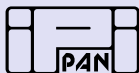


Identyfikacja terminów wielowyrazowych w tekście TermoPL

Małgorzata Marciniak, Agnieszka Mykowiecka,
Piotr Rychlik

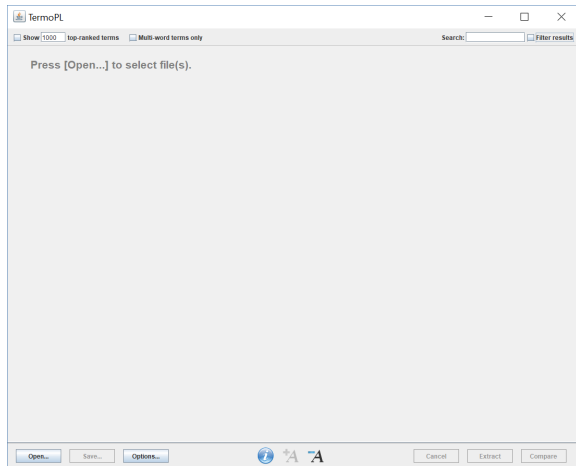


INSTITUTE OF COMPUTER SCIENCE
POLISH ACADEMY OF SCIENCES
ul. J. K. Ordona 21, 01-237 Warszawa

Warsztaty CLARIN-PL, Toruń 16-17 XI 2018

instalacja: <http://zil.ipipan.waw.pl/TermoPL>

uruchamianie: termopl.bat zamiast TermoPL.jar (.bat rozszerza dostępną pamięć)



- kolokacje (kryteria leksykalne/korpusowe)
klatka piersiowa konwalia majowa
- nazwy własne (kryteria dziedzinowe)
New York Nowa Zelandia, Rio de Janeiro
- **terminy dziedzinowe** (kryteria dziedzinowe)
skrzynia biegów, podatek od osób fizycznych, komedia romantyczna
- grupy składniowe ("słabe" kryteria lingwistyczne)
[The debate]_{NP} and [vote]_{NP} [in Cyprus' parliament]_{NP} has now been [postponed]_V until [18:00 local time (16:00 GMT)]_{NP} on [Tuesday]_{NP}
- frazy (kryteria lingwistyczne)
[[The debate]_{NP} and [vote]_{NP}]_{NP-SUBJ} [in [Cyprus'_{ADJ-MOD} parliament]_{NP}]_{Prep-MOD} [has [now]_{ADV-MOD} been postponed]_{VP} ...

Kryteria ośrodka kolokacji:

student

Maksymalny odstęp: 0

Zachowaj szyk:

Kryteria kolokatu:

Część mowy: Przym./Imiesł.

Kontekst z lewej: 1

Kontekst z prawej: 1

Wielkość próbek: 10000

Min. współwystąpien: 5

Ośrodek kolokacji występuje w korpusie 5153 razy.

Podajemy wyniki na podstawie wszystkich współwystąpień.

Znaleziono 56 kolokacji spełniających zadane kryteria.

#	Kolokacja	Pasujące współwystąpienia (kliknij na frekwencję, aby wyświetlić przykłady)	Ogółem	Chi ²
1.	wrocławski	wrocławski__student (7), student__wrocławskiej (4), student__wrocławskiego (3), student__wrocławski (1),	15	2,707,982.49
2.	amerykański	amerykański__student (8), student__amerykański (1), student__amerykańskiego (1),	10	1,805,321.65
3.	dziesiąty	dziesiąty__student (5),	5	902,660.82
4.	pomysłowy	pomysłowy__student (4), student__pomysłowy (1),	5	902,660.82
5.	niepełnosprawny	niepełnosprawny__student (5), niepełnosprawnych__student (1),	6	649,910.99
6.	poznański	student__poznańskiej (4), poznańskiego__student (1),	5	451,325.41
7.	gdański	student__gdańskiej (4), gdańska__student (1), gdańsku__student (1), gdański__student (1),	7	442,294.71
8.	upiec	upieczony__student (10),	10	21,969.57
9.	zaoczny	student__zaocznego (3), student__zaocznej (2), student__zaoczny (2), zaoczny__student (1),	8	20,727.62
10.	dwudziestoletni	dwudziestoletni__student (6),	6	13,360.91
11.	prawy	student__prawa (95), prawa__student (1),	96	4,909.26
12.	wieczny	wieczny__student (11),	11	3,006.99
13.	letni	letni__student (11),	11	2,966.79
14.	biedny	biedny__student (12),	12	2,846.13
15.	czwarty	student__czwartego (12), czwarty__student (6),	18	2,223.5

Cel

wydobycie specyficznej terminologii z tekstów dotyczących wybranej dziedziny.

