

Russian Computational Linguistics: Topical Structure in 2007-2017 Conference Papers

Amir Bakarov Andrey Kutuzov Irina Nikishina

Dialogue
Conference on Computational Linguistics and Intelligent Technologies

02.06.2018

Aim and Contribution

Part of a larger project to analyze Russian NLP publishing activity (RusNLP).

- ▶ **Analysis of Russian NLP publications topical structure and topical drift** (for papers written in English):
 - ▶ between venues;
 - ▶ diachronically (2007-2017).
- ▶ A large **dataset of papers published at Russian NLP venues**
 - ▶ considerable amount of **metadata**, available in a machine-readable format

Setup

- ▶ 3 primary Russian NLP venues;
 1. AINL;
 2. AIST;
 3. Dialogue.
- ▶ More than 1700 papers were manually assessed to parse metadata;
- ▶ LDA with varying hyperparameters was trained on this dataset to extract the topics.

Web service is available at <http://nlp.rusvectors.org> (work in progress!).

Primary Findings

Not surprising:

1. Dialogue is a more linguistically-oriented event, while AIST and AINL are more computer science oriented;
2. The trending topics of the Russian NLP in general are determined by the central topics of 'Dialogue';
3. The overall topical structure is shifting towards machine learning domain.

Come to Our Poster to Learn More!

Russian Computational Linguistics: Topical Structure in 2007-2017 Conference Papers

Amir Bakarov¹ Andrey Kuzrov¹ Inna Nikishina^{1,2}

¹National Research University Higher School of Economics,
²Federal Research Center "Computer Sciences and Control" of the Russian Academy of Sciences,
Moscow, Russia

¹University of Oulu,
Oulu, Norway

¹National Research University Higher School of Economics,
²Novosibirsk Institute for System Programming of the Russian Academy of Sciences,
Moscow, Russia

amirb@econ.yandex.ru, andk@it1.uio.no, inna.nikishina@phs.hse.ru

1. Contributions

Part of a larger project to analyze Russian NLP publishing activity (RusNLP)

- Analysis of Russian NLP publications topic structure and topical shift (the papers written in English)
- Online system
- Visualization (2017-17)
- A large dataset of papers published at Russian NLP events
- A considerable amount of statistics available in our research forum

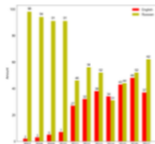


Figure 1. Yearly number of Dialog papers in Russian and English.

2. Approach

Latent Dirichlet Allocation (LDA) with 5 topics

3. Web service (link to preprint)

to <http://rlp.metacrawler.org>
Research trends topics based on RusNLP

Year Topics

Year	Topics
2007	1. grammar, language, use, Slavic, Croatian, Serbian
2011	1. standard, adjective, degree, functional, construction, purpose 2. word, language, rule, text, structure, result 3. rule, word, language, school, state, letter 4. clear, emotional, clause, expression, derivative, word 5. word, text, algorithm, result, collection, chat
2015	1. feature, word, result, result, word, post 2. paragraph, feature, task, result, sentence, word 3. word, result, task, understanding, result, sentence 4. feature, word, use, bridge, language, domain
2017	1. feature, word, result, result, word, post 2. paragraph, feature, task, result, sentence, word 3. word, result, task, understanding, result, sentence 4. feature, word, use, bridge, language, domain



Conference Topics

- Dialog**
- subject, verb, case, language agreement, discourse
 - speech, language, case, word, Russian, text
 - translation, Russian, text, language, result, form
 - word, result, result, feature, Russian, discourse
 - text, language, Russian, Czech, English, relation
- AST**
- topic, result, word, topic, feature, discourse
 - word, analysis, morphological, language, error, Russian
 - word, sentence, result, result, text, noun
 - word, sentence, feature, use, result, based
 - word, language, clause, question, type, word
- ANL**
- word, result, task, domain, result, paragraph
 - system, domain, text, service, result, user
 - feature, word, semantic, sentence, result, paragraph
 - result, speech, language, word, recognition, sentence
 - feature, algorithm, domain, feature, number, result

Table 1. 5 most publications clustered by 5 words for each topic



Figure 2. Topical similarity for papers in different years.

5. Systems

- Dialoger <http://www.dialog-21.ru/>
- AST (RusNLP only): <http://nlp.metacrawler.org/>
- ANL (RusNLP only): <http://nlp.metacrawler.org/>



Figure 3. The yearly number of papers included for each system

6. Primary findings

For example:

- Identical topic in a main language only is not event, while AST and ANL are more complex sentence oriented
- The broader topics of the Russian NLP is great as demonstrated by the result topic of Dialoger
- The overall topical structure is shifting from its main focus (learning domain)

Acknowledgements

This study is sponsored by numerous NLPs and The Project is the only ones for Russian academic community to carry on, thanks to the state from Russian after other main projects frequently break result of academic researcher

References

1. Amir Bakarov, Andrey Kuzrov, Inna Nikishina. Russian Computational Linguistics: Topical Structure in 2007-2017 Conference Papers. In: Proceedings of the 2017 Conference on Empirical Natural Language Processing (EMNLP), 2017, pp. 1-11.

Thank you for your attention!

Amir Bakarov, Andrey Kutuzov, Irina Nikishina

amirbakarov@gmail.com, andreku@ifi.uio.no, irina.nikishina@mail.ru

<http://bakarov.github.io>

