

# Wordnet-based Evaluation of Large Distributional Models for Polish

Maciej Piasecki, Gabriela Czachor, Arkadiusz Janz, Dominik Kaszewski, Paweł Kędzia

G4.19 Research Group, Department of Computational Intelligence  
Wrocław University of Science and Technology, Wrocław, Poland

& CLARIN-PL [clarin-pl.eu](http://clarin-pl.eu)



Politechnika  
Wroclawska



# Agenda

- Wordnet-based tests for Distributional Semantics
- Synonymy tests
- Cut-off rendering tests
- Experiments
  - Corpora and preprocessing
  - Word embedding models tested
  - Tests based on *pWordNet*
  - Analogy tests (not wordnet-based)
- Results
- Conclusions and suggestions



# Wordnet-based tests for Distributional Semantics

- Background
  - Distributional Semantics (DS) is focused on describing semantic associations between words on the basis of their distributional patterns in corpora
  - Corpora → *Measures of Semantic Relatedness*  $\rightsquigarrow$ ? *Measure of Semantic Similarity*
  - *Word embeddings* are based on predicting a word occurrence in a context (mostly a sequence) of other words
  - A large wordnet is built on knowledge originating from humans
- Goals
  - to construct large scale test datasets for word embeddings on the basis of a large wordnet
  - to evaluate and compare different word embeddings extracted from a very large corpus of Polish
  - to publish: tests and word embeddings



# Wordnet-based DS Evaluation

- Wordnet-based Similarity Measures → Correlation of similarity rankings
  - but this comparison depends on a particular wordnet-based similarity measure applied
- *Wordnet-based Synonymy Test (WBST)*
- *Wordnet-based Cut-off Rendering Test (WBCR)*



# Wordnet-based Synonymy Tests

- Proposed by (Freitag et al., 2005) following TOEFL synonymy tests:
  - for a *question word*  $x$
  - an  $n$ -tuple is automatically generated:  
 $\mathbf{D} = \langle d_1, \dots, d_n \rangle$ ,  
such that one the elements:
    - $d_i$  is the correct *answer* – synonymous with  $x$
    - all other  $d_j \neq d_i$  are *detractors*, i.e. false answers, not synonymous with  $x$
  - Elements of  $\mathbf{D}$  and the position of the correct answer are randomly selected
- *Pros*: very large tests can be generated enabling very intensive testing
- *and cons*: numerous singleton synsets, too easy detractors



# Wordnet-based Synonymy Tests

## Hypernymy-expanded WBST (HWBST)

- Answers for singleton synsets are selected from their hypernym synsets
- Such hypernyms are excluded from possible detractors
- Examples of QA tuples:
  - *majątek* 'property':  
{*okręt* 'ship', *uszanowanie* 'respect', ***mienie*** 'property, assets', *żywot* '≈life'}
  - *student* 'student':  
{*momencik* '≈an indefinitely short time', *łysina* 'bald spot', *skażenie* 'contamination', ***żak*** '≈student'}



# Wordnet-based Synonymy Tests

## Extended WBST (EWBST) (1)

- Idea: higher probability for the selection of detractors from synsets semantically similar to the question words
- EWBST consists of pairs:  $\langle x_I, \mathbf{D}_I \rangle$ , where
  - $x_I$  is a question word,
  - $\mathbf{D}_I = \langle d_1, \dots, d_n \rangle$  such that
    - $d_i$  is the correct answer, i.e. a synonym or hypernym of  $x_I$ , as in HWBST,
    - $d_j \in \mathbf{D}_I \wedge d_j \neq d_i$  are selected randomly from the whole wordnet but with the probability correlated to the *wordnet-based similarity measure*  $WSM(d_j, x_I)$ .

# Wordnet-based Synonymy Tests

## Extended WBST (EWBST) (2)

- WSM based on the normalised length of a shortest path in the wordnet graph (Agirre and Edmonds, 2006)

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

- $w_1$  and  $w_2$  are lemmas,
- $path(w_1, w_2)$  is the shortest path in the extended hypernymy graph between two synsets including  $w_1$  and  $w_2$ ,
- $D_m$  is the maximum depth of the extended hypernymy graph.
- Modified

$$WSM_a(w_1, w_2) = \max\left(-\log \frac{path(w_1, w_2)}{2D_a}, 0\right) \quad (2)$$

