

RUSSIAN COLLOCATION EXTRACTION BASED ON WORD EMBEDDINGS

Enikeeva E. V. (protoev@yandex.ru)

Mitrofanova O. A. (o.mitrofanova@spbu.ru)

Saint Petersburg State University

Collocation (extraction)

- Unformalized definition
- No classification according to meaning or semantic structure
- Simple frequency-based methods

raw word combination lists

Collocation

- «A collocation is a succession of two or more words that must be learnt as an integral whole and not pieced together from its component parts». (Palmer 1933)
- distinguished from idiom:
 - еловая, сосновая, кедровая, etc. шишка 'fir, pine, cedar, etc. cone' vs. важная шишка 'boss'
 - бить тревогу, рекорд, поклоны, etc. 'sound the alarm, beat the record, beat bows, etc.' vs. бить баклуши 'twiddle'

Collocation

«A collocation **AB** of language **L** is a semantic phraseme of **L** such that its signified 'X' is constructed out of the signified of the one of its two constituent lexemes — say, of **A** — and a signified 'C' [$X = A \oplus C$] such that the lexeme **B** expresses 'C' contingent on **A**» (Mel'čuk 1998)

⇒ **Lexical Function**

Lexical function

- «A lexical function **f** associates with a word **w**, called its argument or key word, the set of words and phrases, which express - contingent on **w** - the meaning of role which corresponds to **f**.»
(Mel'čuk 1984)
- lexical relation
- both paradigmatic and syntagmatic

Lexical Function

Examples:

- **Magn** - high degree

высокое напряжение 'high voltage', **значительная высота**
'considerable height'

- **Вон** - standard praise for **w**

аккуратно резать 'cut neatly', **комфортабельное судно**
'comfortable ship'

- **Oper_i** - an action for which the the i-th participant of situation is expressed by the subject

Oper₁: **оказывать сопротивление** 'show resistance'

Oper₂: **встречать сопротивление** 'meet resistance'

Method

Rodríguez-Fernández, S., Anke, L., Carlini, R., Wanner, L.
(2016) Semantics-driven recognition of collocations
using word embeddings.

training and test samples drawn from **Macmillan
Collocations Dictionary**

assumption: headword and collocate embeddings
should be trained on different corpora

Method

Argument matrix: $A_T = [a_{t_1}, \dots, a_{t_n}]$

Collocate matrix: $C_T = [c_{t_1}, \dots, c_{t_n}]$

A linear transformation matrix from training set .

$$A_T \Psi_T = C_T$$

approximated using SVD to minimize the sum:

$$\sum_{i=1}^{|T|} \left\| \Psi_T a_{t_i} - c_{t_i} \right\|^2$$

Test settings

Training and test collocations extracted from

- SynTagRus Treebank (<http://www.ruscorpora.ru/instruction-syntax.html>)

Reference data:

- Verbal collocations of Russian abstract nouns dictionary (http://dict.ruslang.ru/abstr_noun.php)

Vector models:

- RusVectōrēs project (<http://rusvectors.org>, version 3)

Test settings

LF	argument	value
<i>OPER1</i>	цель 'aim'	иметь 'have'
<i>MAGN</i>	каблук 'heel'	высокий 'high'
<i>CAUSFUNC0</i>	соревнование 'competition'	проводить 'hold'
<i>FUNC0</i>	открытие 'opening'	состояться 'be held'
<i>INCEPROPER1</i>	работа 'work'	приступать 'start'
<i>OPER2</i>	правка 'correction'	подвергаться 'undergo'
<i>REAL1-M</i>	ракета 'rocket'	запускать 'launch'
<i>REAL1</i>	средства 'means'	расходовать 'spend'
<i>INCEPFUNC0</i>	речь 'conversation'	заходить 'turn to'

Test settings

9 lexical functions

10 headwords for each LF

10 top-ranked collocates for each headword

Metrics: precision(10), recall, MRR

Filtering: UD POS-tags, NPMI score

Test settings

- **M1** – baseline: $\cos(c_i, a_i - c_i + a_j)$, where (a_i, c_i) is an example collocation for a given LF and a_j is a test headword;
- **M2** –baseline filtered by POS tags and NPMI scores;
- **M3** – proposed model, same vector spaces for headwords and collocates trained on RNC;
- **M4** – M3 filtered by POS tags and NPMI scores;
- **M5** –M3, collocate vectors from Russian Wikipedia;
- **M6** – M5 filtered by POS tags and NPMI scores.

