the tagger.



Fig.1. Screenshot of the tagged text

CoreNLP - one of the main advantages of this tagger is its visual presentation. So due to this peculiarity all the sentences in the text beginning on a new line are also numbered. This

feature simplifies the word search. Moreover, this immediately allows the reader to spot/identify/ the sentence types presented in the text. At that tags for tokens are highlighted in colour, which makes it more convenient to define parts of speech.

Nevertheless it should be noted that the texts tagging is done in English making the analysis on the one hand more simple, but imprecise at the same time.



Fig.2. Screenshot of the tagged text

SpaCy

When working with this tagger, you need to take into account the fact that its list of tags in English is the same as that one of the CoreNLP tagger on the one hand, at the same time the total error amount is bigger than the number in the above described part- of-speech taggers.
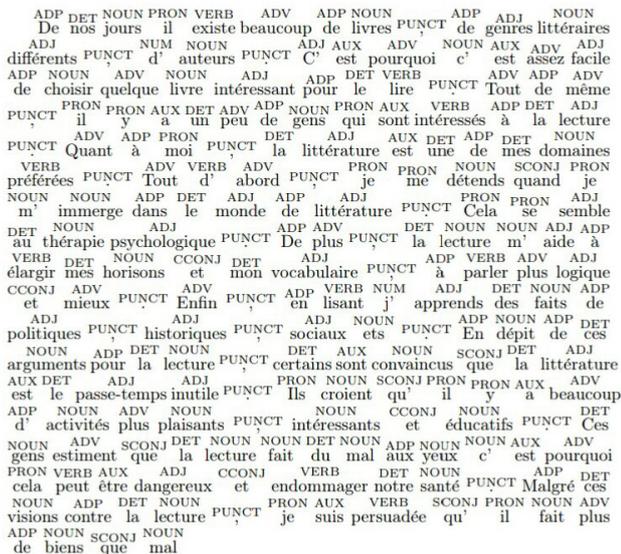


Fig.3. Screenshot of the tagged text

It should be also noted that the number of tokens in the same text varied depending on the used part-of-speech tagger. The data is presented in diagram 1.
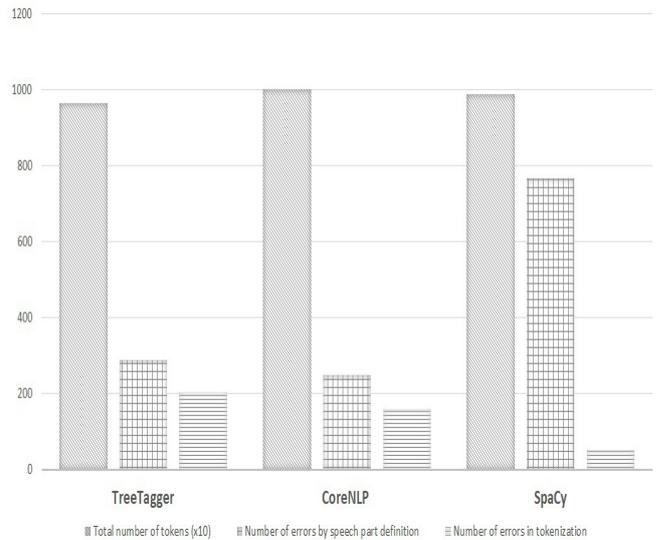


Fig.3. Diagram I

So we found out the indicator of the total number of tokens in the TreeTagger was the lowest for the errors in tokenization are often encountered due to abbreviated articles and pronouns which are written with an apostrophe followed by a word ending in a vowel and which are not separated into certain tokens. As you can see in Table II, TreeTagger, compared to CoreNLP and SpaCy, makes the most of word separating errors (2.1% versus 1.6% and 0.5%).

CoreNLP, on the contrary, splits parts of merged articles into different tokens, therefore, the total number of tokens is much higher than in other taggers. In SpaCy, the percentage of errors in the division of words with the total number of tokens is the smallest, based on it, we can conclude that this indicator is more accurate in this tagger. Thus, if tokenization is the prevailing goal of the text analysis, then SpaCy can be effective in the process of choosing a part-of-speech tagger.

TABLE II .POS TAGGERS COMPARISON

| Percent /indicators (of the total number of tokens) | TreeTagger | CoreNLP | SpaCy |
|---|---|---|---|
| Percent of errors by speech part definition | 3% | 2,5% | 7,8% |
| Percent of errors in tokenization | 2,1% | 1,6% | 0,5% |
| Total percent of errors | 5,1% | 4,1% | 8,3% |

Therefore, with this error percent criterion (determining parts of speech in 3 taggers) TreeTagger and CoreNLP have

similar indicators: 3% and 2.5%. Whereas in SpaCy error percent is 7.8%.

The lowest percent of the total number of errors and the number of tokens, 4.1% is characteristic for CoreNLP tagger, as for the highest one it is for SpaCy namely 8,3% ; and as to TreeTagger this percent is 5.1%. Consequently, the TreeTagger and CoreNLP part-of-speech taggers are more accurate in defining part of speech and allocating tokens.

We have compiled the following diagram III and TABLE III with marked erorrs using one of the texts (essays) analyzed through all three taggers as an example.

