

# Evaluation of Czech Distributional Thesauri

Pavel Rychlý

Natural Language Processing Centre  
Faculty of Informatics, Masaryk University

December 7, 2019

# Sketch Engine Thesaurus

Lemma	Score	Freq
<a href="#">king</a>	0.242	16,899
<a href="#">prince</a>	0.213	6,355
<a href="#">charles</a>	0.189	8,952
<a href="#">elizabeth</a>	0.177	3,567
<a href="#">edward</a>	0.176	6,484
<a href="#">mary</a>	0.173	6,870
<a href="#">gentleman</a>	0.171	6,274
<a href="#">lady</a>	0.170	11,905
<a href="#">husband</a>	0.167	11,669
<a href="#">sister</a>	0.167	8,062
<a href="#">mother</a>	0.164	27,536
<a href="#">princess</a>	0.160	2,944
<a href="#">father</a>	0.159	23,824
<a href="#">wife</a>	0.157	18,308
<a href="#">brother</a>	0.155	11,049
<a href="#">henry</a>	0.151	6,699
<a href="#">daughter</a>	0.150	11,216
<a href="#">anne</a>	0.149	4,386

**queen** (*noun*) British National Corpus (BNC) freq = **7,872** (70.10 per million)



# Thesaurus evaluation

Gold standard

Source	Most similar words to <i>queen</i>
serelex	king, brooklyn, bowie, prime minister, mary, bronx, rolling stone, elton john, royal family, princess
Thesaurus.com	monarch, ruler, consort, empress, regent, female ruler, female sovereign, queen consort, queen dowager
SkE on BNC	king, prince, charles, elizabeth, edward, mary, gentleman, lady, husband, sister, mother, princess, father
SkE on enTenTen08	princess, prince, king, emperor, monarch, lord, lady, sister, lover, ruler, goddess, hero, mistress, warrior
word2vec on BNC	princess, prince, Princess, king, Diana, Queen, duke, palace, Buckingham, duchess, lady-in-waiting, Prince
powerthesaurus.org	empress, sovereign, monarch, ruler, czarina, queen consort, king, queen regnant, princess, rani, queen regent

# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden

# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden
- syntactic: *good* is to *best* as *smart* is to *smarter*
- semantic: *Paris* is to *France* as *Tokyo* is to *Japan*
- agreement by humans:

# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden
- syntactic: *good* is to *best* as *smart* is to *smarter*
- semantic: *Paris* is to *France* as *Tokyo* is to *Japan*
- agreement by humans:  
Berlin –

# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden
- syntactic: *good* is to *best* as *smart* is to *smarter*
- semantic: *Paris* is to *France* as *Tokyo* is to *Japan*
- agreement by humans:  
Berlin – Germany

# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden
- syntactic: *good* is to *best* as *smart* is to *smarter*
- semantic: *Paris* is to *France* as *Tokyo* is to *Japan*
- agreement by humans:  
Berlin – Germany  
London –



# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden
- syntactic: *good* is to *best* as *smart* is to *smarter*
- semantic: *Paris* is to *France* as *Tokyo* is to *Japan*
- agreement by humans:  
Berlin – Germany  
London – England / Britain / UK ?

# Analogy queries

- evaluation of word embeddings (word2vec)
- "  $a$  is to  $a^*$  as  $b$  is to  $b^*$ ", where  $b^*$  is hidden
- syntactic: *good* is to *best* as *smart* is to *smarter*
- semantic: *Paris* is to *France* as *Tokyo* is to *Japan*
- agreement by humans:  
Berlin – Germany  
London – England / Britain / UK ?
- best match for linear combination of vectors:  
$$\arg \max_{b^* \in V} \cos(b^*, a^* - a + b)$$

# Problems of analogy queries

- Pair of words does not define an exact relation
- Berlin – Germany: capital, biggest city
- in what time?
- Canberra

# Problems of analogy queries

- Pair of words does not define an exact relation
- Berlin – Germany: capital, biggest city
- in what time?
- Canberra, Rome

# Problems of analogy queries

- Pair of words does not define an exact relation
- Berlin – Germany: capital, biggest city
- in what time?
- Canberra, Rome
- rare words/phrases
- Baltimore – Baltimore Sun: Cincinnati –