Założenia:

- korzystamy z danych otagowanych morfologicznie,
- **identyfikacja terminów**: gramatyka opisującą składnię interesujących nas fraz,
- **szeregowanie terminów**: informacje o częstości występowania w tekście, C-value,
- uwzględniamy frazy zagnieżdżone,
- wykorzystujemy lematyzację do grupowania form fleksyjnych.

Definicja słownikowa

Wyraz albo połączenie wyrazowe o specjalnym, konwencjonalnie ustalonym znaczeniu naukowym lub technicznym; (Doroszewski)

Definicja robocza

Fraza rzeczownikowa, która w tekstach dziedzinowych występuje dostatecznie często by przypuszczać, że opisuje pojęcie istotne dla dziedziny. Częstość tej frazy w tekstach spoza dziedziny jest niższa.

"student", NKJP 1.2mln

student (122 wystąpienia)
student medycyny
student wszystkich uczelni Szczecina
student Wyższej Szkoły Zarządzania
student Akademii Ekonomicznej
student uczelni artystycznych
młody student ekonomii
student filologii rosyjskiej
student nauk społecznych
student siłaczy (student siłacz)
zaprzyjaźniony student
student architektury
student ASP
student AWF
afrykański student
student politologii
dzisiejszy student
student warszawski

student (1089)

liczba studentów

pożyczka studencka

zrzeszenie studentów

wzrost liczby studentów

zrzeszenie studentów polskich

zwiększyć liczbę studentów

kształcenie studentów

student studium zaocznego

student studium

parlament studenta

student pierwszego roku

przejęcie majątku zrzeszenia studentów polskich

kredyt studencki

- Zgromadzenie tekstów dziedzinowych.
- Wstępna analiza lingwistyczna — tagowanie (przypisanie formy podstawowej, części mowy oraz charakterystyki morfologicznej), można w tym celu użyć Korpusomatu.
- Identyfikacja fraz — kandydatów na terminy.
- Szeregowanie fraz.
- Selekcja fraz.

poz	termin	[forma podstawowa]	wsp.C	dl	#	#wew	kont.
1	pan	[pan]	240.3	1	2405	273	170
2	człowiek	[człowiek]	177.5	1	1777	692	387
3	czas	[czas]	154.8	1	1550	537	315
4	praca	[praca]	136.5	1	1367	758	373
5	osoba	[osoba]	118.5	1	1187	464	241
6	życie	[życie]	112.9	1	1131	534	292
7	sprawa	[sprawa]	111.5	1	1117	534	321
8	Polska	[polska]	110.7	1	1109	286	177
9	miejsce	[miejsce]	103.3	1	1035	483	263
10	dziecko	[dziecko]	102.3	1	1025	329	209
11	unia europejski	[unia europejska]	97.3	2	99	46	27

klikamy **Open**, a następnie wybieramy plik (NKJP1m_hash.txt)
klikamy **Extract**

klikamy **Open**, a następnie wybieramy plik (NKJP1m_hash.txt)
klikamy **Extract**


TermoPL - NKJP1m_hash_NPMI3

Show 1000 top-ranked terms Multi-word terms only Search: student

#	Rank	Term	C-value	Length	Freq_s	Freq_in	Cor
1	1	pan	240,34	1	2405	273	170
2	2	człowiek	177,52	1	1777	692	387
3	3	czas	154,83	1	1550	537	315
4	4	praca	136,5	1	1367	758	373
5	5	osoba	118,51	1	1187	464	241
6	6	życie	112,92	1	1131	534	292
7	7	sprawa	111,53	1	1117	534	321
8	8	polska	110,74	1	1109	286	177
9	9	miejsce	103,32	1	1035	483	263