TABLE III

| Word | Errors by speech part definition | Correction | Total number | Errors in tokenization | Total number |
|---|---|---|---|---|---|
| **Treetagger:** | | | | | |
| d'auteurs | adj | prp+nom | 12 | d'auteurs | 8 |
| c'est | nam | pro:per+ver:pres | | c'est(x2) | |
| c'est | ver:subi | pro:per+ver:pres | | d'abord | |
| interressant | ver:ppre | adj | | m'immerge | |
| lire | nom | ver:infini | | m'aide | |
| m'immerge | nom | pro:per+ver:pres | | j'apprends | |
| m'aide | nom | pro:per+ver:pres | | qu'il | |
| (en) lisant | adj | ver:ppre | | d'activités | |
| j'apprends | nom | pro:per+ver:pres | | | |
| qu'il | nom | kon+pro:per | | | |
| d'activités | ver:pper | prp+nom | | | |
| c'est | nom | pro:per+ver:pres | | | |
| **CoreNLP:** | | | | | |
| nos | det | pron | 13 | au therapie | 4 |
| d' | det | adp | | des faits | |
| C' | propn | pron | | du mal | |
| est X4 | aux | verb | | aux yeux | |
| c' | noun | pron | | | |
| quelque | det | adj | | | |
| intéressé | verb | adj | | | |
| j' | det | pron | | | |
| apprends | noun | verb | | | |
| convaincus | verb | adj | | | |
| sont X2 | aux | verb | | | |
| notre | det | pron | | | |
| persuadé | verb | adj | | | |
| **SpaCy** | | | | | |
| genres | adj | noun | 32 | - | 0 |
| littéraires | noun | adj | | | |
| d' | num | adp | | | |
| C' | adj | pron | | | |
| C'(x2) | noun | pron | | | |
| choisir | noun | verb | | | |
| quelque | adv | pron | | | |
| (pour) le (lire) | det | | | | |
| Tout | adv | pron | | | |
| lecture | adj | noun | | | |
| Quant | adv | adp | | | |
| littérature(x3) | adj | noun | | | |
| d' | verb | adp | | | |
| abord | adv | noun | | | |
| détends | noun | verb | | | |
| m'(x2) | noun | pron | | | |
| immerge | noun | verb | | | |
| monde | adj | noun | | | |
| semble | adj | verb | | | |

III. CONCLUSION

Due to our research analysis of TreeTagger, CoreNLP and SpaCy as part-of-speech taggers for processing texts in French we chose the CoreNLP tagger as the most effective tool in use for our project. CoreNLP proved to have almost the same number of errors in determining parts of speech, moreover, it

made fewer errors in separating words and it turned out to be more illustrative at work tan other tools.

Finally, it can be summarized as follows: using TreeTagger, CoreNLP and SpaCy universal taggers it is necessary to take into account that every one of them has its merits and shortcomings in terms of various kinds of errors. So when choosing a tagger the latter must always be carefully considered.

In 2020 the Department of German and French Languages of the Institute of Foreign Languages and the Center of Artificial Intelligence of Petrozavodsk State University initiative was the starting point for learner texts corpus creation project in German and French languages with the subsequent development of various virtual applications based on it. This project on the one hand, provided teachers of a foreign language with efficient tools of checking a large number of written works, on the other hand, it gave the students the opportunity to improve their writing skills while completing student assignments [5].
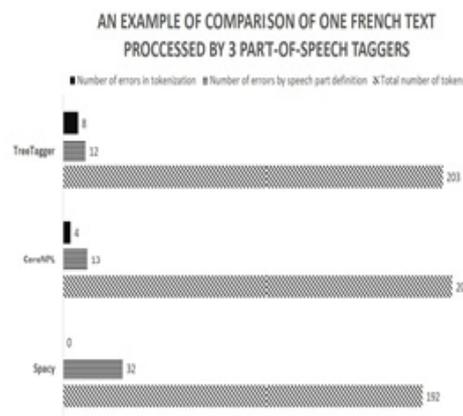


AN EXAMPLE OF COMPARISON OF ONE FRENCH TEXT PROCCESSED BY 3 PART-OF-SPEECH TAGGERS

Fig.4. Diagram III

IV ACKNOWLEDGMENT

REFERENCES

[1] Nagel O.V. "Crpus Linguistics and its use in computer-based language teaching"//Language and culture – 2008 - № 4. Web : HTTPS://WWW.ELIBRARY.RU/ITEM.ASP?ID=11990991
[2] Maltseva M.S. Corpus technologies in foreign languages teaching methodology// Socio-economic phenomena and processes- 2011-№8. Web: HTTPS://CYBERLENINKA.RU/ARTICLE/N/KORPUSNYE-TEHNOLOGII-V-METODIKE-PREPODAVANIYA-INOSTRANNYH-YAZYKOV
[3] Julien Fago. Conception et réalisation d'une chaîne de traitement automatique des langues adaptée à des projets littéraires. Web : HTTPS://DUMAS.CCSD.CNRS.FR/DUMAS-02987314/DOCUMENT
[4] French TreeTagger part-of-speech tagset Web https://www.sketchengine.eu/french-treetagger-part-of-speech-tagset/
[5] I.A. Kotyurova, "Sozdanie korpusov uchebnykh tekstov kak razvivayushheesya napravlenie korpusnoj lingvistiki", *International Scientific Journal*, №5, 2020, pp. 100-109.