- $D_a$  is an average depth – promotes closer synsets





# Wordnet-based Synonymy Tests

## Extended WBST – Examples of QA tuples

- *majątek* 'property':  
⟨**mienie** 'property, assets', *banknot* 'banknote', *bon* 'voucher', *wyrównanie* 'compensation'⟩
- *student* 'student':  
⟨*aspirant* '≈candidate', *licencjat* 'bachelor's degree', **żak** '≈student', *lektor* 'lector'⟩



# Wordnet-based Cut-off Rendering Test

- Idea: to expand tests on other relations than synonymy and hyper/hyponymy
- 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 selected wordnet relations.
- Different *path definitions* can be used
- Evaluated MSR is used to reconstruct the extracted bag-of-words
  - 1 for a word  $x$  the *k-nearest neighbours* list  $k\text{-NNL}(x)$  of the words most related to  $x$  according to MSR
  - 2 for the assumed  $k$ , the top  $k$  words from the list are collected as a reconstructed bag-of-words,
  - 3 and compared with the wordnet-based bag-of-words



# Experiments

## Corpora and preprocessing

- Wordnet
  - p1WordNet 3.1 – a very large wordnet of Polish
  - 190,853 lemmas, 284,925 lexical units, 219,380 synsets and  $\approx 650,000$  relations
  - expresses very good coverage of words in large corpora
- p1WordNet Corpus 10.0 (p1WNC) of Polish
  - more than 4 billion words: several corpora supplemented with text acquired from the Web, only text in Polish, automated elimination of duplicates
- Corpora created from the Polish Wikipedia (of  $\approx 600M$  words)

**p1WNC-lem** morphosyntactically tagged, strings:  
“lemma:grammatical class” were in the input to  
*word2vec* (Mikolov et al., 2013)

**p1WNC-multi** Proper Names and multiword expressions (60k) from  
p1WordNet 3.1 merged as single tokens



# Experiments

## Word embedding models tested (1)

- Models generated by *word2vec* – *Gensim* library implementation
  - ① vector size: 100, 300 and 1000,
  - ② algorithm type: *Skip-gram*, *CBOW ns* (negative subsampling) and *CBOW hs* (with hierarchical softmax).
  - ③ tested models:  
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
  - ④ minimal frequency of tokens:  $\geq 8$  (`min_count=8`)
  - ⑤ freely available:  
<https://clarin-pl.eu/dspace/handle/11321/442>



# Experiments

## Word embedding models tested (2)

- *word2vec* models from literature
  - (Rogalski & Szczepaniak, 2016) built on Wikipedia
    - CBOW and Skip-gram models with negative sampling and the vector size: 300
    - text to lower case, numbers were divided into separate digits, and some non-text elements were deleted
  - (Mykowiecka et al., 2017) ) on National Corpus of Polish
    - '*ncp-lemmas*' or '*ncp-forms*' – full data
    - "restricted data sets": only nouns, adjectives, adverbs, verb forms, and abbreviations
    - Skip-gram and CBOW architectures and the vector size of 100 and 300
- *fastText* models (words represented as n-grams)
  - (Bojanowski et al., 2016) Skip-gram models, vector size the vector size 300 for many languages on the basis of Wikipedia
  - *fastText.plWNC*: Skip-gram 300 models with min. word frequencies of 5, 20 and 50 built on the plWNC 10.0 Corpus



# Experiments

## Tests based on *pWordNet* (1)

- Wordnet-based Synonymy Tests
  - WBST, HWBST and EWBST
  - three versions corresponding to the minimal frequency of words in pWNC 10.0: 30, 200 and 1000
  - e.g.
    - EWBST(min. 1000) includes 19,996 question – answers pairs,
    - HWBST (min. 30) includes 48,263 pairs,
    - and WBST(min. 1000) includes 9,100 pairs (singleton synsets omitted)
  - freely available:  
<https://clarin-pl.eu/dspace/handle/11321/446>



# Experiments

## Tests based on *pWordNet* (2)

- Wordnet-based Cut-off Rendering Tests
  - three versions for to the minimal frequency of words in pWNC 10.0: 30, 200 and 1000
  - smaller numbers of bag of words, but still large data sets
  - types of paths for indirect links to the problem lemma  $x$ 
    - Cnt** – only direct relation links (synset or lexical), including synonymy
    - CntH** – **Cnt** expanded with all indirect hyponyms and hypernyms of  $x$  up to the path length 3.
    - CntHC** – **CntH** expanded with all  $k = m + n$  *cousins* of  $x$  with  $k = 3$
  - freely available:  
<https://clarin-pl.eu/dspace/handle/11321/446>



