

Wordnet-based Evaluation of Large Distributional Models for Polish

anonymised ...

anonymised ...

Abstract

The paper presents construction of large scale test datasets for word embeddings on the basis of a large wordnet, apply them for evaluation of word embedding models and next to analyse and compare the usefulness of different word embeddings extracted from a very large corpus of Polish. In addition, several large word embeddings models built on the basis of a very large Polish corpus are presented.

1. Introduction

Distributional Semantics (DS) is focused on describing semantic associations between words on the basis of their distributional patterns in texts by applying statistical methods. DS methods are used to extract from corpora different kinds of the *Measures of Semantic Relatedness* (MSR). An MSR can cover the whole range of semantic relations, from topic or domain based till lexico-semantic relations. For many applications it is desirable to obtain an MSR which is close to a *Measure of Semantic Similarity* (MSS), i.e. a measure which assigns the highest values to words associated by linguistic lexico-semantic relations. Recently, word embeddings have become one of the best tools of DS. However, word embeddings, e.g. (Mikolov et al., 2013), are based on predicting word occurrence in a context (mostly a sequence) of other words. This aspect of co-occurrence prediction in a local context can influence an MSR built on the basis of word embeddings. An MSS can be an important source of knowledge supporting wordnet development, e.g. (Piasecki et al., 2009). However, the question is how to evaluate to which extent the given MSR resembles an proper MSS? Experiments with the participation of humans are laborious, costly and the datasets created as a result are of limited size. It is hard to construct an evaluation by application in a way revealing the properties of a potential MSS.

A large wordnet is built on knowledge originating from humans. It includes directly the knowledge about lexico-semantic relations and offers an opportunity to build large scale, realistic tests. Our goal is to construct large scale test datasets for word embeddings on the basis of a large wordnet, apply them for evaluation of word embedding models and next to analyse and compare the usefulness of different word embeddings extracted from a very large corpus of Polish. Finally, we want to publishing word embedding models of known properties built on the basis of a very large corpus of Polish.

2. Related Works

MSR evaluation methods can be roughly divided into intrinsic and extrinsic. The former are based on the direct evaluation of the MSR properties, e.g. by assessment by humans or comparison with a gold standard. The latter is based on applying an MSR as knowledge source in some NLP application.

Typical datasets used in the intrinsic evaluation are small, e.g. (Rubenstein and Goodenough, 1965), WS-353

(Finkelstein et al., 2002) and most of all 10 data sets discussed in (Faruqui and Dyer, 2014), where only two of them include ≈ 2000 and ≈ 3000 word pairs. They were used in many tests, in fact overused. Small sizes of these datasets make performing proper evaluation more difficult, e.g. because of the lack of the common partitioning into training, tuning and testing parts.

Very often datasets for MSR evaluation are collected during experiments based on testing human judgement in reaction to some prompting signal, which is close to reaction to a stimuli. For instance (Auguste et al., 2017) measured the correlation between the reaction times in the context of priming with ranking based on word embeddings. This is slightly different situation than analysis of lexical meanings during language utterance interpretation especially textual utterance. MSR is extracted from a text corpus, and it is more natural to evaluate it against language resources. Moreover, (Faruqui et al., 2016) noticed that the distinction between similarity and relatedness is not well defined and consistently expressed in most popular test datasets.

(Schnabel et al., 2015) evaluated systematically different DS models, but finally all tests were based on data collected in crowdsourcing experiments using Amazon Turk. (Jastrzebski et al., 2017) performed "evaluation focused on data efficiency" with respect to 4 categories, namely: "Similarity, Analogy, Sentence and Single word". In the case of similarity, which is most interesting for us, they used only well known data sets for English. For each type of dataset different combinations of preprocessing and classification algorithms were used. It is worth to notice, that the cost of preparing larger datasets for other language than English is quite substantial. This is one of the reasons that there are not many approaches for other languages, with notable exceptions e.g. (Hartmann et al., 2017) for Portuguese. In our case we want to explore construction of large datasets on the basis of an already existing wordnet. As we are interested in supporting wordnet development, so comparison with data collected in experiments with humans is not necessarily the best solution.