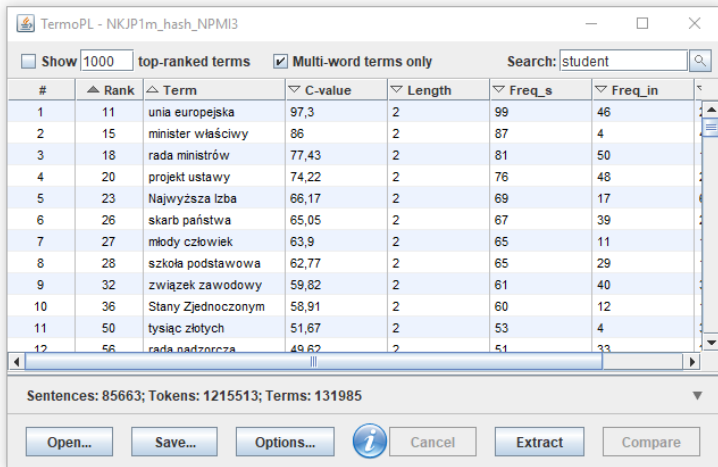
Forms:
 pracy[326,434], praca[20,16], prace[120,51], praca[41,29], prac[14,86], prace[54,46], Prace[3,7], Praca[20,7], Pracy[1,62],
 pracach[5,12], Pracę[2,0], pracami[0,3], PRACE[1,0], pracom[0,1], pra[0,2], Praca[1,0], Prac[0,2], Pracami[1,0]

Sentences: 85663; Tokens: 1215513; Terms: 131985

Open... Save... Options...  Cancel Extract Compare

klikamy 'Multi-word terms only'

klikamy 'Multi-word terms only'




TerminoPL - NKJP1m_hash_NPMI3

Show 1000 top-ranked terms Multi-word terms only Search: student

#	Rank	Term	C-value	Length	Freq_s	Freq_in
1	11	unia europejska	97,3	2	99	46
2	15	minister właściwy	86	2	87	4
3	18	rada ministrów	77,43	2	81	50
4	20	projekt ustawy	74,22	2	76	48
5	23	Najwyższa Izba	66,17	2	69	17
6	26	skarb państwa	65,05	2	67	39
7	27	młody człowiek	63,9	2	65	11
8	28	szkoła podstawowa	62,77	2	65	29
9	32	związek zawodowy	59,82	2	61	40
10	36	Stany Zjednoczonym	58,91	2	60	12
11	50	tysiąc złotych	51,67	2	53	4
12	56	rada nadzorcza	49,62	2	51	33

Sentences: 85663; Tokens: 1215513; Terms: 131985

Open... Save... Options...  Cancel Extract Compare

wpisujemy tekst w polu 'Search' i naciskamy 'Enter'

wpisujemy tekst w polu 'Search' i naciskamy 'Enter'

TermoPL - NKJP1m_hash_NPM13

Show 1000 top-ranked terms Multi-word terms only 28 matches Search: student Filter results

#	Rank	Term	C-value	Length	Freq_s	Freq_in	Context #
1	2052	student medycyny	2	2	2	0	0
2	2052	student wszystkich uczelni Szczecina	2	4	1	0	0
3	2052	student Wyższej Szkoły Zarządzania	2	4	1	0	0
4	2179	student Akademii Ekonomicznej	1,58	3	1	0	0
5	2179	student uczelni artystycznych	1,58	3	1	0	0
6	2179	młody student ekonomii	1,58	3	1	0	0
7	2179	student filologii rosyjskiej	1,58	3	1	0	0
8	2301	student siłaczy	1	2	1	0	0
9	2301	zaprzyjaźniony student	1	2	1	0	0
10	2301	student architektury	1	2	1	0	0
11	2301	student ASP	1	2	1	0	0
12	2301	student AWF	1	2	1	0	0
13	2301	afrykański student	1	2	1	0	0
14	2301	student politologii	1	2	1	0	0
15	2301	wszystek student	1	2	1	0	0
16	2301	dzisiejszy student	1	2	1	0	0
17	2301	student warszawski	1	2	1	0	0
18	2301	student SZSP	1	2	1	0	0
19	2301	student izraelski	1	2	1	0	0
20	2301	student historii	1	2	1	0	0
21	2301	student fizyki	1	2	1	0	0
22	2301	student UMCS	1	2	1	0	0