# Experiments

## Analogy tests

- Most popular technique of the evaluation of word embeddings
  - testing MSR ability of reflecting word analogies
  - analogy consists of 2 pairs of words in a similar relation
  - MSR is used to find the best fitting lemma  $d$  in  $(\vec{b} + \vec{c}) - \vec{a} = \vec{d}$
- Limitations
  - small size of a dataset – typically 200-300
  - potential polysemy of lemmas in pairs
- dataset of (Mykowiecka et al., 2017): ~200? analogy pairs from (?????) manually translated to Polish (but out of context)



# Results – selected

## Wordnet-based Synonymy Tests

Vector size	Min freq.	Model	WBST	HWBST	EWBST
1000	1000	<i>w2w-plWNC-multi-skipg-ns</i>	<b>92.43</b>	89.00	<b>63.97</b>
		<i>w2w-plWNC-multi-cbow-hs</i>	91.54	<b>89.34</b>	63.21
		<i>w2w-plWNC-multi-cbow-ns</i>	91.68	89.31	62.99
	200	<i>w2w-plWNC-multi-skipg-ns</i>	92.52	89.80	<b>62.51</b>
		<i>w2w-plWNC-multi-cbow-hs</i>	<b>92.71</b>	90.11	60.94
		<i>w2w-plWNC-multi-cbow-ns</i>	92.58	<b>90.11</b>	60.97
	30	<i>w2w-plWNC-multi-skipg-ns</i>	90.43	88.84	<b>58.92</b>
		<i>w2w-plWNC-multi-cbow-hs</i>	<b>92.56</b>	90.05	57.35
		<i>w2w-plWNC-multi-cbow-ns</i>	92.51	<b>90.07</b>	57.30
1000	1000	pl-embeddings-cbow	71.63	69.36	43.71
		pl-embeddings-skip	76.30	74.54	47.16
		fastText.wiki.pl	<b>80.01</b>	<b>78.17</b>	<b>52.42</b>
	200	pl-embeddings-cbow	71.79	69.46	42.31
		pl-embeddings-skip	76.89	74.65	45.53
		fastText.wiki.pl	<b>80.11</b>	<b>79.16</b>	<b>51.40</b>
	30	pl-embeddings-cbow	71.49	70.35	41.85
		pl-embeddings-skip	77.41	75.69	45.28
		fastText.wiki.pl	<b>81.44</b>	<b>80.27</b>	<b>51.39</b>



# Results – selected

## Wordnet-based Cut-off Rendering Test

### Cut-off Precision

Model	k NN	10			100		
	Min. f.	Cnt	CntH	CntHC	Cnt	CntH	CntHC
w2w- <i>pIWNC-multi-cbow-hs</i>	1000	13.42	15.12	<b>35.67</b>	<b>3.31</b>	<b>4.29</b>	<b>17.04</b>
w2w- <i>pIWNC-multi-cbow-ns</i>	1000	<b>13.62</b>	<b>15.16</b>	34.25	3.30	4.22	15.96
w2w- <i>pIWNC-multi-skipg</i>	1000	12.35	13.47	28.07	2.66	3.18	10.12
ft- <i>pIWNC-multi-skipg</i>	1000	8.74	9.24	15.72	2.59	3.00	8.14
w2w- <i>pIWNC-lem-cbow-hs</i>	1000	12.86	14.26	33.38	3.11	3.93	15.75
w2w- <i>pIWNC-lem-cbow-ns</i>	1000	9.65	10.58	25.40	2.17	2.60	9.71
w2w- <i>pIWNC-lem-skipg</i>	1000	11.61	12.61	27.15	2.47	2.92	9.82
ft- <i>pIWNC-lem-skipg</i>	1000	7.39	7.72	13.31	2.25	2.54	7.25