3. Wordnet-based Evaluation

In many approaches a wordnet was used to generate a wordnet-based measure of semantic similarity that was next used to assess the correlation between it and an MSR, e.g. (Lin, 1998). It was assumed that similarity rankings generated by the two measures should be similar. However, there are many wordnet-based similarity measures of

different properties and some of them depend on the additional knowledge like the information about the frequency of word senses. Thus, the result of the comparison can be different depending on the wordnet-based similarity measure applied and in all cases not straightforward in interpretation. We want to follow a different approach and to explore two methods that are free of these problems.

3.1. Synonymy tests

(Freitag et al., 2005) proposed wordnet-based synonymy test (WBST) in which for a *question word* x is automatically generated n -tuple:

$\mathbf{D} = \langle d_1, \dots, d_n \rangle$, such that one of the elements: d_i is a proper *answer*, i.e. it is synonymous with x and belongs to the same synset as x , and all other $d_j \neq d_i$ are *detractors*, i.e. false answers that are not synonymous with x . Elements of \mathbf{D} and the position of the answer are randomly selected. MSR is used to select for the problem word x the correct answer

In the case of some wordnets, e.g. plWordNet many synsets are singletons and include only one word. Thus they would be excluded from the test, and this can bias the evaluation result. To prevent this, in *Hypernymy-expanded WBST* (HWBST) answers for singleton synsets are selected from their hypernym synset, and in the same time these hypernyms are excluded from possible detractors. For a large wordnet, WBST and HWBST can include many thousands of question – answer pairs enabling very intensive testing of an MSR and partitioning the set in many different ways, e.g. test vs train, frequent vs infrequent or according to the domains of words.

Because detractors in WBST and HWBST are selected completely randomly, the majority of them comes from the parts of the wordnet that are very remote in relation to the question word. Thus these types of tests are relatively easy to be solved on the basis of an MSR. In order to make the test harder we need to select detractors in such a way that words from synsets semantically similar to the question words have a higher probability of being selected than words from synsets of small similarity. This version of the test is called *Extended WBST* (EWBST) (Piasecki et al., 2009). EWBST consists of pairs $\langle x_l, \mathbf{D}_l \rangle$, where x_l is a question word and $\mathbf{D}_l = \langle d_1, \dots, d_n \rangle$ is a sequence of possible answers such that of d_i is the correct answer, i.e. a synonym or hypernym of x_l , as in HWBST, while the rest of $d_j \in \mathbf{D}_l \wedge d_j \neq d_i$ are selected randomly from the whole wordnet but with the probability correlated to the *wordnet-based similarity measure* (WSM) between d_j and x_l . Any WSM can be used to generate EWBST, in the experiments presented in this work, we use a simple measure (1) proposed in (Agirre and Edmonds, 2006) based on the normalised length of a shortest path in the wordnet graph, that can be computed without knowing the frequency of senses:

$$WSM(w_1, w_2) = -\log \frac{\text{path}(w_1, w_2)}{2D} \quad (1)$$

In (1), w_1 and w_2 are words, $\text{path}(w_1, w_2)$ is the shortest path in the extended hypernymy graph between two synsets including, respectively: w_1 and w_2 , and

D is the maximal depth of the extended hypernymy graph. The graph was built from hypernymy relations and type/instance relations. In addition, as plWordNet hypernymy is not a single-rooted structure, we added to the graph several SUMO (Pease, 2011) concepts as top level nodes on the basis of mapping of plWordNet hypernymy root synsets onto SUMO concepts.

The idea of EWBST is to make detractors more similar to the correct answer and more difficult to be properly distinguished from the correct answer on the basis of MSR values.

3.2. Cut-off rendering tests

WBST-family tests illustrate the ability of an MSR to distinguish between words whose senses are located in different parts of the wordnet graph, while EWBST gives also insights into the sensitivity to small local differences. However, WBST-family tests concentrate on synonymy and hypernymy, as these two relations are mostly used in selection of the correct answers and detractors. However, from a good MSS we can also expect an ability to express other types of lexico-semantic relations. This can be measured with the help of a simple *Wordnet-based Cut-off Rendering* test (WBCR). In WBCR for each question word x a bag-of-words of words is generated in which they come from:

- the synset S_x of x
- and synsets S_i connected directly and also indirectly to S_x by wordnet relation.