Forms:

Studenci medycyny[1,0], studenta medycyny[1,0]

klikamy 'trójkącik' przy nazwie kolumny, np. 'Length'


klikamy 'trójkącik' przy nazwie kolumny, np. 'Length'

TermoPL - NKJP1m_hash_NPMI3

Show 1000 top-ranked terms Multi-word terms only Search: student

#	Rank	Term	C-v...	L...	Freq_s	Freq...
2	1762	pełna realizacja rządowego programu wspierania ośrodków wczesnej rehabilitacji dzieci niepełnosprawnych	3	10	1	0
3	1762	najciekawszy wątek piątkowej rozprawy lustracyjnej posła SLD Tadeusza Matyka	3	9	1	0
4	1762	obowiązek uprzedniego uzyskiwania zezwolenia właściwego terenowego organu administracji państwowej	3	9	1	0
5	1762	protokół siedemdziesiątego pierwszego posiedzenia Senatu Rzeczypospolitej Polskiej czwartej kadencji	3	9	1	0
6	1762	przewodniczący łódzkiego oddziału Krajowego Związku Zawodowego Pracowników Ratownictwa Medycznego	3	9	1	0
7	1762	sala intensywnej terapii Oddziału Leczenia Stanów Nagłych Instytutu Pediatrii	3	9	1	0
8	1762	absolwentka Szkoły Oficerskiej Centralnego Ośrodka Szkolenia Służby Więziennej	3	8	1	0
9	1762	czas wladania wielkiego mistrza Zakonu Henryka von Hohenlohe	3	8	1	0
10	1762	długoletni dyrygent związkowy Związku Kół Śpiewaczych Śląska Opolskiego	3	8	1	0
11	1762	dyplom University of Cambridge Certificate in Advanced English	3	8	1	0
12	1762	główna teza programu polityczno-gospodarczego Grupy Inicjatywnej Partii Robotniczej	3	8	1	0
13	1762	onć tenoroczno. Międzynarodowego Festiwalu Sztuki Autorów Zdział Filmowych	3	8	1	0

Sentences: 85663; Tokens: 1215513; Terms: 131985

Open... Save... Options...  Cancel Extract Compare

Zapisanie listy terminów



TermoPL - NKJP1m_hash_NPMI3

Show 1000 top-ranked terms Multi-word terms only

#	Rank	Term
1	11	unia europejska
2	15	minister właściwy
3	18	rada ministrów
4	20	projekt ustawy
5	23	Najwyższa Izba
6	26	skarb państwa
7	27	młody człowiek
8	28	szkoła podstawowa
9	32	Stany Zjednoczonym
10	33	związek zawodowy
11	50	tysiąc złotych
12	57	rada nadzorcza
13	58	zabranie głosu
14	60	II wojna światowa
15	62	porządek dzienny
16	66	klub parlamentarny
17	68	druga strona
18	76	Paweł II
19	77	pan posła
20	78	miejsce pracy
21	81	wojna światowa
22	85	działalność gospodarcza

Forms:
 Studenci medycyny[1,0], studenta medycyny[1,0]

Options

Filters Grammar Search Save

Select fields to be saved:

- #
- Rank
- Term (simplified form)
- Term (base form)
- C-value
- LL/TFITF/CSmw/TW
- Length
- Freq_s
- Freq_in
- Context #
- Save all forms (if available)

