(B)

2017　2018

Corpora on Demand: Scalable Methods of Obtaining Linguistic Data

Corpora on Demand: Scalable Methods of Obtaining Linguistic Data

Drozd, Aleksandr

1,400,000

(

)　　http://vecto.space/

This research project contributed to the development of methods of filtering textual data and evaluating text representations.

We have developed a set of evaluation benchmarks for text representations. These benchmarks can be used to estimate quality of text corpora, as well as methods for constructing representations themselves and corresponding hyperparameters.　We have implemented heuristics for filtering good quality text fragments and instrumental scripts to help parsing different formats of archived internet documents. Since internet documents might contain copyrighted and private information, we could not publish the raw corpora. Instead we are publishing all our source codes to extract raw texts from different archive formats as well as models (such as word embeddings) that we have trained on large scale texts. All resources and source codes are available at the project website http://vecto.space

natural language processing

corpora　evaluation　text representations

At the beginning of research project most of experiments in computational linguistics and natural language processing were done on a couple of easy-to-get corpora, mainly dumps of Wikipedia articles. Another problem is that NLP models were often compared using different corpora and on a small set of probing tasks.

The purpose of this research project was to help NLP and Cl practitioners to get training data easily and evaluate how particular choice of training corpora affects performance of downstream tasks.

In this research project we had to address both scientific and engineering challenges. We integrated existing and created new methods for filtering text and evaluation of text representations. Besides observing metrics like accuracy for evaluation tasks, we were also trying to investigate which linguistic phenomena are involved in the pipeline of each task and how downstream results reflect what upstream stages like model and source corpus capture in terms of linguistic nuances. We were trying to find and integrate existing solutions when possible and have made ourselves software implementations for novel ideas, including those proposed by us in a context of this project.

We have developed novel methods for evaluating text representations, including new algorithms to solve analogical reasoning task, new datasets for Japanese, English and Russian languages and reference models for various diagnostic tasks. These evaluations methods can be used to benchmark various stages of NLP pipelines, including choice or methods to create source corpora. We have created an approach to preserve metadata at all stages of NLP pipelines, so that impact of individual stages can be easily evaluated. We have created helper scripts to extract raw texts from sources like Wikipedia, project Gutenberg collection of books and Common Crawl archives. We have integrated some of existing solutions for detecting language, parsing HTML pages etc.

We also investigated linguistic-theoretical issues related to collection of textual data and evaluation of text representation and pre-trained model in transfer-learning scenarios.

We have integrated all our results in an open-source software library called Vecto. The library and resources are available at the http://vecto.space/ website. In Vecto library we dedicated special attention to the problem of reproducibility of results. We have implemented metadata collection and assimilation in all stages of supported NLP pipelines, including choice and collection of source corpora. Resources we published include pre-trained embeddings, reference implementation of evaluation models, datasets etc. User-base of vecto is continuously growing: we are regularly receiving feature requests, bug reports etc.

We could not publish the corpora we collected to avoid risk of copyright violation complains. Instead we have released models and embeddings which were pre-trained using corpora we have collected.

We have organized a tutorial ("Compositionality in the Age of Word Embeddings", collocated with LREC 2018 conference) and promoted some of our results and software in it. We have also organized a mini-symposium on convergence of HPC and AI, co-located with SIAM PP 2018 conference, in which one of the sessions was dedicated to the problem of scaling NLP models.

Collaborations established during this research continue after the expiration of the funding and we aspire to extend and improve existing results.

1) Marzena Karpinska, Bofang Li, Anna Rogers, and <u>Aleksandr Drozd</u>. Subcharacter information in japanese embeddings: when is it worth it? In In Proceedings of the Workshop on Relevance of Linguistic Structure in Neural Architectures for NLP (RELNLP) 2018, 28–37. ACL, 2018.

2) Bofang Li, <u>Aleksandr Drozd</u>, Tao Liu, and Xiaoyong Du. Subword-level composition functions for learning word embeddings. In Proceedings of The 2nd Workshop on Subword and Character level models in NLP (SCLeM), 38–48. ACL, 2018.

3) Bofang Li, Tao Liu, Zhe Zhao, Buzhou Tang, <u>Aleksandr Drozd</u>, Anna Rogers, and Xiaoyong Du. Investigating different syntactic context types and context representations for learning word embeddings. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, 2421–2431. 2017.

4) Anna Rogers, <u>Aleksandr Drozd</u>, and Bofang Li. The (too many) problems of analogical reasoning with word vectors. In Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (* SEM 2017), 135–148. 2017.

http://vecto.space

( 1)

8

( 2)