S_x and S_i are indirectly connected if there is a path in the graph of wordnet relations such that it consists of a proper sequence of wordnet relations. Depending on the type of relations allowed for direct and indirect connections, as well as the assumed patterns for the paths and their maximal length, we can define different types of bags-of-words. Next, the evaluated MSR is used to reconstruct the extracted bag-of-words:

1. for the problem word x a ranking list of the words most related to x on the basis of the MSR values is generated; such a list will be called the *k-nearest neighbours* list (henceforth *k>NNL*) of x .
2. for the assumed k , the top k words from the list are collected as a reconstructed bag-of-words,
3. the reconstructed bag-of-words for x is compared with the wordnet-based bag-of-words, and precision, recall and F-measure are calculated.

This simple test is meaningful only for large, comprehensive wordnets or wordnets describing well some selected domains. However, WBCR has very simple interpretation and can be easily tuned to different subsets or domains of words and senses.

4. Experiments

During experiments, we built several word embeddings models from the largest corpus of Polish available. Next

we evaluated them in several tests based on plWordNet 3.1 (i.e. the most contemporary version) and compared with other word embedding models for Polish extracted from smaller corpora and published in the web.

4.1. Corpora and preprocessing

As a basis for the experiments we selected plWordNet 3.1 – a very large wordnet of Polish including $\approx 190,500$ different words, described by $\approx 282,500$ senses, more than 217,000 synsets and more than 750,000 relation links. plWordNet has been built by corpus-based wordnet developed method (Maziarz et al., 2013) and expresses very good coverage of words in large corpora (Maziarz et al., 2016).

For the extraction for word embeddings we used plWordNet Corpus 10.0 (plWNC) of Polish, which includes more than 4 billion words¹. It is also probably the largest corpus of Polish built in a controlled way and was used during the plWordNet development.

plWNC was used in the experiments in two versions of preprocessing:

plWNC-lem the corpus was first morphosyntactically tagged and lemmatised with the help of WCRFT2 tagger (Radziszewski, 2013; Radziszewski and Warzocha, 2014); strings: lemmas;grammatical class were in the input to *word2vec* (Mikolov et al., 2013).

plWNC-multi in the morpho-syntactically tagged plWNC Proper Names and multiword expressions described in plWordNet 3.1 were merged to single tokens.

plWNC-multi was prepared with the help of *Liner2* (Marcinićzuk et al., 2013) tool for recognition and classification of PNs. plWordNet 3.1 includes almost 60,000 Polish MWEs represented as lexical units and described by lexicalised morpho-syntactic constraints that allow for their efficient and accurate recognition in tagged texts (Kurc et al., 2012). We represent Proper Names (one and multiword, including many common words) and multiword expressions as single tokens in *plWNC-multi* in order to block the interpretation of their components as individual words. Components of PNs and MWEs can have very specific meanings (e.g. in non-compositional MWEs) that can influence the resulting word embeddings.

Corpora created from the Polish Wikipedia data alone (of $\approx 600M$ words) were used in two experiments reported in the literature. We evaluated these published word embedding models against our tests, too see Sec. 5.

¹It consists of IPI PAN Corpus (Przepiórkowski, 2004), the first annotated corpus of Polish, National Corpus of Polish (Przepiórkowski et al., 2012), Polish Wikipedia (from 2016), *Rzeczpospolita* Corpus (Weiss, 2008) – corpus of electronic editions of a Polish newspaper from the years 1993–2003, supplemented with text acquired from the Web – only text with small percentage of words unknown to a very comprehensive morphological analyser Morfeusz 2.0 (Woliński, 2014) were included; duplicates were automatically eliminated from the merged corpus.

4.2. Word embedding models tested

For the generation of word2vec models *Gensim* library was used (Rehurek and Sojka, 2010). On the basis of the set of 6 parameters, we selected during pre-experiments 9 different types of models to be evaluated experimentally, i.e. the following combinations:

1. vector size: 100, 300 and 1000,
2. algorithm type: *Skip-gram*, *CBOW ns* (with negative subsampling) and *CBOW hs* (with hierarchical softmax).

Thus, we tested: *Skip-gram* 100, *Skip-gram* 300, *Skip-gram* 1000, *CBOW ns* 100, *CBOW ns* 300, *CBOW ns* 1000, *CBOW hs* 100, *CBOW hs* 300 and *CBOW hs* 1000. In all models the minimal frequency of tokens (i.e. tagged lemmas and/or PN and MWE tokens) was set to ≥ 8 (`min_count=8`).