Select all fields Select fields for corpora comparing Clear

Cancel OK

Open... Save... Options... Cancel Extract Compare

- rzeczownik, akronim lub skrót rzeczownika:
 - *podatek, angiografia,*
 - *PKB, USG*
 - *ust.(awa),*
- rzeczownik z przymiotnikiem (który wystąpił po lub rzadziej przed rzeczownikiem):
 - *stosunki gospodarcze,*
 - *granulocyty obojętnochłonne;*
- sekwencja rzeczownika z rzeczownikiem w dopełniaczu:
 - *udar_{n,nom} mózgu_{n,gen};*
 - *kodeks_{n,nom} pracy_{n,gen};*
- kombinacja powyższych dwóch struktur:
 - *europejski_{adj} rynek_{n,nom} usług_{n,gen} finansowych_{adj},*
 - *wodonercze niewielkiego stopnia dolnego układu podwójnego nerki prawej;*

- fraza rzeczownikowa modyfikowana frazą przyimkową:
 - *wierzytelność podatnika wobec skarbu państwa,*
 - *podatek dochodowy od osoby fizycznej;*
 - *poziom hormonów we krwi;*
- można uwzględnić koordynację:
 - *bezsorna i wymagalna wierzytelność podatnika wobec skarbu państwa,*
 - *zapalenie mózgu i rdzenia,*
 - *oddział alergologii, endokrynologii i pediatrii ogólnej.*

NPP : \$*NAP* *NAP_GEN**;
NAP[*agreement*] : *AP** *N* *AP**;
NAP_GEN[*case* = *gen*] : *NAP*;
AP : *ADJ* | *ADJA* *DASH* *ADJ* | *PPAS*;
N[*pos* = *subst*, *ger*];
ADJ[*pos* = *adj*];
ADJA[*pos* = *adja*];
PPAS[*pos* = *ppas*];
DASH[*form* = "-"];

jama ustny [jama ustna] 241

śluzówka jama ustny [śluzówka jamy ustnej] 79

jama ustny czysty [jama ustna czysta] 4

suchość jama ustny [suchość jamy ustnej] 4

błona śluzowy jama ustny [błona śluzowa jamy ustnej] 3

grzybica śluzówka jama ustny [Grzybica śluzówek jamy ustnej] 2

pielęgnacja jama ustny [pielęgnacji jamy ustnej] 2

płukanie jama ustny [Płukanie jamy ustnej] 2

grzybica jama ustny [Grzybica jamy ustnej] 2

pędzlować jama ustny [pędzlowanie jamy ustnej] 2

sanacja jama ustny [sanacji jamy ustnej] 1

silny grzybica jama ustny [silna grzybica jamy ustnej] 1

jama ustny pleśniawka [jama ustnej pleśniawki] 1

zapalenie jama ustny [zapalenie jamy ustnej] 1

dno jama ustny [dno jamy ustnej] 1

Zwykle nie traktujemy jako terminów dziedzinowych fraz składających się z:

- słów wskazujących na określenie czasu, jak np: *miesiąc, dzień*;
- nazwy dni i miesięcy, np: *styczeń, poniedziałek*;
- przymiotników wymagających kontekstu do interpretacji np: *inny, niektóry, jakiś, pewien*.

Chcemy też wykluczyć przyimki złożone:

- [*w kierunku*] zapalenia nerek → *kierunek zapalenia nerek*;
- [*pod postacią*] podatku VAT → *postać podatku VAT*;
- [*pod kątem*] diagnostyki obrazowej → *kąt diagnostyki obrazowej*;
- [*pod kątem*] prostym → *kąt prosty*.

Wykluczenie niektórych słów – opcje



TermoPL - NKJP1m_hash_NPMI3

Show 1000 top-ranked terms Multi-word terms only 30 matches Search: student Filter results

#	Rank	Term
1	1	pan
2	2	czlowiek
3	3	czas
4	4	praca
5	5	osoba
6	6	zycie
7	7	sprawa
8	8	polska
9	9	miejsce
10	10	dziecko
11	11	unia europejska
12	12	kto
13	13	pani
14	14	dom
15	15	minister właściwy
16	16	miasto
17	17	prawo
18	18	rada ministrów
19	19	świat
20	20	projekt ustawy
21	21	poseł
22	22	państwo
23	23	Najwyższa Izba
24	24	ręka
25	25	minister
26	26	skarb państwa

Options

Filters Grammar Search Save

Stop Words Compound Prepositions Common Terms

Check for stop words:

- to
- co
- ten
- taki
- który
- niektóry
- każdy
- możliwy
- niniejszy
- jak
- jaki
- pewien
- sam
- swój
- raz
- wtedy
- gdy
- miesiąc
- rok
- dzień
- styczeń
- luty
- marzec
- kwiecień
- maj

Load... and merge Save... Clear

Cancel OK

Open... Save... Options... Cancel Extract Compare

	pojedyncza	mnoga
nom	<i>przewlekły nieżyt żołądka</i>	<i>przewlekłe nieżyty żołądka</i>
gen	<i>przewlekłego nieżytu żołądka</i>	<i>przewlekłych nieżytów żołądka</i>
dat	<i>przewlekłemu nieżytowi żołądka</i>	<i>przewlekłym nieżytom żołądka</i>
acc	<i>przewlekły nieżyt żołądka</i>	<i>przewlekłe nieżyty żołądka</i>
inst	<i>przewlekłym nieżytem żołądka</i>	<i>przewlekłymi nieżytami żołądka</i>
loc	<i>przewlekłym nieżycie żołądka</i>	<i>przewlekłych nieżytach żołądka</i>

Wykorzystujemy uproszczoną formę podstawową:

- *przewlekły nieżyt żołądka* → *przewlekły nieżyt żołądek*;
- *ostra niewydolność nerek* → *ostry niewydolność nerka*.

Taką samą uproszczoną formę podstawową mają:

- frazy w liczbie mnogiej i pojedynczej np. *zapalenie ucha* i *zapalenie uszu*, uproszczona: *zapalenie ucho*;
- przymiotniki w różnych stopniach (mały, mniejszy) np. *miednica mała* (częściej *mała miednica* — opisuje rozmiar) podczas gdy *miednica mniejsza* (określenie anatomiczne), uproszczona: *miednica mały*;
- pozytywne i zanegowane imiesłowy przymiotnikowe . *powiększony*/*niepowiększony* mają formę podstawową *powiększyć_{inf}*;
- gerundia i imiesłowy mają bezokoliczniki jako formy podstawowe:
 - *usunięcie_{ger} kamienia_{subst:gen}* — operacja,
 - *usunięty_{ppas} kamień_{subst:nom}* — opis kamienia,forma uproszczona: *usunąć_{inf} kamień_{subst}*.

<i>planowa</i>	<i>operacja</i>	<i>przepukliny</i>	<i>pachwinowej</i>	<i>lewostronnej</i>
	<i>operacja</i>	<i>przepukliny</i>	<i>pachwinowej</i>	<i>lewostronnej</i>
<i>planowa</i>	<i>operacja</i>	<i>przepukliny</i>	<i>pachwinowej</i>	
	<i>operacja</i>	<i>przepukliny</i>	<i>pachwinowej</i>	
		<i>przepuklina</i>	<i>pachwinowa</i>	<i>lewostronna</i>
	<i>lewostronna</i>	<i>przepuklina</i>	<i>pachwinowa</i>	
		<i>przepuklina</i>	<i>pachwinowa</i>	<i>prawostronna</i>
		<i>przepuklina</i>	<i>pachwinowa</i>	<i>obustronna</i>
	<i>prawostronna</i>	<i>przepuklina</i>	<i>pachwinowa</i>	
	<i>uwięźnięta</i>	<i>przepuklina</i>	<i>pachwinowa</i>	<i>prawostronna</i>

Metody liczenia kontekstów (ograniczamy do jednego słowa):

- 1 liczba różnych kontekstów liczona po obu stronach razem;
- 2 suma różnych kontekstów po obu stronach;
- 3 maksimum z kontekstów liczonych z lewej i prawj strony osobno.

Konteksty dla frazy: *przepuklina pachwinowa*:

- 1 'operacja'–'lewostronny', 'operacja'–[pusty],
[pusty]–'lewostronny', 'lewostronny'–[pusty],
[pusty]–'prawostronny', [pusty]–'obustronny',
'prawostronny'–[pusty], 'uwięźnięty'–'prawostronny';
- 2 'operacja', 'lewostronny', 'prawostronny', 'obustronny',
'uwięźnięty';
- 3 'operacja', 'lewostronny', 'prawostronny', 'uwięźnięty' (lewych
o jeden więcej).