# Problems of analogy queries

- Pair of words does not define an exact relation
- Berlin – Germany: capital, biggest city
- in what time?
- Canberra, Rome
- rare words/phrases
- Baltimore – Baltimore Sun: Cincinnati – Cincinnati Enquirer

# Outlier detection

- list of words
- find the one which is not part of the cluster
- examples:
  - red, blue, green, dark, yellow, purple, pink, orange, brown

# Outlier detection

- list of words
- find the one which is not part of the cluster
- examples:
  - red, blue, green, dark, yellow, purple, pink, orange, brown
  - t-shirt, sheet, dress, trousers, shorts, jumper, skirt, shirt, coat



# Evaluating Outlier Detection

- original data set by Camacho-Collados, Navigli
- 8 pairs of 8 words in a cluster and 8 outliers
- $8 \times 8 = 64$  queries
- Accuracy – the percentage of successfully answered queries,
- Outlier Position Percentage (OPP) Score – average percentage of the right answer (Outlier Position) in the list of possible clusters ordered by their compactness

# Problems of original data set

- English only
- needs extra knowledge
  - Mercedes Benz, BMW, Michelin, Audi, Opel, Volkswagen, Porsche, Alpina, Smart

# Problems of original data set

- English only
- needs extra knowledge
  - Mercedes Benz, BMW, Michelin, Audi, Opel, Volkswagen, Porsche, Alpina, Smart
  - (Bridgestone, Boeing, Samsung, Michael Schumacher, Angela Merkel, Capri, pineapple)

# Problems of original data set

- English only
- needs extra knowledge
  - Mercedes Benz, BMW, Michelin, Audi, Opel, Volkswagen, Porsche, Alpina, Smart
  - (Bridgestone, Boeing, Samsung, Michael Schumacher, Angela Merkel, Capri, pineapple)
  - Peter, Andrew, James, John, Thaddaeus, Bartholomew, Thomas, Noah, Matthew

# Problems of original data set

- English only
- needs extra knowledge
  - Mercedes Benz, BMW, Michelin, Audi, Opel, Volkswagen, Porsche, Alpina, Smart
  - (Bridgestone, Boeing, Samsung, Michael Schumacher, Angela Merkel, Capri, pineapple)
  - Peter, Andrew, James, John, Thaddaeus, Bartholomew, Thomas, Noah, Matthew
  - January, March, May, July, Wednesday, September, November, February, June

# Problems of original data set

- English only
- needs extra knowledge
  - Mercedes Benz, BMW, Michelin, Audi, Opel, Volkswagen, Porsche, Alpina, Smart
  - (Bridgestone, Boeing, Samsung, Michael Schumacher, Angela Merkel, Capri, pineapple)
  - Peter, Andrew, James, John, Thaddaeus, Bartholomew, Thomas, Noah, Matthew
  - January, March, May, July, Wednesday, September, November, February, June
  - tiger, dog, lion, cougar, jaguar, leopard, cheetah, wildcat, lynx
- mostly proper names (7 out of 8)

# New data set

- 5 languages: Czech, Slovak, English, German, French
- 48 clusters (8 words + 8 outliers)

# New data set – example

Colors		Electronics	
Czech	English	Czech	English
červená	red	televize	television
modrá	blue	reproduktor	speaker
zelená	green	notebook	laptop
žlutá	yellow	tablet	tablet
fialová	purple	mp3 přehrávač	mp3 player
růžová	pink	mobil	phone
oranžová	orange	rádio	radio
hnědá	brown	playstation	playstation
dřevěná	wooden	blok	notebook
skleněná	glass	sešit	workbook
temná	dark	kniha	book
zářivá	bright	CD	CD
pruhovaný	striped	energie	energy
puntíkový	dotted	světlo	light
smutná	sad	papír	paper
nízká	low	ráno	morning



- 9 clusters only, 72 queries

	OOP	Accuracy
Czes2	92.2	70.8
czTenTen12	93.4	79.2
csTenTen17	94.3	81.9
czTenTen12 (fasttext)	97.7	87.5
Czech Common Crawl	98.1	95.8