(Rogalski and Szczepaniak, 2016) first preprocessed a text corpus based on the Polish Wikipedia² by changing the text to lower case, numbers were divided into separate digits, and some non-text elements were deleted. Next two word embedding models were constructed: *CBOW* and *Skip-gram* models with negative sampling and the vector size: 300. The extracted models are publicly available in the internet³ and following the original names they will be called in the experiments, respectively: *pl-embeddings-cbow* and *pl-embeddings-skip*.

(Bojanowski et al., 2016) built *Skip-gram* models⁴ using *fastText* technique with the vector size 300 for many languages on the basis of Wikipedia data. For the extraction of the models a novel method in which “each word is represented as a bag of character n-grams”, cf (Bojanowski et al., 2016), was applied. It was designed for languages with richer inflection and was meant to better deal with a large number of word forms in such languages. Their model will be simply called *fastText.wiki.pl* in the experiments.

4.3. Tests

4.3.1. Wordnet-based Synonymy Tests

All three types of tests, namely: *WBST*, *HWBST* and *EWBST* were generated on the basis of noun part of plWordNet 3.1 in three versions corresponding to the minimal frequency of words in plWNC 10.0: 30, 200 and 1000, i.e. in a given tests all question, answer and detractor words had to express the predefined minimal frequency in the corpus. However, still the generated tests are very large e.g. *EWBST*(min. 1000) includes 19,996 question – answers pairs, *HWBST* (min. 30) includes 48,263 pairs, *WSBT*, smaller because singleton synsets are omitted, and *WBST*(min. 1000) includes 9,100 pairs – the smallest set.

²<https://pl.wikipedia.org>

³http://publications.it.p.lodz.pl/2016/word_embeddings/

⁴<https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md>

Vector size	Min freq.	Model	WBST	HWBST	EWBST
1000	1000	Skip-gram	92.43	89.00	63.97
		CBOw hs	91.54	89.34	63.21
		CBOw ns	91.68	89.31	62.99
	200	Skip-gram	92.52	89.80	62.51
		CBOw hs	92.71	90.11	60.94
		CBOw ns	92.58	90.11	60.97
	30	Skip-gram	90.43	88.84	58.92
		CBOw hs	92.56	90.05	57.35
		CBOw ns	92.51	90.07	57.30
300	1000	Skip-gram	90.81	88.24	62.50
		CBOw hs	90.32	88.12	61.00
		CBOw ns	90.70	88.49	62.13
	200	Skip-gram	91.81	89.36	61.24
		CBOw hs	91.46	89.29	59.45
		CBOw ns	91.11	89.50	60.76
	30	Skip-gram	90.99	89.43	58.25
		CBOw hs	91.36	89.41	55.97
		CBOw ns	91.35	89.79	57.50
100	1000	Skip-gram	88.84	86.01	59.42
		CBOw hs	87.71	86.14	58.26
		CBOw ns	88.14	86.71	59.34
	200	Skip-gram	89.78	87.53	58.52
		CBOw hs	88.97	87.33	56.75
		CBOw ns	89.05	87.57	58.12
	30	Skip-gram	89.79	88.21	55.99
		CBOw hs	89.44	87.62	53.52
		CBOw ns	89.63	88.13	55.27
	1000	pl-embeddings-cbow	71.63	69.36	43.71
		pl-embeddings-skip	76.30	74.54	47.16
		fastText.wiki.pl	80.01	78.17	52.42
	200	pl-embeddings-cbow	71.79	69.46	42.31
		pl-embeddings-skip	76.89	74.65	45.53
		fastText.wiki.pl	80.11	79.16	51.40
	30	pl-embeddings-cbow	71.49	70.35	41.85
		pl-embeddings-skip	77.41	75.69	45.28
		fastText.wiki.pl	81.44	80.27	51.39

Table 1: WBST-like tests generated from noun in plWordNet 3.1 and applied to word embedding models extracted from *plWNC-multi*.

4.3.2. Wordnet-based Cut-off Rendering tests

As in the case of the WBST-like tests, the cut-off tests were generated on the basis of nouns in plWordNet 3.1 and in three main versions with respect to the minimal frequency of nouns in plWNC 10.0: 30, 200 and 1000 (the numbers of bag of words are smaller than the number of pairs in WBST-like tests but similarly large).

The wordnet context of a problem word x , which was represented as a bag of words was defined in three different ways:

Cnt – all words linked to x by direct relation links, i.e. from synsets linked directly to the synset of x and also by direct lexical relations to one of the x senses; it also includes synonyms of x .