NPMI wykorzystujemy do oceny siły powiązania pomiędzy słowami.

referencja bibliograficzna

Gerlof Bouma, 2009, *Normalized (pointwise) mutual information in collocation extraction.*, w: *Proceedings of the Biennial GSCL Conference 2009*, strony 31—40.

Przykład

infekcja górnych dróg oddechowych

Noun; Adj; Noun; Adj;

infekcja | górnych dróg | oddechowych

infekcja górny droga oddechowy

bigram	NPMI
infekcja górny	0.66
górny droga	0.79
droga oddechowy	0.95

Poprawne gramatycznie podfrazy	Podfrazy z wykorzystaniem NPMI
'infekcja' 'górnny' 'droga' 'oddechowy'	'infekcja' 'górnny' 'droga' 'oddechowy'
infekcja górnych dróg oddechowych	infekcja górnych dróg oddechowych
infekcja górnych dróg	—
infekcja	infekcja
górne drogi oddechowe	górne drogi oddechowe
górne drogi	—
drogi oddechowe	drogi oddechowe
drogi	drogi

*prawidłowa*_{adj} *mikroflora*_{noun} *górných*_{adj} *dróg*_{noun} *oddechowych*_{adj}
—> *prawidłowa mikroflora* oraz *górne drogi oddechowe*

*częste*_{adj} *infekcje*_{noun} *górných*_{adj} *dróg*_{noun} *oddechowych*_{adj} —>
częste modyfikuje całą frazę *infekcje górnych dróg oddechowych*

Modyfikacja:

- szukamy najślabszej pozycji pozwalającej podzielić frazę na dwie podfrazy rzeczownikowe;
- jeśli różnica pomnięcy nastłabszym miejscem podziału a tym dzielącym na dwie frazy rzeczownikowe jest mniejsza od ustalonego progu to preferujemy podział na dwie frazy rzeczownikowe.

Cel

Na podstawie porównania wyników ekstrakcji terminologii dla dwóch korpusów mają być wskazane frazy:

- bardziej specyficzne dla innej dziedziny (porównanie z terminologią wydobytą z innego korpusu dziedzinowego)
- terminy ogólne np. ” *własny sposób, lewa strona, trudne zadanie* (porównanie z korpusem języka ogólnego).

Zaimplementowane metody wykorzystują:

- Log-Likelihood (LL – logarytm wiarygodności): na ile różni się częstość konkretnego terminu w dwóch porównywanych korpusach;
- Term Frequency Inverse Term Frequency (TFITF): łączy częstość występowania w korpusie dziedzinowym z odwrotną częstością występowania w korpusie ogólnym (liczoną jako stosunek wielkości korpusu do częstości badanego terminu);
- Contrastive Selection of Multi-Word Terms(CSmw): dla terminów wielowyrazowych, uwzględnia zarówno częstość występowania pełnych terminów, ale też częstość występowania słów stanowiących element główny badanej frazy.

TermoPL - autor_hash_NPMI3

Show 1000 top-ranked terms Multi-word terms only

#	Rank	Term
1	1	prawo autorskie
2	2	prawo
3	3	prawo własności
4	4	utwór muzyczny
5	5	wyłączne prawo
6	6	autor
7	7	publiczne wykonywanie
8	8	utwór
9	9	prawo własności intelektualnej
10	10	własność
11	11	statut Anny
12	12	prawo wyłączne
13	13	monopol autorski
14	14	własność literacka
15	15	konwencja berneńska
16	16	prawo autora
17	17	własność intelektualna
18	18	dobra niematerialne
19	19	prawo naturalne
20	20	książka
21	21	prawo osobiste
22	22	twórca
23	23	ochrona
24	24	publiczne wykonywanie utworów muzycznych
25	25	stan rzeczy
26	26	ustawa

Options

Filters Grammar Search Save

Compare corpora using:

- Corpora-comparing log-likelihood (LL),
- Term Frequency Inverse Term Frequency (TFITF),
- Contrastive Selection of multi-word terms (CSmw),
- Term Weight (TW = $0.9 \times DR + 0.3 \times DC$),
- for terms with C-value greater than
- for terms with frequency greater than

