A Distributional Multi-word Thesaurus in Sketch Engine

Miloš Jakubíček, Pavel Rychlý



Lexical Computing jak@fi.muni.cz milos.jakubicek@sketchengine.eu



NLP Centre, Masaryk University pary@fi.muni.cz pavel.rychlyk@sketchengine.eu

RASLAN 2019, 7. 12. 2019

Miloš Jakubíček, Pavel Rychlý

Lexical Computing & Masaryk University

Outline



2 Conclusions

Miloš Jakubíček, Pavel Rychlý A Distributional Multi-word Thesaurusin Sketch Engine Lexical Computing & Masaryk University

Motivation

- word sketch and thesaurus in Sketch Engine
- word sketch: both single words and multi-words
- thesaurus: only single words so far ⇒ extension to multi-words coming

Word sketch

modifiers of "resource"		
natural	•••	
natural resources		
human	•••	
human resources		
water	•••	
water resources		
limited	•••	
limited resources		
valuable	•••	
a valuable resource		
online	•••	
online resources		
financial	•••	
financial resources		

nouns modified by "resource"			
management	•••		
resource management			
allocation	•••		
resource allocation			
center	•••		
Resource Center			
centre	•••		
Resource Centre			
efficiency	•••		
resource efficiency			
conservation	•••		
resource conservation			
guide	•••		
Resource Guide			

verbs with "resourc object	e" as
allocate	•••
manage	•••
access	•••
provide	•••
utilize	•••
conserve	•••
share	•••
pool	•••
invest	•••
need resources needed to	•••

verbs with "resource" as subject		
include	•••	
resources including		
permit	•••	
resources permit		
exist	•••	
resources exist		
provide	•••	
resource provides		
relate	•••	
•		
resources relating to		
cover		
cover resources covering	•••	
cover resources covering allow	•••	

Miloš Jakubíček, Pavel Rychlý

Thesaurus

test (noun) Alternative PoS: <u>verb</u> (freq: 941,372) enTenTen [2012] freq = <u>1,915,482</u> (147.70 per million)



Miloš Jakubíček, Pavel Rychlý

Lexical Computing & Masaryk University

Outlir 0 Introduction 00000000000 Conclusions O

Thesaurus



Miloš Jakubíček, Pavel Rychlý

Lexical Computing & Masaryk University

Thesaurus



Miloš Jakubíček, Pavel Rychlý

Lexical Computing & Masaryk University

Outline	Introduction	Conclusions
0	0000000000	0

$$Dist(w_1, w_2) = \frac{\sum_{(r,c) \in ctx(w_1) \cap ctx(w_2)} AS_{(w_1,r,c)} + AS_{(w_2,r,c)} - (AS_{(w_1,r,c)} - AS_{(w_2,r,c)})^2 / 50}{\sum_{i \in ws_1} AS_i + \sum_{i \in ws_2} AS_i}$$

The term $(AS_i - AS_j)^2/50$ is subtracted in order to give less weight to shared triples, where the triple is far more salient with w_1 than w_2 or vice versa. We find that this contributes to more readily interpretable results, where words of similar frequency are more often identified as near neighbours of each other.

Extensions to multi-words

The most difficult thing is to say what a multi-word is. The rest is basically straighforward.

bucket as noun 153,053× 👻

₹			D	×	←	
	verbs with "buc	ket" as	object		modifie	ers of "
	kick kicked the bucket	925 t	7.8		bucket + kick	Ø
	empty	420	7.5	\odot	bucket + kick	>
	sweat sweating buckets	228	7.4	•=	kick	Ø

...

<s> I think all of Drogheda was listening to the radio when the draw was being made and I thought **somebody** was **pulling** my **leg** when we were drawn against them.

<s>"I thought **somebody** was **pulling** my **leg**, but one morning, just after it snowed, I saw footprints going to the tower and I thought it was one of the guys, but I went in there to do some work and nobody was there. </s>

<s>"Both Andy and I were taken aback when we first received the initial email from him, to the extent that we thought it was **somebody pulling** our **leg**. </s>

<s> She has no problem if **somebody pulls** her **leg** and she will give right back to you. </s>

<s> Ok, somebody 's pulling my leg, right?! </s>

WORD SKETCH

English Web 2015 (enTenTen15) 🔍

	WS collocations	Frequency	Score
1	$pull-v$ particles after "%w" with object $out\-x$	127,030	9.25
2	pull-v particles after "%w" out-x	65,648	7.66
3	pull-v particles after "%w" with object $off-x$	57,617	9.57
4	pull-v pronominal subjects of "%w" $h-d$	48,229	5.29
5	pull-v pronominal objects of "%w" it-d	48,039	5.93
6	pull-v pronominal subjects of "%w" i-d	47,720	4.32
7	$pull-v$ particles after "%w" with object $up\-x$	42,505	7.6

Multi-word thesaurus

salient multi-word sketches == multi-words

same scoring

Miloš Jakubíček, Pavel Rychlý A Distributional Multi-word Thesaurusin Sketch Engine

\cap		÷			r	1
\circ	u	Ľ	1	1	1	1
0						

Table: Thesaurus items for the phrase "kohl helmut chancellor" on the $\mathsf{BNC}.$

score	frequency	item
1.00	755	kohl chancellor helmut
0.88	1790	kohl helmut
0.75	7307	kohl chancellor
0.34	606	kohl
0.20	140	mitterrand president
0.18	536	chancellor kohl
0.17	340	bush president us
0.17	153	bush us president
0.17	20	re-unification
0.17	116	bush us
0.16	370	bush george president
0.16	283	bush president george
0.16	116	clinton president

Miloš Jakubíček, Pavel Rychlý

Lexical Computing & Masaryk University

Conclusions

- initial approach complete, to be implemented within Sketch Engine
- need to disregard ordering of multi-word items
- evaluation and comparison (as in preceding talk)