CntH – **Cnt** expanded with all indirect hyponyms and hypernyms of x up to the hypernymy and hyponymy paths of the maximal length 3.

CntHC – **CntH** expanded with all $k = m + n$ cousins of x with $k = 3$, i.e. words from synsets accessible from the synsets of x by hyper/hyponymy paths of up to m hypernymy and n hyponymy links.

Thus, **Cnt** measures the ability of an MSR to find words in very close relations (e.g. as a potential tool supporting build description of x senses), **CntH** illustrates the use of the MSR as a tool supporting construction of hyper/hyponymy structures and **CntHC**. All cut-off tests were applied to the k -best neighbours lists with $k \in \{10, 100\}$ generated for nouns from plWordNet.

5. Results

The results of the tests in Tab. 1 illustrate well the differences in the difficulty of the tests: WBST is the simplest one, EWBST the hardest. The difference between EWBST and the two other tests is striking in all experiments. The difficulty of EWBST can be tuned by changing

Cut-off Precision							
k NN		10			100		
Model	Min. freq.	Cnt (%)	CntH (%)	CntHC (%)	Cnt (%)	CntH (%)	CntHC (%)
CBOW hs	1000	10.08	11.38	27.38	2.47	3.24	11.67
CBOW ns	1000	10.10	11.29	26.21	2.48	3.17	10.95
Skip-Gram	1000	9.09	10.01	22.02	2.05	2.48	7.72
CBOW hs	200	8.46	9.48	24.43	2.00	2.58	10.33
CBOW ns	200	8.20	9.08	23.10	1.96	2.46	9.64
Skip-Gram	200	7.74	8.46	19.94	1.69	2.02	7.06
CBOW hs	30	6.74	7.51	20.63	1.59	2.03	8.74
CBOW ns	30	6.56	7.23	19.89	1.55	1.93	8.32
Skip-Gram	30	6.04	6.58	16.43	1.32	1.57	5.94

Cut-off Recall							
k NN		10			100		
		Cnt (%)	CntH (%)	CntHC (%)	Cnt (%)	CntH (%)	CntHC (%)
CBOW hs	1000	8.59	4.84	2.98	17.52	10.42	7.45
CBOW ns	1000	8.32	4.57	2.76	17.26	9.90	7.05
Skip-Gram	1000	7.65	4.20	2.46	15.04	8.46	5.61
CBOW hs	200	8.69	5.12	3.47	17.28	10.61	8.19
CBOW ns	200	7.96	4.56	3.06	16.36	9.66	7.60
Skip-Gram	200	7.88	4.55	2.95	15.04	8.78	6.40
CBOW hs	30	7.71	4.62	3.27	15.66	9.73	7.86
CBOW ns	30	7.18	4.21	2.99	14.81	8.89	7.38
Skip-Gram	30	6.74	3.92	2.64	13.07	7.72	5.94

F measure							
k NN		10			100		
		Cnt (%)	CntH (%)	CntHC (%)	Cnt (%)	CntH (%)	CntHC (%)
CBOW hs	1000	9.28	6.79	5.38	4.33	4.94	9.10
CBOW ns	1000	9.12	6.51	5.00	4.33	4.80	8.58
Skip-Gram	1000	8.31	5.92	4.42	3.61	3.83	6.50
CBOW hs	200	8.57	6.65	6.07	3.59	4.15	9.13
CBOW ns	200	8.08	6.07	5.41	3.50	3.92	8.50
Skip-Gram	200	7.81	5.92	5.13	3.04	3.29	6.71
CBOW hs	30	7.20	5.72	5.64	2.88	3.36	8.28
CBOW ns	30	6.86	5.32	5.20	2.81	3.17	7.82
Skip-Gram	30	6.37	4.92	4.55	2.40	2.60	5.94

Table 2: Wordnet-based Cut-off Rendering tests generated from nouns in plWordNet 3.0 and applied for word embedding models extracted from *plWNC-multi*, where kNN is the length of the k nearest neighbour lists.

the wordnet-base similarity measure it is based on and the dependency between the similarity measure and the distribution of the probability of detractor word selection.

Skip-gram model is better than CBOW in most of the cases and in the other cases the differences are small. Also among the models from the literature, models based on Skip-gram scheme, including *fastText.wiki.pl* (which is a Skip-gram model too) express higher results. This is especially visible in the case of the mode difficult EWBST test. The *wiki.pl* was superior among the models built only on the data from Wikipedia, i.e. several times smaller than *plWNC 10.0*. However, all models built on much smaller corpus produced much worse results. We tested also models based on *plWNC-lem* version of the large corpus (to be presented in the full paper) and all models were slightly but significantly worse in the WBST-family of tests.