Use C-values instead of frequencies

Contrastive terms: KSIAZKAPUBL_hash_NPMI3 Select

Use NPMI method to search for nested terms

- Method 1
- Method 2
- Method 3 Preference factor: %

Try to divide the phrase into subphrases so as to at least one of them satisfies the grammar rules. Prefer cases where both phrases obtained after splitting are accepted by the given grammar. This preference is expressed by the "Preference factor". Choose the weakest possible connection point according to NPMI value to do the split. Continue this process for the resulting subphrases. If the phrase cannot be split in such a way, use Method 2.

Cancel OK

Open... Save... Options... Cancel Extract Compare

Porównanie dwóch list terminów



TermoPL - autor_hash_NPMI3 vs. KSIAZKAPUBL_hash_NPMI3 [LL]

Show 1000 top-ranked terms Multi-word terms only 1 match Search: student Filter results

#	Rank	Term	C-value	LL	Length	Freq_s	Freq_in	i
43	43	interes	29,2	13,33	1	294	234	11
44	44	jedna strona	29	0,6	2	30	2	2
45	45	organizacja zbiorowego zarzadzania	28,53	-	3	19	6	6
46	46	monopol	28,28	45,69	1	285	222	10
47	47	rzecz	26,48	0,04	1	268	165	51
48	48	sposob	25,65	0,05	1	258	140	91
49	49	prawo autorskie majatkowe	25,36	-	3	17	6	6
50	50	monopol eksploatacyjny	25	-	2	26	6	6
51	51	Stany zjednoczone	23,92	-	2	25	13	12
52	52	XVIII wiek	23,5	-	2	26	10	4
53	53	CO	23,34	-	1	235	24	15
54	54	dzieło sztuki	23	14	2	24	8	8
55	55	osoba	22,74	1,11	1	229	89	54
56	56	rynek wydawniczy	22,14	-	2	24	13	7
57	57	kompozycja muzyczna	22	-	2	23	3	3
58	58	ochrona prawna	21,91	-	2	23	12	11
59	59	praca	21,85	0,09	1	220	107	71
60	60	prawo zakazowe	21	-	2	22	9	9
61	61	komercyjna eksploatacja	20,5	-	2	22	18	12
62	62	domena publiczna	20,25	-	2	22	7	4
63	63	zasada	20,25	2,31	1	204	129	85
64	64	mechaniczna reprodukcja	20	-	2	22	12	6
65	64	XX wiek	20	-	2	21	3	3
66	64	wielka Brytania	20	-	2	21	2	2
67	65	stacja radiowa	19,8	38,59	2	21	12	10
68	66	...	19,7	19,74	1	199	117	74

Open... Save... Options... Cancel Extract Compare

Cel ekstrakcji terminologii:

wydobycie specyficznej terminologii z tekstów dotyczących wybranej dziedziny.

Co otrzymujemy:

- listę terminów w formie podstawowej lub w uproszczonej formie podstawowej;
- listę możemy sortować według:
 - adekwatności frazy jako terminu;
 - częstości w tekstach;
 - długości frazy
- możliwość wyboru fraz wielowyrazowych;
- możliwość zebrania wystąpień fraz w korpusie wraz z częstościami;
- możliwość porównań dwóch zbiorów terminów.

Opracowany w ramach projektu Clarin.PL

- Java Runtime Environment w wersji 7 lub nowszej;
- Wymaga Morfeusza 2 do wygenerowania formy podstawowej z uproszczonej formy;
- Wymaga otagowanego i ujednoznaczonego korpusu danych w jednym z formatów:
 - NKJP;
 - XCES;
 - XML-owe wyjście z Korpusomatu;
 - zapis uproszczony: token # lemat # tag.
- na wyjściu: lista uporządkowanych terminów (w uproszczonych formach lub zrekonstruowanych formach podstawowych wraz z formami znalezionych fraz).

Przydatne adresy internetowe

- <http://clarin-pl.eu/en/uslugi/>
- <http://zil.ipipan.waw.pl/TermoPL>