# Results – selected

## Wordnet-based Cut-off Rendering Test

### Cut-off Recall

		10			100		
k NN		Cnt	CntH	CntHC	Cnt	CntH	CntHC
w2w- <i>plWNC-multi-cbow-hs</i>	1000	<b>10.33</b>	<b>7.10</b>	<b>3.42</b>	<b>20.83</b>	<b>15.69</b>	<b>8.61</b>
w2w- <i>plWNC-multi-cbow-ns</i>	1000	10.09	6.84	3.24	20.27	14.84	8.16
w2w- <i>plWNC-multi-skipg</i>	1000	9.24	6.26	2.91	17.22	12.20	6.26
ft- <i>plWNC-multi-skipg</i>	1000	7.33	4.87	2.18	17.54	12.22	5.80
w2w- <i>plWNC-lem-cbow-hs</i>	1000	8.74	6.05	2.85	17.67	13.03	7.03
w2w- <i>plWNC-lem-cbow-ns</i>	1000	6.71	4.61	2.18	13.20	9.46	4.99
w2w- <i>plWNC-lem-skipg</i>	1000	8.19	5.64	2.60	15.12	10.82	5.41
ft- <i>plWNC-lem-skipg</i>	1000	5.92	4.04	1.82	14.88	10.40	4.85



# Results – selected

## Wordnet-based Cut-off Rendering Test

F measure

		F measure					
		10			100		
	k NN	Cnt	CntH	CntHC	Cnt	CntH	CntHC
<i>w2w-pIWNC-multi-cbow-hs</i>	1000	<b>11.67</b>	<b>9.66</b>	<b>6.23</b>	<b>5.72</b>	<b>6.74</b>	<b>11.44</b>
<i>w2w-pIWNC-multi-cbow-ns</i>	1000	11.59	9.42	5.92	5.68	6.57	10.80
<i>w2w-pIWNC-multi-skipg</i>	1000	10.57	8.55	5.27	4.61	5.05	7.73
<i>ft-pIWNC-multi-skipg</i>	1000	7.97	6.38	3.83	4.51	4.82	6.77
<i>w2w-pIWNC-lem-cbow-hs</i>	1000	10.41	8.49	5.25	5.29	6.04	9.72
<i>w2w-pIWNC-lem-cbow-ns</i>	1000	7.91	6.42	4.02	3.73	4.08	6.59
<i>w2w-pIWNC-lem-skipg</i>	1000	9.60	7.80	4.75	4.24	4.60	6.98
<i>ft-pIWNC-lem-skipg</i>	1000	6.57	5.30	3.20	3.90	4.09	5.81



# Results – selected

## Wordnet-based Cut-off Rendering Test

F measure

		F measure					
k NN		10			100		
		Cnt	CntH	CntHC	Cnt	CntH	CntHC
<i>ncp-forms-rest-cbow-ns</i>	1000	6.13	5.13	3.26	3.30	3.95	6.88
<i>ncp-lemmas-all-cbow-hs</i>	1000	9.62	<b>8.08</b>	<b>5.07</b>	<b>4.91</b>	<b>5.78</b>	<b>9.46</b>
<i>ncp-lemmas-all-cbow-ns</i>	1000	<b>9.72</b>	7.91	4.80	5.01	5.74	9.04
<i>ncp-lemmas-all-skipg-hs</i>	1000	8.64	7.08	4.36	4.30	4.88	7.84
<i>ncp-lemmas-all-skipg-ns</i>	1000	8.42	6.88	4.18	3.75	4.07	6.20
<i>ncp-lemmas-rest-cbow-hs</i>	1000	9.91	8.27	5.05	4.99	5.82	9.20
<i>ncp-lemmas-rest-cbow-ns</i>	1000	10.03	8.13	4.80	5.06	5.75	8.77
<i>ncp-forms-rest-cbow-ns</i>	200	5.29	4.51	3.35	2.62	3.14	6.67
<i>ncp-lemmas-all-cbow-hs</i>	200	<b>8.71</b>	<b>7.52</b>	<b>5.54</b>	4.00	<b>4.75</b>	<b>9.44</b>
<i>ncp-lemmas-all-cbow-ns</i>	200	8.55	7.13	5.12	<b>4.01</b>	4.61	8.98
<i>ncp-lemmas-all-skipg-hs</i>	200	8.01	6.79	4.96	3.58	4.11	8.13
<i>ncp-lemmas-all-skipg-ns</i>	200	7.30	6.11	4.29	3.00	3.28	6.00
<i>ncp-lemmas-rest-cbow-hs</i>	200	9.02	7.75	5.56	4.07	4.80	9.20
<i>ncp-lemmas-rest-cbow-ns</i>	200	8.87	7.38	5.17	4.07	4.65	8.79