Contrary to the synonymy tests, in the case of WBCR evaluations of the models generated from *plWNC-multi*

presented in Tab. 2, we can notice that CBOW models are superior in all cases in comparison to Skip-gram models. It means that Skip-gram models are better in describing differences between word meanings, while CBOW enable broader exploration of potential lexico-semantic relations. However, relatively good precision mean that instances of lexico-semantic relations receive higher values. Definitely the results of the test are negatively biased by lacking relation instances in plWordNet. This kind of tests and evaluations can be used as diagnostic tool to spot the subdomains in a wordnet that are potentially not well enough described by relation links.

We can observe also the application hierarchical softmax consistently produces better results in all frequency ranges. However, hierarchical softmax should result in better estimation of the representation.

We evaluated also word embedding models extracted from *plWNC-lem*, i.e. a version without folding PNs and

Cut-off Precision							
k NN		10			100		
Model	Min. freq.	Cnt (%)	CntH (%)	CntHC (%)	Cnt (%)	CntH (%)	CntHC (%)
pl-embeddings-cbow	1000	4.97	6.10	20.98	1.37	1.94	10.51
pl-embeddings-skip	1000	3.72	3.87	7.41	1.15	1.33	5.49
fastText.wiki.pl	1000	4.03	4.24	7.24	1.52	1.78	6.04
pl-embeddings-cbow.	200	3.92	4.81	17.77	1.08	1.52	8.78
pl-embeddings-skip	200	3.42	4.05	13.58	1.02	1.34	6.65
fastText.wiki.pl	200	3.90	4.07	7.33	1.31	1.51	5.76
pl-embeddings-cbow	30	3.28	4.03	15.56	0.90	1.27	7.67
pl-embeddings-skip	30	2.99	3.55.	12.56.	0.88	1.16	6.14
fastText.wiki.pl	30	3.72	3.87.	7.41.	1.15	1.33	5.49

Cut-off Recall							
k NN		10			100		
		Cnt (%)	CntH (%)	CntHC (%)	Cnt (%)	CntH (%)	CntHC (%)
pl-embeddings-cbow	1000	2.42.	1.59.	0.89.	5.93.	4.22.	2.71
pl-embeddings-skip	1000	3.14	1.97	1.07	7.67	4.86	3.02
wiki.pl	1000	2.43.	1.48.	0.72.	7.56.	4.64.	2.44
pl-embeddings-cbow	200	2.19	1.46	0.88	5.37.	3.90	2.66
pl-embeddings-skip	200	2.09	1.37	0.79	5.42.	3.82	2.45
wiki.pl.	200	2.90	1.80	0.94	7.78	4.86	2.83
pl-embeddings-cbow	30	1.99	1.34	0.83	4.93	3.60	2.53
pl-embeddings-skip	30	2.04	1.37	0.82	5.27	3.78	2.50
fastText.wiki.pl	30	3.14	1.97	1.07	7.67	4.86	3.02

F measure							
k NN		10			100		
		Cnt (%)	CntH (%)	CntHC (%)	Cnt (%)	CntH (%)	CntHC (%)
pl-embeddings-cbow	1000	3.26	2.53	1.71	2.23	2.66	4.31
pl-embeddings-skip	1000	3.40	2.61	1.88	2.00	2.08	3.89
fastText.wiki.pl	1000	3.03.	2.19	1.30	2.53	2.57	3.48
pl-embeddings-cbow	200	2.81	2.24	1.68	1.80	2.18	4.08
pl-embeddings-skip	200	2.59	2.05	1.50	1.71	1.98	3.58
fastText.wiki.pl	200	3.32	2.50	1.67	2.24	2.31	3.79
pl-embeddings-cbow	30	2.48.	2.01	1.58	1.52	1.87	3.81
pl-embeddings-skip	30	2.43	1.97	1.54	1.50	1.77	3.56
fastText.wiki.pl	30	3.40	2.61	1.88	2.00	2.08	3.89

Table 3: Wordnet-based Cut-off Rendering tests generated from nouns in plWordNet 3.0 and applied for word embedding models published in the Web extracted from the Polish Wikipedia, where kNN is the length of the k nearest neighbour lists.