Results (precision)

LF	<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>M6</i>
<i>OPER1</i>	0.11	0.31	0.10	0.14	0.37	0.63
<i>MAGN</i>	0.23	0.28	0.24	0.28	0.63	0.84
<i>CAUSFUNC0</i>	0.10	0.33	0.22	0.23	0.54	0.64
<i>FUNC0</i>	0.21	0.40	0.29	0.33	0.42	0.42
<i>INCEPOPER1</i>	0.10	0.38	0.64	0.64	0.15	0.15
<i>OPER2</i>	0.17	0.28	0.12	0.11	0.29	0.39
<i>REAL1-M</i>	0.20	0.66	0.24	0.26	0.40	0.52
<i>REAL1</i>	0.15	0.37	0.32	0.33	0.66	0.66
<i>INCEPFUNC0</i>	0.13	0.28	0.24	0.23	0.35	0.43

Results (recall)

LF	<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>M6</i>
<i>OPER1</i>	0.50	0.50	0.48	0.65	0.57	0.65
<i>MAGN</i>	0.68	0.70	0.70	0.70	0.83	0.78
<i>CAUSFUNC0</i>	0.33	0.45	0.87	0.9	0.80	0.80
<i>FUNC0</i>	0.70	0.80	0.90	0.81	0.58	0.58
<i>INCEPOPER1</i>	0.40	0.55	0.50	0.50	0.50	0.50
<i>OPER2</i>	0.55	0.60	0.53	0.52	0.75	0.70
<i>REAL1-M</i>	0.55	0.70	0.87	0.87	0.70	0.75
<i>REAL1</i>	0.40	0.45	0.73	0.70	0.60	0.60
<i>INCEPFUNC0</i>	0.50	0.70	0.72	0.67	0.82	0.77

Results (MRR)

LF	<i>M1</i>	<i>M2</i>	<i>M3</i>	<i>M4</i>	<i>M5</i>	<i>M6</i>
<i>OPER1</i>	0.22	0.55	0.11	0.48	0.34	0.73
<i>MAGN</i>	0.30	0.68	0.30	0.68	0.50	0.90
<i>CAUSFUNC0</i>	0.11	0.43	0.37	0.76	0.41	0.82
<i>FUNC0</i>	0.30	0.82	0.47	0.89	0.36	0.64
<i>INCEPOPER1</i>	0.11	0.60	0.01	0.64	0.15	0.48
<i>OPER2</i>	0.19	0.64	0.23	0.48	0.37	0.66
<i>REAL1-M</i>	0.08	0.77	0.36	0.70	0.34	0.86
<i>REAL1</i>	0.23	0.50	0.37	0.70	0.30	0.59
<i>INCEPFUNC0</i>	0.10	0.64	0.37	0.73	0.42	0.72

Results (discussion-1)

- NPMI filtering discards relevant examples

MAGN(*медаль*) = *золотой* ...

- training on different corpora — lower scores for specific LFs
- frequent LFs — more correct collocates in top-10

Results (discussion-2)

MAGN(довод)= решительный, убедительный,
основательный, веский, главный,
бесспорный, достаточный...

OPER1(домино) = играть, поиграть, стучать,
резаться, сыграть, игра, бильярд,
футбол...

INCEPFUNCO (день) = наступать, наставать,
начинаться, приходиться, прийти,
намечаться, длиться,
заканчиваться...

Results (discussion-3)

INCEPOPER1(азарт) = приходить, *игра*, увлекаться...
входить

OPER2(арест) = подвергаться, находиться,
сидеть подвергать, брать, миновать,
попадать...

FUNC0 (дорога) = идти, пойти, тянуться, лежать,
проходить плестись, тащиться ...

Thank you!