# Results – selected

## Wordnet-based Cut-off Rendering Test

F measure

		F measure					
		10			100		
	k NN	Cnt	CntH	CntHC	Cnt	CntH	CntHC
pl-embeddings-cbow	1000	<b>3.79</b>	<b>3.31</b>	<b>2.15</b>	2.32	<b>2.94</b>	<b>5.08</b>
pl-embeddings-skipg	1000	3.35	2.82	1.80	2.15	2.56	4.20
fastText.wiki.pl	1000	3.52	2.83	1.70	<b>2.63</b>	2.81	4.12
pl-embeddings-cbow	200	3.26	2.90	2.11	1.86	2.38	<b>4.79</b>
pl-embeddings-skipg	200	3.01	2.63	1.89	1.77	2.15	4.17
fastText.wiki.pl	200	<b>3.82</b>	<b>3.14</b>	<b>2.16</b>	<b>2.31</b>	<b>2.48</b>	4.45
pl-embeddings-cbow	30	2.87	2.58	1.97	1.58	2.03	4.44
pl-embeddings-skipg	30	2.80	2.50	1.93	1.55	1.90	4.11
fastText.wiki.pl	30	<b>3.89</b>	<b>3.24</b>	<b>2.41</b>	<b>2.06</b>	<b>2.22</b>	<b>4.53</b>



# Results – selected

## Analogy Tests

Model	VS	Score	Model	VS	Score
<i>w2w-pIWNC-multi-cbow-hs</i>	100	40.82	<i>w2w-pIWNC-multi-cbow-ns</i>	300	<b>57.14</b>
<i>w2w-pIWNC-lem-cbow-ns</i>	100	47.96	<i>w2w-pIWNC-lem-skipg</i>	300	<b>60.20</b>
<i>ft-pIWNC-multi-skipg-mC20</i>	300	<b>53.30</b>	<i>ft-pIWNC-lem-skipg-mC20</i>	300	54.23
<i>ft-pIWNC-multi-skipg-mC50</i>	300	50.75	<i>ft-pIWNC-lem-skipg-mC50</i>	300	<b>59.28</b>
<i>ncp-lemmas-all-300-cbow-ns</i>	300	<b>57.95</b>	<i>ncp-forms-all-300-cbow-ns</i>	300	43.18
<i>ncp-lemmas-all-300-skipg-ns</i>	300	54.36	<i>ncp-forms-all-300-skipg-ns</i>	300	<b>46.82</b>
<i>ncp-lemmas-rest-300-cbow-ns</i>	300	<b>59.49</b>	<i>ncp-forms-rest-300-cbow-ns</i>	300	43.64



## Conclusions and Further Works

- EWBST is the hardest synonymy tests and its difficulty can be tuned with the help of a wordnet-based similarity measure
- Skip-gram model is better than CBOW according to WBST and EWBST
  - only better performance of CBOW-ns in HWBST can be attributed to a kind of generalisation caused by hypernyms in answers
  - also models from literature based on Skip-gram scheme , including *fastText.wiki.pl* express higher results
- CBOW models are superior in all cases in comparison to Skip-gram models in WBCRT
  - models with merged MWEs and PNs are better than those based on lemmas
- Skip-gram models are better in describing meaning differences, while CBOW enable broader exploration of potential lexico-semantic relations





## Conclusions and Further Works

- WCBRT can be used also as a diagnostic tool to spot worse described subdomains
- Hierarchical softmax consistently produces better results in all frequency ranges
- Smaller corpora
  - all results are much worse than those obtained on pIWNC 10 corpus
  - the models behave in a slightly different way in WBCR tests
  - Skip-gram models express higher recall, especially fastText Skip-gram with sub-word representation
- In analogy tests Skip-gram model built on pIWNC 10 is still the best one, but the difference to models built on much smaller NCP is minimal
  - the analogy tests of include mostly general and frequent words
  - the differences are small only for models based on the restricted version of NCP

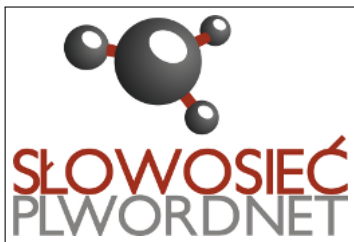


## Conclusions and Further Works

- A large comprehensive wordnet can be successfully used as a basis for two different types of MSR evaluation methods
- The datasets are enough large to conveniently partitioned according to the frequency criteria of semantic criteria
- the datasets and tests are based on human decisions expressed in the wordnet structure
- We plan
  - to develop a wordnet-based test that has properties of contextual tests
  - and tests covering all four PoS



Thank you very much for your attention!



<http://clarin-pl.eu>

<http://nlp.pwr.edu.pl>

<http://plwordnet.pwr.edu.pl>