MWEs into single tokens. WBCR tests showed higher precision by about 2%, but significantly lower recall and F-measure (the results will be presented in the full paper). It can be probably caused by the access to additional co-occurrences that are hidden in *plWNC-multi* in the PN and MWE tokens. However, models based on *plWNC-multi* offer a unique opportunity of obtaining good distributional description of PNs and MWEs.

Quite surprisingly, we can observe in Tab. 3 that models built on a smaller corpus of Wikipedia behave in a different way in WBCR tests than those constructed on a very large corpus. In Tab. 3Skip-gram models express higher recall, fastTest Skip-gram with subword representation have much higher recall for words with lower frequency. However, all results obtained on the Polish Wikipedia are worse than those in Tab. 2 generated from a very large corpus (including also the Wikipedia data). It means that for WBCR tests covering a spectrum of relations, larger data used re-

sult in the improvement of the model.

In the full paper we will also present evaluation of models for other Parts of Speech.

6. Conclusions

We showed that a large comprehensive wordnet can be successfully used as a basis for two different types of MSR evaluation methods, namely the family of Wordnet-based Synonymy Tests and Wordnet-based Cut-off Rendering tests. In both types of tests very large datasets can be generated allowing for very intensive testing and high statistical significance of the test results. The datasets are enough large to conveniently partitioned according to the frequency criteria of semantic criteria. In fact the datasets and tests are based on human decisions expressed in the wordnet structure. Both tests describe the ability of an MSR to be used as a basis for developing a lexico-semantic language resource.

WBST-family tests focus on the ability of an MSR to distinguish between different lexical meanings, while WBCR is sensitive more to representation of different types of wordnet relations by an MSR. As a result both types of tests are quite complementary. Moreover, by changing the similarity and context definitions in EWBST we can obtain tests of different difficulty.

In the further work, we develop a wordnet-based test that has properties of contextual tests, e.g. which is similar to Stanford contextual word similarity dataset (SCWS) (Huang et al., 2012).

We will also expand the presented evaluation to the dataset covering all four PoS, namely nouns, adjectives, verbs and adverbs.

The constructed word embedding models and evaluation datasets have been published on open licences under the link: *anonymised*

7. References

- Agirre, Eneko and Philip Edmonds (eds.), 2006. *Word Sense Disambiguation: Algorithms and Applications*. Springer.
- Auguste, Jeremy, Arnaud Rey, and Benoit Favre, 2017. Evaluation of word embeddings against cognitive processes: primed reaction times in lexical decision and naming tasks. In *Proceedings of the 2nd Workshop on Evaluating Vector-Space Representations for NLP*. Copenhagen, Denmark: ACL.
- Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov, 2016. Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- Faruqui, Manaal and Chris Dyer, 2014. Community evaluation and exchange of word vectors at `wordvectors.org`. In *Proc. of ACL: System Demo*.
- Faruqui, Manaal, Yulia Tsvetkov, Pushpendre Rastogi, and Chris Dyer, 2016. Problems with evaluation of word embeddings using word similarity tasks. In *Proceedings of the 1st Workshop on Evaluating Vector Space Representations for NLP*. Berlin, Germany: ACL.
- Finkelstein, Lev, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppín., 2002. Placing search in context: The concept revisited. *ACM Transactions on Information Systems*, 20(1).
- Freitag, Dayne, Matthias Blume, John Byrnes, Edmond Chow, Sadik Kapadia, Richard Rohwer, and Zhiqiang Wang, 2005. New experiments in distributional representations of synonymy. In *Proceedings of the Ninth Conference on Computational Natural Language Learning (CoNLL-2005)*. Ann Arbor, Michigan: Association for Computational Linguistics.
- Hartmann, Nathan, Erick R. Fonseca, Christopher Shulby, Marcos Vinícius Treviso, Jessica Rodrigues, and Sandra M. Aluísio, 2017. Portuguese word embeddings: Evaluating on word analogies and natural language tasks. *CoRR*, abs/1708.06025.
- Huang, Eric H, Richard Socher, Christopher D Manning, and Andrew Y Ng., 2012. Improving word representations via global context and multiple word prototypes. In *Proceedings of ACL*.
- Jastrzebski, Stanislaw, Damian Lesniak, and Wojciech Marian Czarnecki, 2017. How to evaluate word embeddings? on importance of data efficiency and simple supervised tasks. *CoRR*, abs/1702.02170.
- Kurc, Roman, Maciej Piasecki, and Bartosz Broda, 2012. Constraint based description of polish multiword expressions. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*. Istanbul, Turkey: European Language Resources Association (ELRA). **10 pkt Web of Science ACL Anthology**.
- Lin, Dekang, 1998. Automatic retrieval and clustering of similar words. In *International Conference On Computational Linguistics (COLING'98). Proceedings of the 17th International Conference on Computational Linguistics*, volume 2. ACL.
- Marcinićzuk, Michał, Jan Kocoń, and Maciej Janicki, 2013. Liner2 – a customizable framework for proper names recognition for Polish. In Robert Bembenik, Lukasz Skonieczny, Henryk Rybinski, Marzena Kryszkiewicz, and Marek Niezgodka (eds.), *Intelligent Tools for Building a Scientific Information Platform*. pages 231–253.
- Maziarz, Marek, Maciej Piasecki, Ewa Rudnicka, and Stan Szpakowicz, 2013. Beyond the transfer-and-merge wordnet construction: plWordNet and a comparison with WordNet. In Ruslan Mitkov, Galia Angelova, and Kalina Boncheva (eds.), *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*. Hissar, Bulgaria: INCOMA Ltd. Shoumen, BULGARIA. **ACL Anthology**.
- Maziarz, Marek, Maciej Piasecki, Ewa Rudnicka, Stan Szpakowicz, and Paweł Kędzia, 2016. plwordnet 3.0 – a comprehensive lexical-semantic resource. In N. Calzolari, Y. Matsumoto, and R. Prasad (eds.), *Proc. of COLING 2016, 26th Inter. Conf. on Computational Linguistics*. ACL.
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean, 2013. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781.
- Pease, Adam, 2011. *Ontology - A Practical Guide*. Articulate Software Press.
- Piasecki, Maciej, Stanisław Szpakowicz, and Bartosz Broda, 2009. *A Wordnet from the Ground Up*. Wrocław: Oficyna Wydawnicza Politechniki Wrocławskiej. **Monografia, 25 pkt**.
- Przepiórkowski, Adam, 2004. *The IPI PAN Corpus: Preliminary version*. Institute of Computer Science, Polish Academy of Sciences.
- Przepiórkowski, Adam, Mirosław Bańko, Rafał L. Górski, and Barbara Lewandowska-Tomaszczyk (eds.), 2012. *Narodowy Korpus Języka Polskiego [Eng.: National Corpus of Polish]*. Wydawnictwo Naukowe PWN.
- Radziszewski, Adam, 2013. A tiered CRF tagger for Polish. In R. Bembenik, Ł. Skonieczny, H. Rybiński, M. Kryszkiewicz, and M. Niezgodka (eds.), *Intelligent Tools for Building a Scientific Information Platform*:

- Advanced Architectures and Solutions*. Springer Verlag.
- Radziszewski, Adam and Radosław Warzocha, 2014. WCRFT2. CLARIN-PL digital repository, <http://hdl.handle.net/11321/36>.
- Řehůřek, Radim and Petr Sojka, 2010. Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. Valletta, Malta: ELRA. <http://is.muni.cz/publication/884893/en>.
- Rogalski, Marek and Piotr S. Szczepaniak, 2016. Word embeddings for the polish language. In L. Rutkowski, M. Korytkowski, R. Scherer, R. Tadeusiewicz, L.A. Zadeh, and J.M. Zurada (eds.), *15th International Conference, ICAISC 2016, Zakopane, Poland, June 12-16, 2016, Proceedings*, volume 9692 of *LNAI, Artificial Intelligence and Soft Computing*. Springer.
- Rubenstein, Herbert and John B. Goodenough, 1965. Contextual correlates of synonymy. *Communications of ACM*, 8(10).
- Schnabel, Tobias, Igor Labutov, David Mimno, and Thorsten Joachims, 2015. Evaluation methods for unsupervised word embeddings. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Lisbon, Portugal: ACL.
- Weiss, Dawid, 2008. Korpus Rzeczpospolitej [Corpus of text from the online edition of “Rzeczpospolita”]. <http://www.cs.put.poznan.pl/dweiss/rzeczpospolita>.
- Woliński, Marcin, 2014. Morfeusz reloaded. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Ninth International Conference on Language Resources and Evaluation, LREC 2014*. Reykjavík, Iceland: ELRA.